# TESTING DECISION TREE ALGORITHM USING TANAGRA APPLICATIONS

**Soni Bahagia Sinaga**
Pogram studi Teknik Informatika
AMIK Stiekom Sumatera Utara Indonesia, 20154
bahagiainaga86@gmail.com

**Abstract**
Decision tree is one algorithm that is used to classify segmentation or grouping which is predictive, Decision Tree Algorithm has the Advantages in processing numerical (continuous) and discrete data, can handle missing attribute values, produces rules that are easily interpreted and the fastest among other algorithms. Prediction accuracy is the ability of the model to be able to predict class labels against new or previously unknown data well. In terms of speed or computational time efficiency needed to create and use a model. The application used is Tanagra because the application is available for Decision tree architecture.
Keywords: Decision Tree, Data Mining, tanagra

## 1. Introduction

With the development of information technology today, the need for accurate information is needed in daily life, so that information will become an important element in the development of society today and the future. However, high information needs are sometimes not balanced by the presentation of adequate information, often the information must still be reviewed from very large amounts of data. The ability of information technology to collect and store various types of data far leaves the ability to analyze, summarize and extract knowledge from data. The traditional method for analyzing existing data, cannot handle large amounts of data.

Utilizing existing data in information systems to support activities taking activities, it is not enough to rely solely on operational data, we need a data analysis to explore the potential of existing information. Decision-makers try to utilize the data warehouse that is already owned to dig up useful information to help retrieve the data needed, this encourages the emergence of new branches of science to overcome the problem of extracting information or patterns that are important or interesting from large amounts of data, called data mining. The use of data mining techniques is expected to provide knowledge previously hidden in the data warehouse so that it becomes valuable information.Decision Tree Method is a method for determining the main factors based on a comparison between one factor with another. In determining a graduation level, of course various considerations need attention.

## 2. Literature Review
### 2.1 Data Mining

Data Mining is a process that employs one or more computer learning techniques to analyze and extract knowledge automatically or a series of processes to explore the added value of a data set in the form of knowledge that has not been known manually. It is worth remembering that the word mining itself means an attempt to get a few valuable items from a large amount of basic material. Therefore Data Mining has long roots in the fields of science such as artificial intelligence, machine learning, statistics, and databases. Data mining is the process of applying this method to data to uncover hidden patterns. In other words Data mining is a process for extracting patterns from data. Data mining is becoming an increasingly important tool for turning that data into information. This is often used in various profile practices, such as marketing, surveillance, fraud detection, and scientific discovery. It has been used for years by businesses, scientists,

and governments to filter out volumes of data such as flight passenger travel records, census data, and supermarket scanner data to produce a market research report. [1]–[3]

The main reason for using data mining is to assist in the analysis of a collection of observational behaviors. The data is susceptible to collinearity because a link is known. The inevitable fact of data mining is that the subsets of data sets analyzed may not represent the entire domain, and therefore may not contain examples of certain critical relationships and behaviors that exist in other parts of the domain. To overcome this kind of problem, analysis can be augmented using trial-based and other approaches, such as Choice Modeling for human-generated data. In this situation, inherent can be controlled correlation to, or removed altogether, during design construction. [4], [5]

## 2.2 Decision Tree

Decision tree is a classification method that uses tree structure representations where each node represents an attribute, its branches represent the values of the attributes, and leaves represent the class. The top node of the decision tree is called root. A decision tree is the most popular classification method used. Besides being relatively fast construction, the results of the model built are easy to understand

In the decision tree there are 3 types of nodes, namely: [6]–[8]

1. Root Node, is the top node, at this node there is no input and can have no output or have more than one output.
2. Internal Node, is a branching node, at this node there is only one input and has a minimum output of two.
3. Leaf node or terminal node is the final node, at this node, there is only one input and has no output

## 2.3 Algorithm C 4.5

C4.5 algorithm is one method for making decision trees based on the training data provided. Some of the developments that have been carried out at C45 include being able to overcome missing values, being able to overcome continuous data, and pruning. Here is the basic algorithm from C4.5: [9]–[12]

1. Build a decision tree from the training set
2. Pruning to simplify the tree.
3. Change the tree generated in several rules. The number of rules is equal to the number of paths that might be built from the root to the leaf node.

To select the attribute as the root, based on the highest gain value of the existing attributes. To calculate the gain the formula is used as stated in formula 1:

$$\text{Gain(S,A)} = \text{Entropy(S)} - \sum_{i-}^{n} \frac{S_i}{S} \text{xEntropy}(S_i) \dots\dots\dots\dots\dots\dots\dots\dots(1)$$

With:
S: The set of cases
A: Attributes
n: Number of attribute attributes A
[Si]: number of cases on partition i
[S]: number of S cases

While the entropy value calculation can be seen in the following formula 2:

$$\text{Entropy(S)} = \sum_{i-1}^{a} -pi * log_2 pi \dots\dots\dots\dots\dots\dots\dots\dots\dots(2)$$

With:
S: Case Set
a: Feature
I: Number of partitions S
Pi: Proportion from Si to S

## 3. Results and Discussion

Software testing is an investigation conducted to obtain information about the quality of the product or service being tested. Software testing also provides an objective and independent perspective on software, which is useful in business operations to understand the level of risk in its implementation. Testing techniques include, but are not limited to the process of executing a part of the program or the entire application to find errors or other defects. The testing in Tanagra 1.4 software can be done in a computer device.

