

Decision Trees in Machine Learning

1 Introduction

Decision trees are a popular method for various machine learning tasks such as classification and regression. They are simple to understand and interpret, making them a useful tool for data analysis.

2 Mathematical Formulation

A decision tree is a flowchart-like structure where an internal node represents a feature (or attribute), the branch represents a decision rule, and each leaf node represents the outcome.

2.1 Entropy and Information Gain

Entropy is a measure of the impurity or randomness in the data. For a binary classification problem, the entropy H is defined as:

$$H(S) = -p_+ \log_2(p_+) - p_- \log_2(p_-)$$

where p_+ is the proportion of positive examples in set S and p_- is the proportion of negative examples in set S .

Information gain is used to decide which feature to split on at each step in building the tree. It is defined as the reduction in entropy:

$$IG(D, a) = H(D) - \sum_{i=1}^n \frac{|D_i|}{|D|} H(D_i)$$

where:

- D is the dataset.
- D_i are the partitions of D after splitting by a feature a .
- $|D|$ and $|D_i|$ are the number of instances in D and D_i , respectively.
- $H(D)$ is the entropy of the dataset D .

2.2 Gini Index

Another metric used for splitting is the Gini index, which measures the impurity of a dataset. For a binary classification problem, the Gini index G is defined as:

$$G(S) = 1 - p_+^2 - p_-^2$$

where p_+ and p_- are the proportions of positive and negative examples in set S , respectively.

3 Structure of Decision Trees

A decision tree consists of three types of nodes:

- **Root Node:** This is the topmost node in a decision tree. It represents the entire dataset and is split into two or more homogeneous sets.
- **Internal Nodes:** These nodes represent the features of the dataset and are used to make decisions based on the values of these features.
- **Leaf Nodes:** These nodes represent the final output or class label. They do not split further.

4 How Decision Trees Work

The process of building a decision tree involves the following steps:

1. **Select the Best Attribute:** Use a metric like information gain or Gini index to select the attribute that best separates the data.
2. **Split the Dataset:** Divide the dataset into subsets based on the selected attribute.
3. **Create Decision Nodes or Leaf Nodes:** If a subset is pure (i.e., all examples belong to the same class), create a leaf node. Otherwise, create an internal node and repeat the process for each subset.

5 Splitting Criteria

The choice of attribute to split on at each step is crucial for the performance of the decision tree. Common splitting criteria include:

- **Information Gain:** Measures the reduction in entropy after a dataset is split on an attribute.
- **Gini Index:** Measures the impurity of a dataset. A lower Gini index indicates a better split.

- **Chi-Square:** Measures the statistical significance of the differences between the observed and expected frequencies of the classes.
- **Reduction in Variance:** Used for regression trees, it measures the reduction in variance after a split.

6 Algorithm for Building a Decision Tree

The algorithm for building a decision tree can be summarized in the following steps:

Algorithm 1 Decision Tree Algorithm

Require: Dataset T , Attribute list A

Ensure: Decision Tree

```

1: function BuildTree( $T$ ,  $A$ )
2: if all examples in  $T$  belong to the same class then
3:   return a leaf node with that class
4: end if
5: if  $A$  is empty then
6:   return a leaf node with the majority class in  $T$ 
7: end if
8: Select the best attribute  $a$  from  $A$  using a splitting criterion (e.g., information gain, Gini index)
9: Create a root node for the tree with  $a$ 
10: for each value  $v$  of attribute  $a$  do
11:   Let  $T_v$  be the subset of  $T$  where attribute  $a$  has value  $v$ 
12:   if  $T_v$  is empty then
13:     Attach a leaf node with the majority class in  $T$  to the root node
14:   else
15:     Attach the subtree BuildTree( $T_v$ ,  $A \setminus \{a\}$ ) to the root node
16:   end if
17: end for
18: return root node

```

The algorithm starts by checking if all examples in the dataset belong to the same class. If so, it returns a leaf node with that class. If the attribute list is empty, it returns a leaf node with the majority class in the dataset. Otherwise, it selects the best attribute to split on, creates a root node, and recursively builds subtrees for each subset of the dataset.

7 Numerical Example

Let's consider a simple numerical example to illustrate how a decision tree is built. Suppose we have the following dataset of weather conditions and the decision to play tennis:

Outlook	Temperature	Humidity	Windy	PlayTennis
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	True	No

7.1 Step 1: Calculate Entropy of the Entire Dataset

The entropy of the entire dataset is calculated as follows:

$$H(S) = -p_+ \log_2(p_+) - p_- \log_2(p_-)$$

where p_+ is the proportion of positive examples (PlayTennis = Yes) and p_- is the proportion of negative examples (PlayTennis = No).

In our dataset, we have 9 "Yes" and 5 "No" examples:

$$p_+ = \frac{9}{14}, \quad p_- = \frac{5}{14}$$

$$H(S) = -\left(\frac{9}{14} \log_2 \frac{9}{14}\right) - \left(\frac{5}{14} \log_2 \frac{5}{14}\right) \approx 0.940$$

7.2 Step 2: Calculate Information Gain for Each Attribute

We calculate the information gain for each attribute to determine the best attribute to split on.

Outlook:

$$H(Sunny) = -\left(\frac{2}{5} \log_2 \frac{2}{5}\right) - \left(\frac{3}{5} \log_2 \frac{3}{5}\right) \approx 0.971$$

$$H(Overcast) = 0 \quad (\text{since all examples are "Yes"})$$

$$H(Rainy) = -\left(\frac{2}{5} \log_2 \frac{2}{5}\right) - \left(\frac{3}{5} \log_2 \frac{3}{5}\right) \approx 0.971$$

$$IG(S, \text{Outlook}) = H(S) - \left(\frac{5}{14} H(Sunny) + \frac{4}{14} H(Overcast) + \frac{5}{14} H(Rainy)\right) \approx 0.246$$

Temperature:

$$H(Hot) = - \left(\frac{2}{4} \log_2 \frac{2}{4} \right) - \left(\frac{2}{4} \log_2 \frac{2}{4} \right) = 1$$

$$H(Mild) = - \left(\frac{4}{6} \log_2 \frac{4}{6} \right) - \left(\frac{2}{6} \log_2 \frac{2}{6} \right) \approx 0.918$$

$$H(Cool) = - \left(\frac{3}{4} \log_2 \frac{3}{4} \right) - \left(\frac{1}{4} \log_2 \frac{1}{4} \right) = 0.811$$

$$IG(S, Temperature) = H(S) - \left(\frac{4}{14} H(Hot) + \frac{6}{14} H(Mild) + \frac{4}{14} H(Cool) \right) \approx 0.029$$

Humidity:

$$H(High) = - \left(\frac{3}{7} \log_2 \frac{3}{7} \right) - \left(\frac{4}{7} \log_2 \frac{4}{7} \right) \approx 0.985$$

$$H(Normal) = - \left(\frac{6}{7} \log_2 \frac{6}{7} \right) - \left(\frac{1}{7} \log_2 \frac{1}{7} \right) \approx 0.592$$

$$IG(S, Humidity) = H(S) - \left(\frac{7}{14} H(High) + \frac{7}{14} H(Normal) \right) \approx 0.151$$

Windy:

$$H(False) = - \left(\frac{6}{8} \log_2 \frac{6}{8} \right) - \left(\frac{2}{8} \log_2 \frac{2}{8} \right) = 0.811$$

$$H(True) = - \left(\frac{3}{6} \log_2 \frac{3}{6} \right) - \left(\frac{3}{6} \log_2 \frac{3}{6} \right) = 1$$

$$IG(S, Windy) = H(S) - \left(\frac{8}{14} H(False) + \frac{6}{14} H(True) \right) \approx 0.048$$

From the calculations, we see that the attribute "Outlook" has the highest information gain. Therefore, we split the dataset based on the "Outlook" attribute.

7.3 Step 3: Repeat the Process for Each Subset

We repeat the process for each subset of the dataset created by the "Outlook" split. This involves calculating the entropy and information gain for each remaining attribute within each subset and selecting the best attribute to split on.

By following these steps recursively, we can build the entire decision tree.

.

.

.

.

.

.

.

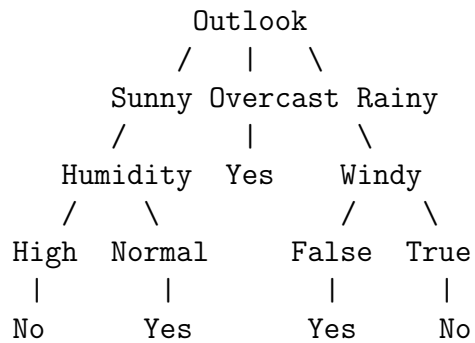
.

.

.

7.4 Resulting Decision Tree

The resulting decision tree for our example dataset is as follows:



8 Conclusion

Decision trees are a fundamental machine learning technique that can be used for both classification and regression tasks. Understanding the mathematical foundations, such as entropy, information gain, and the Gini index, is crucial for effectively implementing and interpreting decision trees.