

Decision Trees: Mathematical Foundations and Implementation

Albert Villanueva

July 9, 2025

Contents

1	Introduction to Decision Trees	2
2	Mathematical Foundations	2
2.1	Impurity Measures	2
2.1.1	Gini Impurity	2
2.1.2	Information Gain and Entropy	3
2.2	Splitting Criteria	4
3	Tree Construction Algorithm	4
3.1	Stopping Criteria	5
4	Overfitting and Regularization	5
4.1	Pre-pruning (Early Stopping)	6
4.2	Post-pruning	6
5	Custom Decision Tree Implementation	6
5.1	Class Structure Overview	6
5.2	Node Class Analysis	7
5.3	Main Classifier Methods	7
5.3.1	Initialization	7
5.3.2	Training Method	7
5.3.3	Tree Construction	7
5.3.4	Best Split Selection	8
5.3.5	Gini Impurity Calculation	8
5.3.6	Prediction Method	8
5.4	Key Implementation Features	8
5.5	Computational Complexity	9
6	Usage Example	9
7	Advantages and Limitations	9
7.1	Advantages	9
7.2	Limitations	9
8	Conclusion	10

1 Introduction to Decision Trees

Decision trees are supervised learning algorithms used for both classification and regression tasks. They create a model that predicts the value of a target variable by learning simple decision rules inferred from the data features. The resulting model is represented as a tree structure where:

- Internal nodes represent tests on attributes
- Branches represent outcomes of tests
- Leaf nodes represent class labels or values

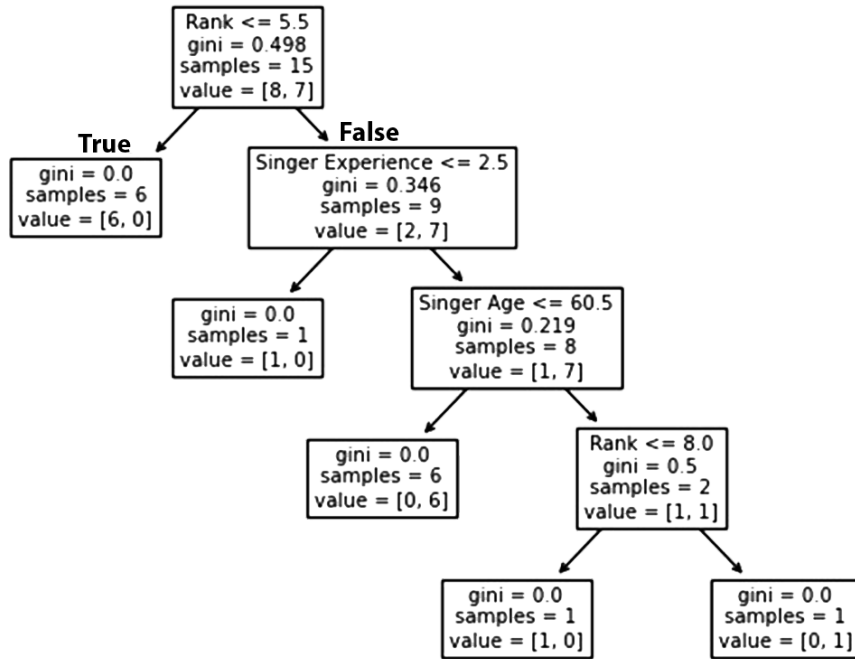


Figure 1: Example of a decision tree structure showing internal nodes with splitting conditions, branches representing outcomes, and leaf nodes with class predictions

2 Mathematical Foundations

2.1 Impurity Measures

The quality of a split in a decision tree is measured using impurity measures. The goal is to find splits that create the most homogeneous (pure) child nodes.

2.1.1 Gini Impurity

The Gini impurity measures the probability of incorrectly classifying a randomly chosen element if it were randomly labeled according to the distribution of labels in the subset.

For a node t with C classes, the Gini impurity is defined as:

$$\text{Gini}(t) = 1 - \sum_{i=1}^C p_i^2 \quad (1)$$

where p_i is the proportion of samples belonging to class i at node t .