1. *Infut File*

Setelah semua kebutuhan *hardware* dan s*oftwa*re dilakukan, selanjutnya melakukan penginputan nilai data siswa kedalam Microsoft Excel.



Figure 1 Student data in Entropy calculations

After the student data is distorted in Microsoft Excel, then saved in Text (Tab Delimited) format.

2. Tanagra Initial Display 1.4

The display of the Tanagra software that has not been inputted with data or is empty can be seen in the following image display:
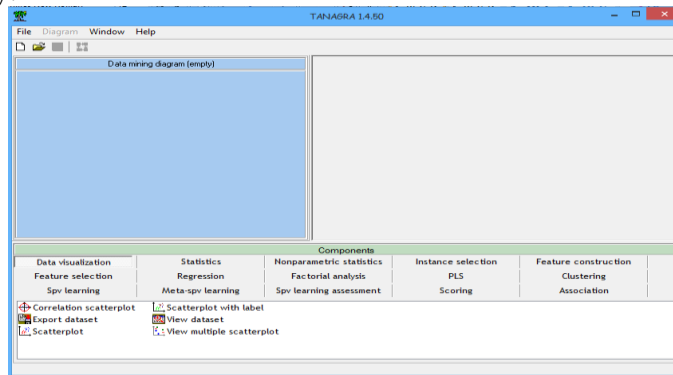


Figure 2 Tanagra Display 1.4

3. Display Download

Before the data is processed, the data must first be selected or downloaded as shown in the following image:
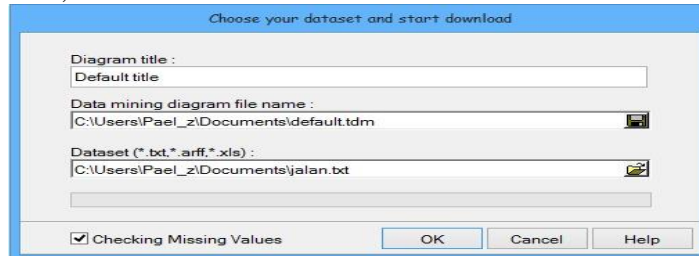


Figure 3 Display Download

4. Display View Database

To display the data that has been selected, the Database view tool is used, while the display is as follows:
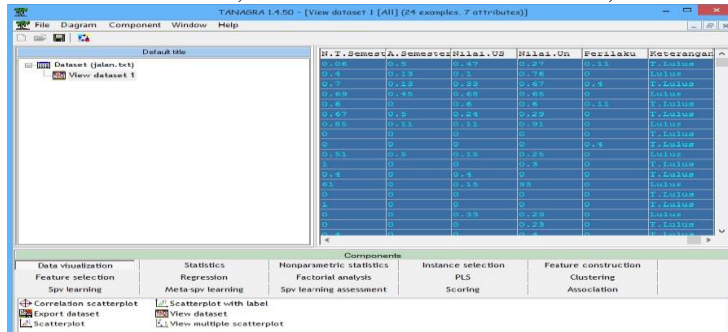


Figure 3 Display data

5. Display Define Attribute Statuses

Define Attribute Statuses are used to input attributes and input Targets. The appearance is as follows:
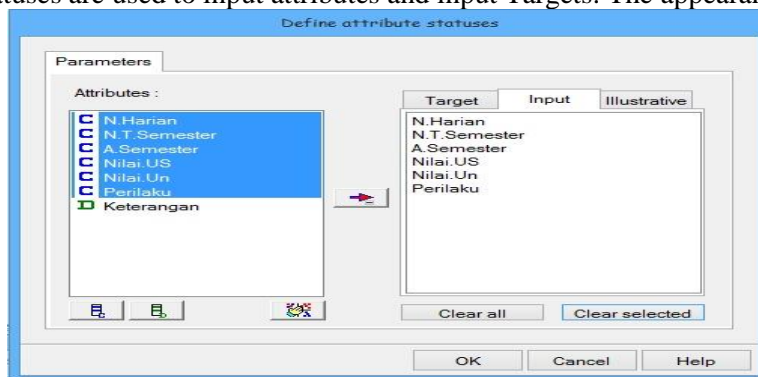


Figure 5 Display Define Attribute Statuses

6. Display C 4.5 Parameters

C 4.5 Parameters are used to input a minimum size of leaves and a Confidancel level, while the appearance is as follows:
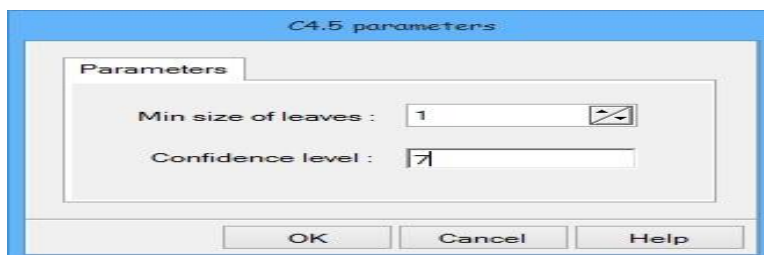


