

Data Mining Machine Learning

Md. Abdullah-Al-Kafi
Lecturer
Daffodil International University

Understanding Decision Trees

A decision tree is a type of algorithm used for making decisions based on data. Think of it like a flowchart where each step (or node) represents a decision point that leads you to the next step, eventually arriving at a final decision (or leaf). Here's a simple explanation:

Key Components of a Decision Tree

- **Root Node:** This is the topmost decision point in the tree. It represents the starting point where the first decision is made.
- **Decision Nodes:** These are points where decisions are made based on data features. Each node represents a question or test about the data.
- **Branches:** These are the possible outcomes of a decision node. Each branch leads to another decision node or a leaf node.
- **Leaf Nodes (Terminal Nodes):** These are the end points of the tree, representing the final decision or outcome after all the decisions have been made.

How Decision Trees Work

1. **Splitting the Data:** The algorithm starts at the root node and splits the data based on the feature that provides the best separation between the different classes (for classification) or outcomes (for regression). The “best” feature is usually determined by some metric like Gini impurity or information gain for classification tasks.
2. **Recursive Process:** The algorithm continues to split the data recursively at each decision node, creating branches, until one of the stopping criteria is met (like reaching a maximum depth, having too few samples to split further, or all the data in a node belonging to a single class).
3. **Final Decision:** When the algorithm reaches a leaf node, it assigns the most common class (for classification) or the average value (for regression) of the data points in that node as the final decision.

Example

Imagine you want to decide whether to play outside or stay indoors based on the weather. Your decision tree might look like this:

- **Root Node:** Is it sunny?
 - **If yes:** Is it warm?
 - * **If yes:** Go play outside! (Leaf Node)
 - * **If no:** Stay indoors. (Leaf Node)
 - **If no:** Is it raining?
 - * **If yes:** Stay indoors. (Leaf Node)
 - * **If no:** Is it windy?
 - **If yes:** Stay indoors. (Leaf Node)
 - **If no:** Go play outside! (Leaf Node)

Advantages and Disadvantages

- **Advantages:**

- Easy to understand and interpret.
- Can handle both numerical and categorical data.
- No need for feature scaling (normalization or standardization).

- **Disadvantages:**

- Can overfit the data, meaning it might learn noise instead of the actual pattern.
- Can become complex and less interpretable if the tree is too deep.

Decision trees are powerful tools for making decisions and are widely used in machine learning for tasks like classification and regression.

Decision Tree Example Using Gini Index

Suppose we have a dataset of 10 samples with two features, **Weather** (Sunny, Overcast, Rainy) and **Temperature** (Hot, Mild, Cool), and a target variable, **Play Tennis** (Yes, No). The dataset is as follows:

Sample	Weather	Temperature	Play Tennis
1	Sunny	Hot	No
2	Sunny	Hot	No
3	Overcast	Hot	Yes
4	Rainy	Mild	Yes
5	Rainy	Cool	Yes
6	Rainy	Cool	No
7	Overcast	Cool	Yes
8	Sunny	Mild	No
9	Sunny	Cool	Yes
10	Rainy	Mild	Yes

We want to determine which feature (Weather or Temperature) is the best first split using the Gini index.

Calculating the Gini Index

The **Gini index** measures the impurity of a dataset; a Gini index of 0 indicates a perfectly pure node (all samples belong to one class).

The formula for the Gini index is:

$$Gini = 1 - \sum (p_i^2)$$

where p_i is the proportion of samples that belong to class i .

Step-by-Step Calculation

1. **Calculate the Gini index for the overall dataset.**

There are 10 samples, with 6 "Yes" and 4 "No".

$$Gini_{parent} = 1 - \left(\frac{6}{10}\right)^2 - \left(\frac{4}{10}\right)^2$$

$$Gini_{parent} = 1 - 0.36 - 0.16 = 0.48$$

2. **Calculate the Gini index for splits on the feature "Weather".**

The "Weather" feature has three categories: Sunny, Overcast, and Rainy.

