# Decision Tree Machine Learning Methods

## 1 Introduction

Decision trees are a popular machine learning method for classification and regression problems. They are easy to understand and interpret, and can be used for both categorical and continuous input and output variables.

In this handout, we will cover the basics of decision tree machine learning methods, including how they work, how to build them, and how to evaluate them.

## 2 Decision Tree Algorithm

The decision tree algorithm works by recursively partitioning the data into subsets based on the values of the input variables. The goal is to create a tree where each internal node represents a decision based on one of the input variables, and each leaf node represents a class or a value for the output variable.

The algorithm uses a greedy approach to select the best variable and value to split the data at each node. The goal is to maximize the information gain, which measures how much the entropy (or impurity) of the data decreases after the split.

### 2.1 Greedy Approach

In decision tree learning, a greedy approach is used to build the tree by selecting the best feature to split on at each node. A greedy algorithm makes the locally optimal choice at each step in the hope of finding a global optimum.

The greedy approach works by selecting the feature that provides the most information gain or reduces the most entropy at each node. This means that the algorithm chooses the feature that best separates the data into different classes.

The downside of the greedy approach is that it may not always lead to the optimal decision tree. By making locally optimal choices at each step, the algorithm may miss a better overall solution. However, the greedy approach is computationally efficient and works well in practice.

In summary, the greedy approach is a decision-making strategy that selects the best option at each step in the hope of finding the best overall solution. In decision tree learning, the greedy approach is used to build the tree by selecting the best feature to split on at each node.

The entropy of a binary classification problem can be calculated as:

$$H(X) = -p_1 \log_2 p_1 - p_2 \log_2 p_2 \tag{1}$$

where $p_1$ and $p_2$ are the proportions of positive and negative examples in the data.

# 3  Entropy in Decision Trees

Entropy is a measure of impurity or uncertainty in a set of examples. In decision trees, entropy is used to determine the best split at each node. A split is considered good if it reduces the entropy of the child nodes compared to the parent node.

## 3.1  What is Entropy?

Entropy is a concept from information theory, which was introduced by Claude Shannon in the late 1940s. Entropy measures the amount of uncertainty or randomness in a system. In decision tree learning, entropy is used to measure the impurity of a set of examples.

The entropy of a binary system can be calculated using the following formula:

$$H(X) = -p_1 \log_2(p_1) - p_2 \log_2(p_2) \tag{2}$$

where $X$ is the set of examples, $p_1$ is the proportion of examples in class 1, and $p_2$ is the proportion of examples in class 2. The logarithm base 2 is used to measure entropy in bits.

## 3.2  Breaking Down the Formula

The entropy formula can be broken down into two parts: the probability of each class and the logarithm of the probability.

The probability of each class is represented by $p_1$ and $p_2$. These probabilities are calculated by dividing the number of examples in each class by the total number of examples.

The logarithm of the probability is represented by $\log_2(p_1)$ and $\log_2(p_2)$. The logarithm is used to measure the amount of uncertainty or randomness in the system. If $p_i$ is close to 0, then $\log_2(p_i)$ will be a large negative number. This means that the entropy will be high, indicating a high level of uncertainty or randomness.

The negative sign in the formula is used to ensure that entropy is always positive. The maximum entropy occurs when $p_1$ and $p_2$ are equal, which means that the system is in a state of maximum uncertainty or randomness.

### 3.2.1 Using Entropy in Decision Trees

In decision trees, entropy is used to measure the impurity of a set of examples. The goal is to find the best split that reduces the entropy of the child nodes compared to the parent node. The split with the highest information gain, which is the difference between the parent node entropy and the weighted average of the child node entropies, is chosen as the best split.

Entropy is a powerful tool in decision tree learning because it can handle both continuous and categorical data. It is also computationally efficient and works well in practice.

In conclusion, entropy is a measure of impurity or uncertainty in a set of examples. It is calculated using the probability of each class and the logarithm of the probability. In decision trees, entropy is used to determine the best split at each node. The split with the highest information gain is chosen as the best split.