Figure 6 Display C 4.5 Parameters

7. Display Supervised Learning

Supervised Learning is the final stage to find out the Decision Tree. The appearance is as follows:
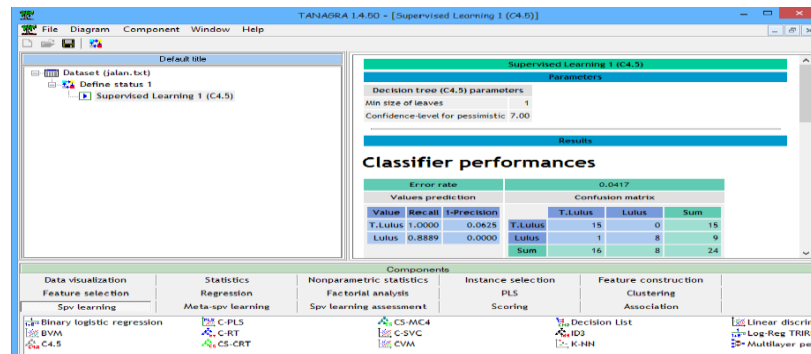
Figure 7 Display of Supervised Learning

8. Display the results of testing C 4.5 algorithm

The results of testing the C 4.5 algorithm can be seen in the following figure, which is also the final result of the process in the Tanagra software. As seen in the following picture:
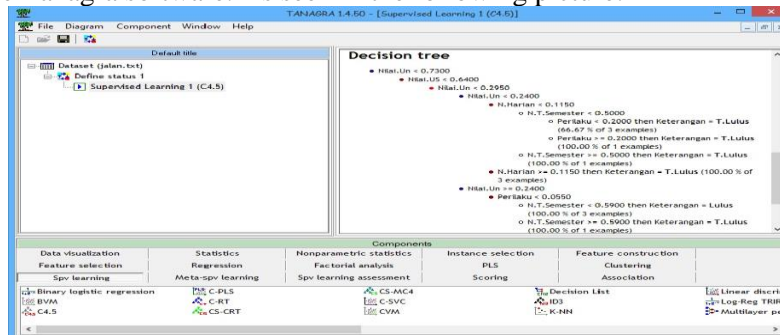


Figure 8 Display of C 4.5 Algorithm with Decision Tree

## 5. Conclusions

The results of data mining using the Decision Tree method is a sequence of activities that support each other in the student assessment process so that it is easier to understand by looking at the stages of the decision tree image.

Reference

[1]     X. Wu, X. Zhu, G. Q. Wu, and W. Ding, "Data mining with big data," *IEEE Trans. Knowl. Data Eng.*, 2014.

[2]     J. Han, M. Kamber, and J. Pei, *Data Mining: Concepts and Techniques*. 2012.

[3]     S. Agarwal, "Data mining: Data mining concepts and techniques," in *Proceedings - 2013 International Conference on Machine Intelligence Research and Advancement, ICMIRA 2013*, 2014.

[4]     D. Mining and T. Mining, "A Programmers Guide to Data Mining," *Text*, 2017.

[5]     J. Apostolakis, "An introduction to data mining," *Struct. Bond.*, 2010.

[6]     Y. Y. Song and Y. Lu, "Decision tree methods: applications for classification and prediction," *Shanghai Arch. Psychiatry*, 2015.

[7]     Y. Ben-Haim and E. Tom-Tov, "A streaming parallel decision tree algorithm," *J. Mach. Learn. Res.*, 2010.

[8]     S. B. Kotsiantis, "Decision trees: A recent overview," *Artificial Intelligence Review*. 2013.

[9]     A. Mukminin and D. Riana, "Komparasi Algoritma C4 . 5 , Naïve Bayes Dan Neural Network Untuk Klasifikasi Tanah," *J. Inform.*, 2017.

[10]   S. Hawani, A. P. Windarto, S. Solikhun, and D. Hartama, "PENERAPAN C 4.5 UNTUK

MENENTUKAN CALON SUAMI TERBAIK DALAM PERNIKAHAN PADA KANTOR KUA SIANTAR MARTOBA PEMATANGSIANTAR," *Jurasik (Jurnal Ris. Sist. Inf. dan Tek. Inform.*, 2017.

[11]   E. R. Ariyanto, Wijanarto, and Sudaryanto, "Klasifikasi Citra Porno dengan Algoritma C 4.5 Berbasis Model Warna YCbCr dan Shape Detector," *Techno.COM*, 2016.

[12]   B. Santoso, "Perancangan Aplikasi Data Mining Penjualan Laptop Pada Sinergi Komputer Lubuklinggau Menggunakan Algoritma C 4.5," *J. Ilm. Betrik*, 2017.