- **Sunny:** 4 samples (2 "No", 2 "Yes")

$$Gini_{Sunny} = 1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2 = 1 - 0.25 - 0.25 = 0.5$$

- **Overcast:** 2 samples (0 "No", 2 "Yes")

$$Gini_{Overcast} = 1 - \left(\frac{2}{2}\right)^2 - \left(\frac{0}{2}\right)^2 = 1 - 1 - 0 = 0$$

- **Rainy:** 4 samples (1 "No", 3 "Yes")

$$Gini_{Rainy} = 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2 = 1 - 0.5625 - 0.0625 = 0.375$$

3. Calculate the weighted Gini index for the "Weather" split.

$$Gini_{Weather} = \left(\frac{4}{10} \times 0.5\right) + \left(\frac{2}{10} \times 0\right) + \left(\frac{4}{10} \times 0.375\right)$$

$$Gini_{Weather} = 0.2 + 0 + 0.15 = 0.35$$

4. Calculate the Gini index for splits on the feature "Temperature".

The "Temperature" feature has three categories: Hot, Mild, and Cool.

- **Hot:** 3 samples (2 "No", 1 "Yes")

$$Gini_{Hot} = 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 = 1 - 0.444 - 0.111 = 0.444$$

- **Mild:** 3 samples (1 "No", 2 "Yes")

$$Gini_{Mild} = 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 = 0.444$$

- **Cool:** 4 samples (1 "No", 3 "Yes")

$$Gini_{Cool} = 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2 = 0.375$$

5. Calculate the weighted Gini index for the "Temperature" split.

$$Gini_{Temperature} = \left(\frac{3}{10} \times 0.444\right) + \left(\frac{3}{10} \times 0.444\right) + \left(\frac{4}{10} \times 0.375\right)$$

$$Gini_{Temperature} = 0.1332 + 0.1332 + 0.15 = 0.4164$$

Conclusion

The feature with the lower weighted Gini index is the better feature for splitting the data first. In this case:

$$Gini_{Weather} = 0.35, \quad Gini_{Temperature} = 0.4164$$

"Weather" has a lower Gini index (0.35), so it is the better feature to split on first.

Math Questions on Decision Trees and Gini Index

1. Calculate the Gini index for a dataset with two classes.

- Suppose you have a dataset with 8 samples: 5 are labeled "Yes" and 3 are labeled "No". Calculate the Gini index for this dataset.

$$\text{Answer: } Gini = 1 - \left(\frac{5}{8}\right)^2 - \left(\frac{3}{8}\right)^2 = 1 - 0.3906 - 0.1406 = 0.4688$$

2. Calculate the weighted Gini index for a split.

- Given a dataset split into two groups: Group 1 has 6 samples (4 "Yes", 2 "No") and Group 2 has 4 samples (2 "Yes", 2 "No"). Calculate the weighted Gini index for this split.

$$\text{Answer: } Gini_1 = 1 - \left(\frac{4}{6}\right)^2 - \left(\frac{2}{6}\right)^2 = 1 - 0.4444 - 0.1111 = 0.4445$$

$$Gini_2 = 1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2 = 1 - 0.25 - 0.25 = 0.5$$

$$\text{Weighted Gini} = \frac{6}{10} \times 0.4445 + \frac{4}{10} \times 0.5 = 0.2667 + 0.2 = 0.4667$$

3. Find the Gini index for a dataset with three classes.

- Suppose you have a dataset with 9 samples: 3 are labeled "A", 3 are labeled "B", and 3 are labeled "C". Calculate the Gini index for this dataset.

$$\text{Answer: } Gini = 1 - \left(\frac{3}{9}\right)^2 - \left(\frac{3}{9}\right)^2 - \left(\frac{3}{9}\right)^2 = 1 - 0.1111 - 0.1111 - 0.1111 = 0.6667$$

4. Compute the Gini index for a perfect split.

- A dataset has 10 samples, 5 labeled "Yes" and 5 labeled "No". If we perfectly split the dataset into two groups of 5 samples each, one group entirely "Yes" and the other entirely "No", what is the Gini index after the split?

$$\text{Answer: } Gini = 0 \text{ (since both groups are pure with only one class)}$$

5. Determine the Gini index for an imbalanced dataset.

- A dataset has 20 samples: 16 are labeled "Positive" and 4 are labeled "Negative". Calculate the Gini index for this dataset.

$$\text{Answer: } Gini = 1 - \left(\frac{16}{20}\right)^2 - \left(\frac{4}{20}\right)^2 = 1 - 0.64 - 0.04 = 0.32$$

Math Questions on Decision Trees Using Entropy

1. Calculate the entropy for a dataset with two classes.

- Suppose you have a dataset with 8 samples: 5 are labeled "Yes" and 3 are labeled "No". Calculate the entropy for this dataset.

$$\text{Entropy} = -\left(\frac{5}{8} \log_2 \frac{5}{8}\right) - \left(\frac{3}{8} \log_2 \frac{3}{8}\right)$$

$$\text{Entropy} = -(0.625 \times \log_2 0.625) - (0.375 \times \log_2 0.375)$$

$$\text{Entropy} = -(0.625 \times -0.678) - (0.375 \times -1.415) = 0.954$$

2. Calculate the weighted entropy for a split.

- Given a dataset split into two groups: Group 1 has 6 samples (4 "Yes", 2 "No") and Group 2 has 4 samples (2 "Yes", 2 "No"). Calculate the weighted entropy for this split.

$$\text{Entropy}_1 = -\left(\frac{4}{6} \log_2 \frac{4}{6}\right) - \left(\frac{2}{6} \log_2 \frac{2}{6}\right)$$

$$\text{Entropy}_1 = -(0.667 \times -0.585) - (0.333 \times -1.585)$$

$$\text{Entropy}_1 = 0.389 + 0.528 = 0.917$$

$$\text{Entropy}_2 = -\left(\frac{2}{4} \log_2 \frac{2}{4}\right) - \left(\frac{2}{4} \log_2 \frac{2}{4}\right)$$

$$\text{Entropy}_2 = -(0.5 \times -1) - (0.5 \times -1)$$

$$\text{Entropy}_2 = 0.5 + 0.5 = 1$$

$$\text{Weighted Entropy} = \frac{6}{10} \times 0.917 + \frac{4}{10} \times 1 = 0.55 + 0.4 = 0.95$$

3. Find the entropy for a dataset with three classes.

- Suppose you have a dataset with 9 samples: 3 are labeled "A", 3 are labeled "B", and 3 are labeled "C". Calculate the entropy for this dataset.

$$\text{Entropy} = -\left(\frac{3}{9} \log_2 \frac{3}{9}\right) - \left(\frac{3}{9} \log_2 \frac{3}{9}\right) - \left(\frac{3}{9} \log_2 \frac{3}{9}\right)$$

$$\text{Entropy} = -(0.333 \times -1.585) - (0.333 \times -1.585) - (0.333 \times -1.585)$$

$$\text{Entropy} = 0.528 + 0.528 + 0.528 = 1.584$$

4. Compute the entropy for a perfect split.

- A dataset has 10 samples, 5 labeled "Yes" and 5 labeled "No". If we perfectly split the dataset into two groups of 5 samples each, one group entirely "Yes" and the other entirely "No", what is the entropy after the split?

$$\text{Entropy} = 0 \text{ (since both groups are pure with only one class)}$$

5. Determine the entropy for an imbalanced dataset.

- A dataset has 20 samples: 16 are labeled "Positive" and 4 are labeled "Negative". Calculate the entropy for this dataset.

$$\text{Entropy} = -\left(\frac{16}{20} \log_2 \frac{16}{20}\right) - \left(\frac{4}{20} \log_2 \frac{4}{20}\right)$$

$$\text{Entropy} = -(0.8 \times -0.322) - (0.2 \times -2.322) = 0.257 + 0.464 = 0.721$$