For a binary split creating left and right child nodes, the weighted Gini impurity is:

$$\text{Gini}_{\text{split}} = \frac{n_{\text{left}}}{n} \cdot \text{Gini}(\text{left}) + \frac{n_{\text{right}}}{n} \cdot \text{Gini}(\text{right}) \quad (2)$$

where n_{left} and n_{right} are the number of samples in the left and right child nodes, and $n = n_{\text{left}} + n_{\text{right}}$.

2.1.2 Information Gain and Entropy

Entropy measures the average amount of information needed to identify the class of a sample:

$$\text{Entropy}(t) = - \sum_{i=1}^C p_i \log_2(p_i) \quad (3)$$

Information gain measures the reduction in entropy after a split:

$$\text{InfoGain}(S, A) = \text{Entropy}(S) - \sum_{v \in \text{Values}(A)} \frac{|S_v|}{|S|} \cdot \text{Entropy}(S_v) \quad (4)$$

where S is the set of samples, A is the attribute, and S_v is the subset of S for which attribute A has value v .

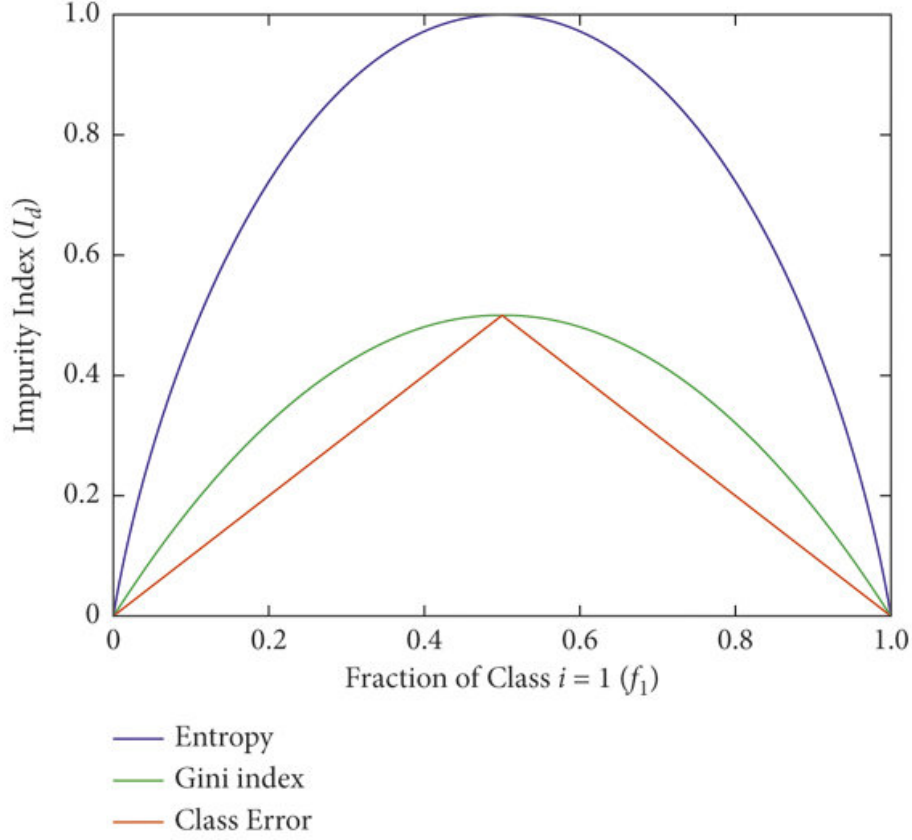


Figure 2: Comparison of Gini impurity and entropy for binary classification. Show curves plotting both measures against the probability of class 1, demonstrating their similar behavior with Gini being slightly faster to compute.

2.2 Splitting Criteria

For continuous features, we need to determine optimal threshold values. For a feature X_j and threshold θ , we create a binary split:

$$\text{Left child: } X_j \leq \theta \quad (5)$$

$$\text{Right child: } X_j > \theta \quad (6)$$

The optimal threshold θ^* minimizes the impurity measure:

$$\theta^* = \arg \min_{\theta} \text{Impurity}_{\text{split}}(\theta) \quad (7)$$

3 Tree Construction Algorithm

The decision tree construction follows a greedy, top-down approach:

Algorithm 1 Decision Tree Construction

```
1: function BuildTree(S, Features, depth)
2: if StoppingCriterion(S, depth) then
3:   return CreateLeaf(S)
4: end if
5: bestFeature, bestThreshold  $\leftarrow$  FindBestSplit(S, Features)
6: if bestFeature is None then
7:   return CreateLeaf(S)
8: end if
9:  $S_{\text{left}} \leftarrow \{(x, y) \in S : x_{\text{bestFeature}} \leq \text{bestThreshold}\}$ 
10:  $S_{\text{right}} \leftarrow \{(x, y) \in S : x_{\text{bestFeature}} > \text{bestThreshold}\}$ 
11: leftChild  $\leftarrow$  BuildTree( $S_{\text{left}}$ , Features, depth+1)
12: rightChild  $\leftarrow$  BuildTree( $S_{\text{right}}$ , Features, depth+1)
13: return CreateNode(bestFeature, bestThreshold, leftChild, rightChild)
```

3.1 Stopping Criteria

Common stopping criteria include:

- Maximum depth reached: $\text{depth} \geq \text{max_depth}$
- Minimum samples for split: $|S| < \text{min_samples_split}$
- Pure node: All samples belong to the same class
- No improvement in impurity measure

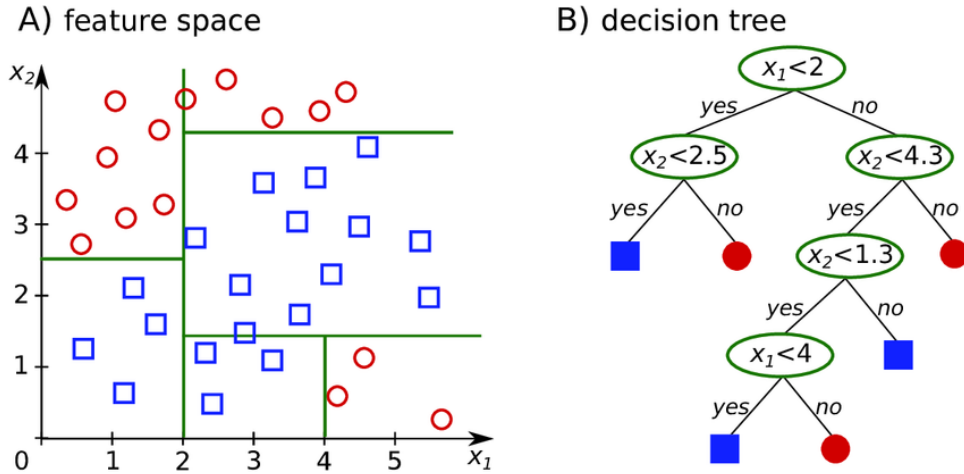


Figure 3: Step-by-step visualization of tree construction showing how the algorithm recursively splits the data space. Display the progressive partitioning of a 2D feature space with decision boundaries at each step.

4 Overfitting and Regularization

Decision trees are prone to overfitting, especially when grown deep. Common regularization techniques include:

4.1 Pre-pruning (Early Stopping)

- Limiting maximum depth
- Setting minimum samples per split/leaf
- Requiring minimum impurity decrease

4.2 Post-pruning

Post-pruning involves growing a full tree and then removing branches that don't provide significant improvement on validation data.

The cost-complexity pruning uses a complexity parameter α :

$$R_\alpha(T) = R(T) + \alpha|T| \quad (8)$$

where $R(T)$ is the error rate of tree T and $|T|$ is the number of leaves.

