



Revision Notes: Decision Trees and Associated Concepts

Introduction to Decision Trees

A **Decision Tree** is a flowchart-like structure used for decision-making and predictive modeling. It comprises nodes that test an attribute, branches that represent the outcome of the test, and leaf nodes that correspond to decisions or classifications. The key concepts associated with decision trees include information gain, entropy, and Gini impurity.

Key Concepts

Entropy

- **Entropy** measures the amount of randomness or uncertainty in data. It quantifies the impurity in a group of examples.
- Formula: $H(D) = -\sum p(x) \log p(x)$, where $p(x)$ is the probability of class x .
- Entropy is highest when the classes are evenly split, and lowest when they are pure (i.e., all one class).

Gini Impurity

- **Gini Impurity** is an alternative to entropy, often preferred for its computational efficiency.
- Formula: $Gini(D) = 1 - \sum(p(x)^2)$.
- A lower Gini impurity indicates a better split in the data.
- It ranges from 0 (perfect purity) to 0.5 (maximum impurity for a binary classification) [\[4:8+source\]](#).

Information Gain

- **Information Gain (IG)** is the reduction in entropy or impurity after a dataset is split on an attribute.



their probability of occurring.

Constructing a Decision Tree

1. **Choose the Best Attribute to Split:** This involves calculating the information gain for each attribute and choosing the one with the highest gain.
2. **Splitting:** Divide the dataset based on the selected attribute to create branches.
3. **Repeat:** The process is repeated recursively for each branch using the subset of data reaching the branch [\[4:0+source\]](#).

Practical Aspects of Decision Trees

Handling Data

- **Standardization:** Not required, as decision trees are insensitive to the scale of features [\[4:18+source\]](#).
- **Encoding Categorical Features:** Necessary to encode categorical variables meaningfully, such as using target encoding.

Decision Tree Hyperparameters

- **Max Depth:** Controls the maximum depth of the tree to prevent overfitting.
- **Min Samples Split:** Minimum number of samples required to split a node.
- **Min Samples Leaf:** Minimum number of samples required to be at a leaf node [\[4:11+source\]](#).

Overfitting and Pruning

- Overfitting occurs if the tree is too complex and captures the noise in the dataset.
- **Pruning** is the technique to reduce overfitting by removing some sections of the tree that provide little power [\[4:13+source\]](#).

Feature Importance

- Decision trees can also be used to evaluate the importance of different features.



each split involving the feature [\[4:13+source\]](#) [\[4:19+source\]](#) .

Algorithm Overview

1. **Initialize:** Calculate the entropy of the entire dataset.
2. **Calculate Gain:** For each attribute, compute the gain of a potential split.
3. **Select Best Split:** Choose the split that maximizes the gain.
4. **Repeat:** Recursively apply the above steps to create sub-trees on each branch.
5. **Stop Conditions:** Maximum tree depth is reached, or further splits are not informative [\[4:7+source\]](#) [\[4:16+source\]](#) .

Example

Consider a dataset with features like age and mileage, where the target is whether a car is sold or not. By calculating information gain for different features and thresholds (e.g., age ≤ 5), the decision tree can be constructed to model the relationship between features and the target variable.

Conclusion

Decision Trees are robust yet simple models used in various domains for classification and regression problems. Understanding their theoretical underpinnings and practical considerations can significantly augment their application in solving real-world problems [\[4:19+source\]](#) .

This concludes the revision notes on decision trees. Make sure to practice coding a decision tree using libraries like `scikit-learn` to cement your understanding further.