# 4 Information Gain in Decision Trees

Information gain is a measure of the reduction in entropy achieved by splitting a set of examples on a particular feature. In decision tree learning, information gain is used to determine the best feature to split on at each node.

## 4.1 What is Information Gain?

Information gain is the difference between the entropy of the parent node and the weighted average of the child node entropies. It measures how much information is gained by splitting the data on a particular feature.

The information gain formula can be expressed as:

$$IG(S, F) = H(S) - \sum_{v \in Values(F)} \frac{|S_v|}{|S|} H(S_v) \tag{3}$$

where $S$ is the set of examples, $F$ is the feature being split on, $Values(F)$ is the set of possible values for feature $F$, $S_v$ is the subset of examples in $S$ that have value $v$ for feature $F$, and $H(S)$ and $H(S_v)$ are the entropy of $S$ and $S_v$, respectively.

## 4.2 Breaking Down the Formula

The information gain formula can be broken down into two parts: the entropy of the parent node and the weighted average of the child node entropies.

The entropy of the parent node, $H(S)$, measures the impurity of the set of examples before the split. The goal is to find the feature that will create child nodes with the lowest possible entropy, which will reduce the overall entropy of the tree.

The weighted average of the child node entropies, $\sum_{v \in Values(F)} \frac{|S_v|}{|S|} H(S_v)$, measures the impurity of the child nodes after the split. The weighted average takes into account the size of each child node, so larger child nodes have a greater impact on the overall entropy of the tree.

The information gain, $IG(S, F)$, is the difference between the entropy of the parent node and the weighted average of the child node entropies. The feature with the highest information gain is chosen as the best feature to split on.

## 4.3   Using Information Gain in Decision Trees

In decision tree learning, information gain is used to determine the best feature to split on at each node. The feature with the highest information gain is chosen as the best feature to split on, as it will create child nodes with the lowest possible entropy.

Information gain is a powerful tool in decision tree learning because it can handle both continuous and categorical data. However, it can be biased towards features with many values or high cardinality, as these features tend to have higher information gain simply because they create more splits.

In conclusion, information gain is a measure of the reduction in entropy achieved by splitting a set of examples on a particular feature. It is calculated as the difference between the entropy of the parent node and the weighted average of the child node entropies. In decision tree learning, information gain is used to determine the best feature to split on at each node.

# 5   Decision Tree Example with NumPy and scikit-learn

Let's illustrate the decision tree algorithm with an example using the scikit-learn library.

First, let's load a dataset and split it into training and testing sets:

Listing 1: Load and split dataset

```
import numpy as np
from sklearn.datasets import load_iris
from sklearn.model_selection import train_test_split

iris = load_iris()
X_train, X_test, y_train, y_test = train_test_split(
    iris.data,
    iris.target,
    test_size=0.2,
    random_state=42)
```

Next, let's create a decision tree classifier and fit it to the training data:

Listing 2: Create and fit decision tree classifier

```python
from sklearn.tree import DecisionTreeClassifier

clf = DecisionTreeClassifier(max_depth=2, random_state=42)
clf.fit(X_train, y_train)
```

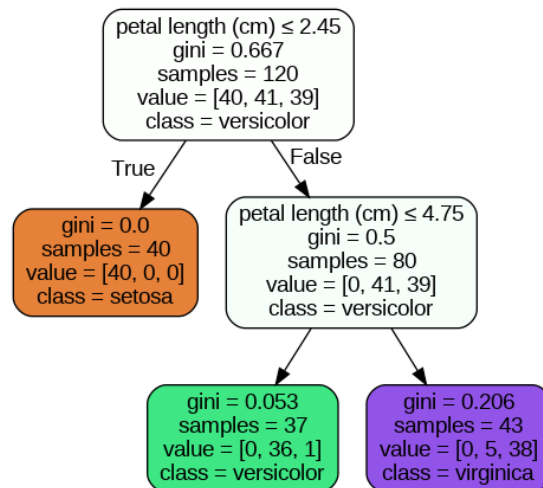We can visualize the decision tree using the Graphviz library:

Listing 3: Visualize decision tree

```python
from sklearn.tree import export_graphviz
import graphviz

dot_data = export_graphviz(clf,
    out_file=None,
    feature_names=iris.feature_names,
    class_names=iris.target_names,
    filled=True, rounded=True,
    special_characters=True)

graph = graphviz.Source(dot_data)
graph
```

The resulting tree should look something like this:



Finally, we can evaluate the performance of the decision tree on the testing data:

Listing 4: Evaluate decision tree

```python
from sklearn.metrics import accuracy_score
```

```
y_pred = clf.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy:", accuracy)
```

For this example, we should get an accuracy of around 0.9666.

## 5.1 Using Decision Trees for Regression

Decision trees can also be used for regression tasks, where the goal is to predict a continuous target variable. In this case, the decision tree algorithm works by recursively partitioning the feature space into regions that have relatively homogeneous values of the target variable. The resulting tree can then be used to make predictions for new data points by traversing the tree from the root to a leaf node and outputting the average target value of the training data points that fall in that region.

### 5.1.1 Different Ways to Work with Regression

There are several ways to work with decision trees for regression tasks, including:

- **Standard Decision Trees:** A standard decision tree is built using the same algorithm as for classification tasks, but instead of minimizing the Gini index or entropy, it minimizes the variance of the target variable within each region. This is done by splitting the data at each node in such a way as to minimize the sum of the squared differences between the target values and their mean within each child node. The resulting tree can then be pruned to reduce overfitting. The variance of the target variable within a node $m$ can be calculated as follows:

$$\mathrm{Var}(y_m) = \frac{1}{N_m} \sum_{i \in \mathrm{samples}_m} (y_i - \bar{y}_m)^2, \tag{4}$$

where $N_m$ is the number of training samples in node $m$, $\mathrm{samples}_m$ is the set of training samples in node $m$, $y_i$ is the target value of the $i$-th training sample, and $\bar{y}_m$ is the mean target value of the training samples in node $m$.

The splitting criterion for a node $m$ can then be defined as follows:

$$\mathrm{argmin}_{j,t} \left[ \mathrm{Var}(y_{\mathrm{left}}) + \mathrm{Var}(y_{\mathrm{right}}) \right], \tag{5}$$

where $j$ is the index of the feature to split on, $t$ is the threshold value for the split, $y_{\mathrm{left}}$ is the target values of the training samples in the left child node, and $y_{\mathrm{right}}$ is the target values of the training samples in the right child node.

- **Gradient Boosted Trees:** Gradient boosted trees (GBTs) are an ensemble learning method that combines multiple decision trees to improve prediction accuracy. In GBTs for regression, each tree is built to predict the residuals of the previous tree, and the final prediction is the sum of the predictions of all trees. GBTs can be tuned using hyperparameters such as the learning rate, the number of trees, and the maximum depth of each tree. The prediction of a GBT model with $T$ trees can be expressed as follows:

$$f(x) = \sum_{t=1}^{T} \gamma_t h_t(x), \tag{6}$$

where $\gamma_t$ is the learning rate for tree $t$, $h_t(x)$ is the prediction of tree $t$ for input $x$, and $x$ is a vector of input features.

### 5.1.2 Regression Metrics

regression metrics can be used to evaluate the performance of decision tree models for regression. Some commonly used metrics include:

- **Mean Squared Error (MSE):** The mean squared error is the average of the squared differences between the predicted and actual target values. It measures the average squared deviation of the predictions from the true values and is commonly used as a loss function for training decision tree models for regression.

  The MSE can be calculated as follows:

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2, \tag{7}$$

where $n$ is the number of samples, $y_i$ is the true target value for the $i$-th sample, and $\hat{y}_i$ is the predicted target value for the $i$-th sample.

- **Mean Absolute Error (MAE):** The mean absolute error is the average of the absolute differences between the predicted and actual target values. It measures the average absolute deviation of the predictions from the true values and is less sensitive to outliers than the MSE.