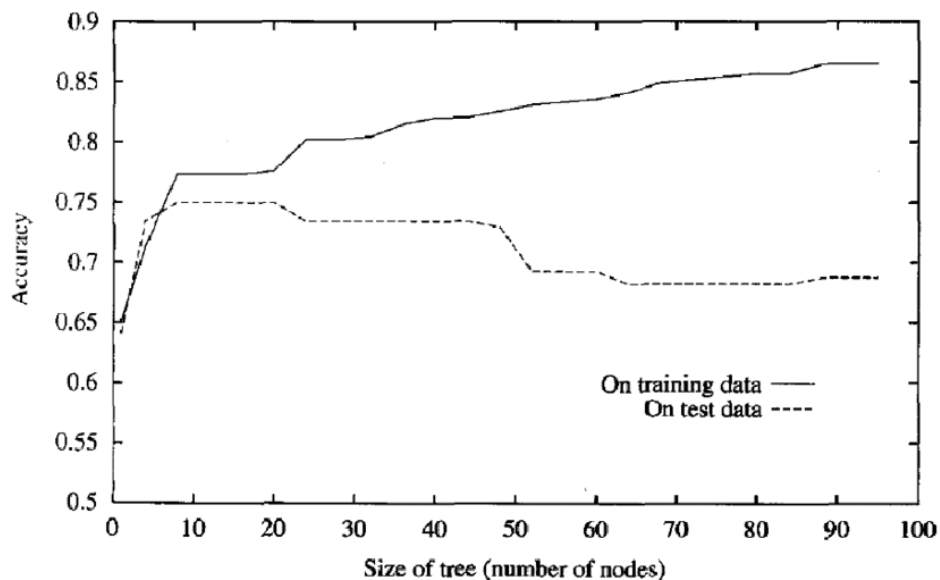


Figure 4: Illustration of overfitting in decision trees. Show training vs validation accuracy curves as tree depth increases, demonstrating how validation accuracy peaks and then decreases while training accuracy continues to improve.

5 Custom Decision Tree Implementation

5.1 Class Structure Overview

The custom implementation consists of two main classes:

- `CustomDecisionTreeNode`: Represents individual nodes in the tree
- `CustomDecisionTree`: The main classifier implementing the decision tree algorithm

5.2 Node Class Analysis

```
1 class CustomDecisionTreeNode:
2     def __init__(self, feature_index=None, threshold=None,
3                   left=None, right=None, value=None):
4         self.feature_index = feature_index # Feature used for
5                                             splitting
6         self.threshold = threshold         # Threshold value for split
7         self.left = left                   # Left child node
8         self.right = right                 # Right child node
9         self.value = value                 # Leaf node prediction
```

The node class implements a binary tree structure where:

- Internal nodes store `feature_index` and `threshold`
- Leaf nodes store the predicted value
- Navigation is handled through `left` and `right` pointers

5.3 Main Classifier Methods

5.3.1 Initialization

```
1 def __init__(self, max_depth=None, min_samples_split=2):
2     self.max_depth = max_depth
3     self.min_samples_split = min_samples_split
4     self.root = None
```

The constructor accepts hyperparameters for regularization:

- `max_depth`: Limits tree depth to prevent overfitting
- `min_samples_split`: Minimum samples required to create a split

5.3.2 Training Method

The `fit` method initiates the recursive tree building process:

```
1 def fit(self, X, y):
2     self.root = self._build_tree(X, y, depth=0)
```

5.3.3 Tree Construction

The `_build_tree` method implements the recursive algorithm:

$$\text{StoppingCondition} = \begin{cases} \text{True} & \text{if } \text{depth} \geq \text{max_depth} \\ \text{True} & \text{if } |\text{unique}(y)| = 1 \\ \text{True} & \text{if } |y| < \text{min_samples_split} \\ \text{False} & \text{otherwise} \end{cases} \quad (9)$$

5.3.4 Best Split Selection

The `_best_split` method finds the optimal feature and threshold by minimizing Gini impurity:

$$(\text{feature}^*, \theta^*) = \arg \min_{\text{feature}, \theta} \text{Gini}_{\text{split}}(\text{feature}, \theta) \quad (10)$$

$$\text{where } \theta \in \{\text{unique values of feature}\} \quad (11)$$

5.3.5 Gini Impurity Calculation

The implementation uses the weighted Gini impurity formula:

```
1 def _gini_index(self, left_y, right_y):
2     m_left, m_right = len(left_y), len(right_y)
3     m = m_left + m_right
4
5     gini_left = 1.0 - sum((np.sum(left_y == c) / m_left) ** 2
6                           for c in np.unique(left_y)) if m_left > 0 else
7                           0
8     gini_right = 1.0 - sum((np.sum(right_y == c) / m_right) ** 2
9                            for c in np.unique(right_y)) if m_right > 0
10                            else 0
11
12     return (m_left / m) * gini_left + (m_right / m) * gini_right
```

5.3.6 Prediction Method

Prediction traverses the tree from root to leaf:

Algorithm 2 Tree Traversal for Prediction

```
1: function Predict(input, node)
2:   if node.value is not None then
3:     return node.value
4:   end if
5:   if input[node.feature_index] ≤ node.threshold then
6:     return Predict(input, node.left)
7:   else
8:     return Predict(input, node.right)
9:   end if
```

