# Efficiency of the Decision Tree Algorithm in Data Analysis

**Xusenov Shoxrux Sherali o'g'li**

Tashkent University of Information Technologies named after Muhammad

al-Kharazmi Faculty of Software Engineering, 4th-year student

xusenovshoxrux7@gmail.com

**Abstract:** This article analyzes the efficiency and application possibilities of the Decision Tree algorithm in the process of data analysis. Nowadays, processing large volumes of data and extracting useful conclusions from them has become one of the most important tasks. The Decision Tree algorithm is one of the widely used methods in the field of Machine Learning, and it provides effective results in classification and prediction tasks. The article discusses the working principles, main characteristics, advantages, and practical application areas of the Decision Tree algorithm.

**Keywords:** Decision Tree, data analysis, machine learning, classification, regression, algorithm, artificial intelligence.

**Introduction**

In recent years, the rapid development of information technologies has led to the emergence of large volumes of data. The process of analyzing these data and extracting useful knowledge from them has contributed to the development of data mining and machine learning fields.

With the help of machine learning algorithms, it is possible to classify, predict, and analyze data. Among such algorithms, the Decision Tree algorithm stands out due to its simple structure, understandable working principle, and high efficiency.

The Decision Tree algorithm represents data in the form of a tree structure. Each node of the tree checks a specific condition, and the branches determine the next decision

based on the result of that condition. This approach is very similar to the human decision-making process and makes the results easier to understand.

Today, the Decision Tree algorithm is widely used in many fields such as finance, medicine, marketing, business analytics, and other areas for data analysis.

### *Concept of the Decision Tree Algorithm*

The Decision Tree is one of the machine learning algorithms used for classification or prediction of data. This algorithm represents data in the form of a tree structure.

A Decision Tree consists of the following elements:

- Root node – the starting point of the tree that divides the data based on the first condition.
- Decision node – a node that divides the data based on a specific condition.
- Leaf node – a node that represents the final outcome.

During the operation of the algorithm, the dataset is divided several times, resulting in the formation of a decision tree. Each split aims to classify the data as accurately as possible.

The Decision Tree algorithm is mainly divided into two types:

1. Classification Tree – used for classifying data into categories.
2. Regression Tree – used for predicting numerical values.

### *Working Principle of the Decision Tree Algorithm*

The Decision Tree algorithm works by splitting data step by step. The algorithm divides the dataset based on different features and selects the best splitting criterion.

The working process of the algorithm includes the following stages:

1. All features in the dataset are analyzed.
2. The most important feature is selected.
3. The data are divided based on this feature.
4. The process is repeated for each subset.
5. The tree is formed until the data are fully classified.

To determine the best split, the algorithm uses special mathematical criteria. The most commonly used criteria include the following:

Entropy represents the level of disorder in the data. If the data belong to the same class, the entropy value will be low.

Information Gain indicates how much the entropy decreases after splitting the data. The feature with the highest Information Gain value is selected for splitting.

The Gini index is also used to evaluate the quality of data splitting. This criterion shows how mixed the data are.

### *Advantages of the Decision Tree Algorithm*

The Decision Tree algorithm is widely used in machine learning due to several important advantages.

Understandable and visual model – The Decision Tree algorithm is represented in the form of a tree. This makes the results easier to understand. Even users without deep knowledge of machine learning can interpret the model.

Low requirement for data preprocessing – Many machine learning algorithms require data normalization or standardization. However, the Decision Tree algorithm usually does not require such preprocessing.

Ability to work with different types of data – The Decision Tree algorithm can work effectively with both numerical and categorical data.

Flexibility – The Decision Tree algorithm can be used to solve different types of problems. It works efficiently in both classification and regression tasks.

### *Disadvantages of the Decision Tree Algorithm*

Like any algorithm, the Decision Tree also has some disadvantages.

First, the algorithm may suffer from the problem of overfitting. If the tree becomes too deep, the model may fit the training data too closely and may not perform well on new data.

Second, even small changes in the data may significantly affect the structure of the tree. This can reduce the stability of the model.

To reduce these problems, pruning techniques are used to simplify the tree.

### *Practical Applications of the Decision Tree Algorithm*

The Decision Tree algorithm is successfully applied in many fields.

In the medical field, the Decision Tree algorithm is used to identify diseases and assist in diagnosis. For example, the probability of a disease can be determined based on a patient's symptoms.

In banking systems, the Decision Tree algorithm is used to evaluate credit risk. Based on a customer's financial information, the bank can decide whether to grant a loan or not.

In marketing, the Decision Tree algorithm is used to analyze customers' purchasing behavior.

In the IT field, the Decision Tree algorithm is used to analyze user behavior and to develop recommendation systems.

### Efficiency of the Decision Tree Algorithm

The efficiency of the Decision Tree algorithm depends on several factors:

- quality of the data
- size of the dataset
- chosen splitting criterion
- depth of the tree

A properly tuned Decision Tree model can achieve a high level of accuracy. In addition, it works faster and requires fewer computational resources compared to many other complex algorithms.

In many machine learning systems, the Decision Tree algorithm is used either as a primary model or in combination with other algorithms. For example, algorithms such as Random Forest and Gradient Boosting are also based on the Decision Tree concept.

### Conclusion

Data analysis is an important component of modern information systems. The Decision Tree algorithm is one of the effective methods widely used in the field of machine learning.

This algorithm allows data to be analyzed in a tree structure, enabling classification and prediction. The simple structure, interpretability, and high efficiency of the Decision Tree algorithm make it applicable in many fields.

However, the algorithm also has some limitations, such as overfitting, which can be reduced using specific techniques. With the further development of machine learning technologies, new and more efficient methods based on the Decision Tree algorithm are expected to emerge in the future.

## References:

1. Han J., Kamber M., Pei J. *Data Mining: Concepts and Techniques.* Morgan Kaufmann, 2012.

2. Mitchell T. *Machine Learning.* McGraw-Hill Education, 1997.

3. Witten I., Frank E., Hall M. *Data Mining: Practical Machine Learning Tools and Techniques.* Morgan Kaufmann, 2016.

4. Géron A. *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow.* O'Reilly Media, 2019.

5. Bishop C. *Pattern Recognition and Machine Learning.* Springer, 2006.

6. Tan P., Steinbach M., Kumar V. *Introduction to Data Mining.* Pearson Education, 2014.

7. Kotu V., Deshpande B. *Data Science: Concepts and Practice.* Morgan Kaufmann, 2018.

8. James G., Witten D., Hastie T., Tibshirani R. *An Introduction to Statistical Learning.* Springer, 2013.