  The MAE can be calculated as follows:

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^{n} |y_i - \hat{y}_i|. \tag{8}$$

- **R-squared ($R^2$):** The R-squared is a measure of how well the model fits the data. It represents the proportion of variance in the target variable that is explained by the model, with higher values indicating better model

fit. R-squared can range from 0 to 1, with 1 indicating a perfect fit and 0 indicating that the model does not explain any of the variance in the target variable.

The R-squared can be calculated as follows:

$$R^2 = 1 - \frac{\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}{\sum_{i=1}^{n}(y_i - \bar{y})^2},\tag{9}$$

where $\bar{y}$ is the mean target value of the training data.

When working with decision trees for regression, it is important to choose an appropriate metric based on the specific problem and data at hand. For example, if the target variable has a large range of values or contains outliers, the MAE or R-squared may be more appropriate than the MSE. Similarly, if the goal is to explain the variation in the target variable, the R-squared may be a better metric to use.

# 6    Conclusion

In this handout, we have covered the basics of decision tree machine learning methods, including how they work, how to build them, and how to evaluate them. We have also provided an example implementation using the scikit-learn library.

Decision trees are a powerful and versatile machine learning method that can be used for a wide range of classification and regression problems. With their intuitive visual representation and ability to handle both categorical and continuous variables, they are a valuable tool in the data scientist's toolbox

## 6.1    Boosted Trees

We have explored the concept of boosted gradient trees, a powerful machine learning technique that combines decision trees with gradient boosting. Boosted trees are able to model complex relationships between input features and target variables, and can achieve high predictive accuracy on a wide range of datasets. By iteratively training weak decision trees on the residual errors of the previous tree, boosted trees are able to gradually improve their predictions and reduce errors.

The residual error is the difference between the predicted values and the actual values in a regression problem. In other words, it represents the amount of error that is left over after the model has made its best attempt to fit the data.

Mathematically, the residual error for a given data point can be calculated as follows:

$$e_i = y_i - \hat{y}_i \tag{10}$$

where $y_i$ is the actual value for the $i$-th data point, and $\hat{y}_i$ is the predicted value. The residual error can be positive or negative, depending on whether the prediction overestimates or underestimates the actual value.

In the context of gradient boosting, the next model in the sequence is trained on the residual errors of the previous model. This allows the subsequent models to focus on the parts of the data that the previous models did not perform well on. By iteratively fitting models to the residual errors, the final ensemble of models can achieve better performance than any individual model in the sequence.

Boosted trees have found many applications in the real world, including:

Predicting customer churn in the telecommunications industry Predicting customer response to marketing campaigns Predicting loan defaults in the finance industry Identifying fraudulent activity in financial transactions Compared to traditional decision trees, boosted trees are less prone to overfitting and can achieve better performance on complex datasets.

Boosted trees have numerous real-world applications, including:

- Fraud detection: Boosted trees can be used to identify fraudulent transactions or activities by analyzing patterns in large datasets.

- Recommender systems: Boosted trees can be used to make personalized recommendations for products, movies, or music based on a user's past behavior and preferences.

- Medical diagnosis: Boosted trees can be used to analyze medical data and help doctors diagnose diseases or conditions.

- Financial modeling: Boosted trees can be used to analyze financial data and make predictions about stock prices, credit risk, or loan defaults.

Compared to other machine learning algorithms, boosted trees have several advantages, including:

- High accuracy: Boosted trees can achieve high predictive accuracy on a wide range of datasets, even when the relationships between input features and target variables are complex and nonlinear.

- Robustness: Boosted trees are relatively robust to noisy or incomplete data, and can handle missing values or outliers without requiring extensive data preprocessing.

- Flexibility: Boosted trees can be used for both regression and classification problems, and can handle both numerical and categorical data.

- Interpretabilty: Boosted trees can provide insights into the relative importance of different input features, making them useful for feature selection or understanding the underlying relationships in the data.

Overall, boosted gradient trees are a powerful and versatile machine learning technique that can be used in a wide range of applications. By combining decision trees with gradient boosting, boosted trees are able to model complex relationships and achieve high predictive accuracy, making them a valuable tool for data scientists and machine learning practitioners.