5.4 Key Implementation Features

1. **Binary Splits Only:** Uses binary splits on continuous features
2. **Gini Impurity:** Uses Gini as the splitting criterion
3. **Greedy Search:** Evaluates all possible thresholds for each feature
4. **Majority Vote:** Leaf nodes predict the most common class
5. **Regularization:** Supports depth limiting and minimum sample constraints

5.5 Computational Complexity

- **Training Time:** $O(n \cdot m \cdot \log n)$ where n is samples and m is features
- **Prediction Time:** $O(\log n)$ for balanced trees, $O(n)$ worst case
- **Space Complexity:** $O(n)$ for storing the tree structure

6 Usage Example

```
1 import numpy as np
2 from sklearn.datasets import make_classification
3
4 # Generate sample data
5 X, y = make_classification(n_samples=1000, n_features=4,
6                           n_classes=2, random_state=42)
7
8 # Split data
9 from sklearn.model_selection import train_test_split
10 X_train, X_test, y_train, y_test = train_test_split(
11     X, y, test_size=0.2, random_state=42)
12
13 # Train custom decision tree
14 clf = CustomDecisionTree(max_depth=5, min_samples_split=10)
15 clf.fit(X_train, y_train)
16
17 # Make predictions
18 y_pred = clf.predict(X_test)
19
20 # Evaluate accuracy
21 accuracy = clf.score(X_test, y_test)
22 print(f"Accuracy: {accuracy:.3f}")
```

7 Advantages and Limitations

7.1 Advantages

- Easy to understand and interpret
- Requires little data preparation
- Handles both numerical and categorical data
- Can model non-linear relationships
- Automatically performs feature selection

7.2 Limitations

- Prone to overfitting
- Can create overly complex trees

- Biased toward features with more levels
- Can be unstable (small data changes can result in different trees)
- Difficulty in optimal tree construction (NP-complete)

Decision Boundaries

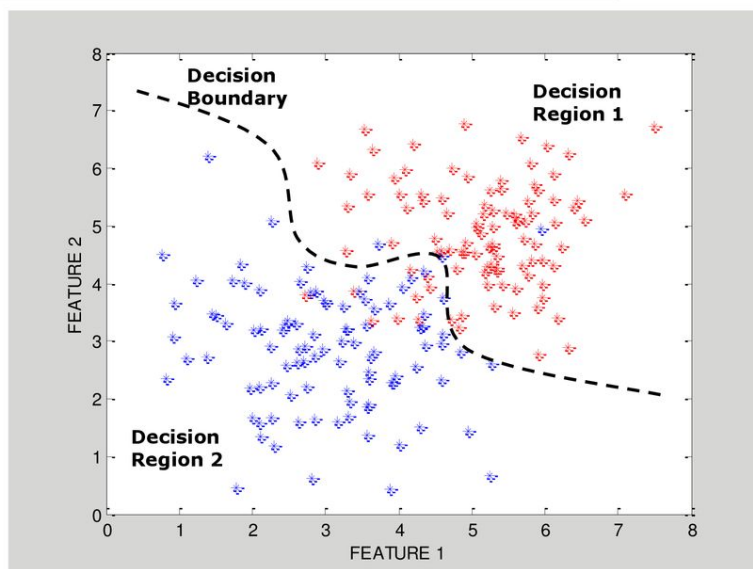


Figure 5: Decision boundaries created by a decision tree on a 2D dataset. Show rectangular regions corresponding to different classes, illustrating how decision trees partition the feature space with axis-parallel splits.

8 Conclusion

Decision trees provide an intuitive and powerful approach to classification problems. The custom implementation demonstrates the core concepts while maintaining simplicity and clarity. The mathematical foundations, particularly the Gini impurity measure and recursive splitting algorithm, form the backbone of more advanced ensemble methods like Random Forests and Gradient Boosting.

Understanding the implementation details helps in:

- Debugging and optimizing tree-based models
- Implementing custom splitting criteria
- Extending to ensemble methods
- Understanding hyperparameter effects

References

- [1] Breiman, L., Friedman, J., Olshen, R., & Stone, C. (1984). *Classification and Regression Trees*. Wadsworth & Brooks/Cole Advanced Books & Software.
- [2] Quinlan, J. R. (1986). Induction of decision trees. *Machine Learning*, 1(1), 81-106.
- [3] Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer.