

DATASET MISSING VALUE HANDLING AND CLASSIFICATION USING DECISION TREE C5.0 AND K-NN IMPUTATION: STUDY CASE CAR EVALUATION DATASET

WAHYU WIDYANANDA^{1*}, MUHAMMAD FAUZAN EDY PURNOMO², MUHAMMAD ASWIN³, PANCA MUDJIRAHARDJO⁴, SHOLEH HADI PRAMONO⁵

¹Department of Electrical Engineering, Brawijaya University, East Java, Indonesia

^{2,3,4,5}Associate Professor, Department of Electrical Engineering, Brawijaya University, Indonesia

E-mail: ¹wahyuwidyandanda.111@gmail.com, ²mfauzanep@ub.ac.id, ³muhaswin@ub.ac.id, ⁴panca@ub.ac.id, ⁵sholehpramono@ub.ac.id

ABSTRACT

Data mining is a data analysis process using software to find certain patterns or rules from a large amount of data which is expected to find knowledge to support decisions. However, missing value presence in data mining often lead to loss of information. Information loss inside dataset such car evaluation can result in poor predictive models. The purpose of this study is to improve the performance of data classification with missing values precisely and accurately using Decision Tree C5.0 and k-NN Imputation. The test method is carried out using the Car Evaluation dataset from the UCI Machine Learning Repository. RStudio and RapidMiner tools were used for testing the algorithm. This study will result in data analysis of the tested parameters to measure the performance of the algorithm. Using test variations: 1. Performance at C5.0, C4.5, and k-NN at 0% missing rate. 2. Performance on C5.0, C4.5, and k-NN at 5-50% missing rate. 3. Performance on C5.0 + k-NNI, C4.5 + k-NNI, and k-NN + k-NNI at 5-50% missing rate. 4. Performance on C5.0 + CMI, C4.5 + CMI, and k-NN + CMI at 5-50% missing rate. The results show that C5.0 with k-NNI produce better classification accuracy than other tested imputation and classification algorithms. For example, for 35% missing in the dataset, this method obtains 93.40% in validation accuracy and 92% accuracy in the test. C5.0 with k-NNI also offers fast processing time compared with others methods.

Keywords: *Missing Value Handling, C5.0, k-NNI, R-Studio, RapidMiner*

1. INTRODUCTION

When consumers consider buying a car, several factors can influence their decision to buy a car. Safety, cost, and luxury are important factors that must be considered in buying a car [1]. Assessing the cost and quality of a new product in the marketing stage of development allows a more accurate prediction of consumer acceptance of the product or service [2]. Collecting data on car purchases regarding these factors is needed to evaluate cars based on consumer interests, which results in a car evaluation dataset. Data mining algorithms have the ability to analyze data in various research fields and classification is one of the main roles in data mining that can be applied to car evaluation datasets to create predictive models based on consumer interest. Several studies have been carried out on car evaluation dataset to make prediction models by applying classification

algorithms [3] [4] [5]. Datasets with missing values are a common problem in data mining, which can lead to loss of information and result in poor predictive models [6]. Decision tree is one of the classification algorithms that can handle missing values during the classification process, besides that there is also a data imputation technique. This is one of the techniques used to handle missing values in a data set.

Various studies have been conducted to solve the problem of missing values in classification using the imputation method. Testing Decision Tree C4.5 without adding more imputation methods resulted in better prediction accuracy than adding Listwise Deletion and Mean Imputation methods to data with missing values [7]. Decision Tree C5.0 is an improvement from the Decision C4.5 algorithm. The Class Mean Imputation method with the Decision Tree C5.0 algorithm performed better than Mean Imputation on a dataset with high categorical

variables [8]. k-Nearest Neighbor Imputation (k-NNI) is an imputation method based on the k-Nearest Neighbor (k-NN) classification algorithm. k-NNI can improve Decision Tree C4.5 performance on small software projects [9] and student records [10].

Although there have been many studies related to missing value datasets using classification and imputation methods, there are no studies on missing values presence for car evaluation. In fact, datasets with missing values can lead to loss of information and result in poor predictive models. Poor predictive models can lead to wrong decisions. The k-NNI approach combined with Decision Tree C5.0 can be used as classification method to improve the prediction accuracy of the Car Evaluation dataset when missing values is presence.

In contrast to previous studies, this study aims to improve the classification performance of the Car Evaluation dataset with missing values presence using Decision Tree C5.0 and k-Nearest Neighbor Imputation (k-NNI). Comparison of other classification algorithms such as Decision Tree C4.5 and k-NN, combined with Class Mean Imputation and k-NNI was also carried out. Parameters used to measure classification performance include performance accuracy and average processing time. This research was conducted to overcome the problem of data classification with missing values to make decisions quickly, precisely, and accurately.

The research specifications are designated to describe the proposed approach and its limitations, due to the fact that it is both a phase and a package of the main design. They are summarized as follows:

- The dataset used in this study is a real-world Car Evaluation dataset from the UCI Machine Learning Repository [11]. This dataset has 1728 samples with six categorical type attributes, including Purchase Price, Maintenance Price, Number of Doors, Person Capacity, Luggage Boot Size, and Estimated Security.
- Missing values in the dataset are artificially generated using the Missing Completely At Random (MCAR) mechanism with ratios of 5%, 10%, 15%, 20%, 25%, 30%, 35%, 40%, 45%, and 50%
- R Studio 1.3.1073 is used for data imputation process while RapidMiner 9.7.002 is used for data classification and analysis process
- This research uses 3 data mining classification algorithms (Decision Tree C5.0, Decision Tree C4.5, k-NN) combined with 2 data imputation algorithms (k-NNI and CMI) in each classification algorithm to offer better comparison on decision support final results.

2. LITERATURE REVIEW

2.1 DATA GROUPING

The grouping of data in data mining is divided into two categories, namely classification and clustering. Classification is a grouping of data that requires training data (supervised). While clustering is data grouping without the need for training data (unsupervised). [12] explained that clustering is the process of dividing unlabelled data into groups of data that have similarities. Each data group (cluster) consists of objects that have similarities with each other and each cluster has dissimilarities with other clusters. Clustering is commonly used in multivariate data analysis.

2.2 DATA MINING

Data mining is a field of several scientific fields that combines techniques from machine learning, pattern recognition, statistics, databases, and visualization for handling problems of retrieving information from large databases [13]. In general, data mining can be grouped into 2 main categories, namely [14]:

- a) Descriptive mining, which is a process to find important characteristics of data in a database. Data mining techniques included in descriptive mining are clustering, association, and sequential mining.
- b) Predictive, which is the process of finding patterns from the data by using several other variables in the future. One of the techniques contained in predictive mining is classification.

2.3 CLASSIFICATION

Classification is the process of finding a model or function that explains or distinguishes a concept or data class, with the aim of being able to estimate the class of an object whose label is unknown [15]. Classification is a learning function that maps (classifies) an element (item) of data into one of several predefined classes. The input data for classification is a collection of records. Each record is known as an instance. Instance is defined by a set of attributes which is defined by x , and a specific attribute that expressed as a class label which is defined by y (also known as a category or target attribute).

Some of the classification techniques used are decision tree classifier, k-nearest neighbors, neural-network, support vector machine, and naive bayes classifier. Each technique uses a learning algorithm to identify the model that provides the most suitable relationship between the attribute set and the class label of the input data. K-nearest neighbor is a

distance-based classification technique [16]. It searches the pattern space from the training samples which is close to unknown samples. This technique produces a better level of accuracy than the naive Bayes technique in classifying Parkinson's disease [17], and is better than the naive Bayes technique and SVM for classifying news [18]. Decision Tree uses several algorithms to generate decision model patterns, including ID3, C4.5, and C5.0. The ID3 algorithm uses information gain calculations in the formation of a decision tree and can only classify data of categorical type. The C4.5 algorithm is the development of the ID3 Algorithm. Unlike the previous algorithm, C4.5 uses a gain ratio calculation in the formation of a decision tree, can classify continuous and categorical data types, and can handle training data sets with missing values [19]. The C5.0 algorithm is a development of the C4.5 algorithm. The C5.0 algorithm again uses information gain calculations to form a decision tree, and can also handle data with missing values. The C5.0 algorithm has a much lower error rate for the prediction case [20].

2.4. MISSING DATA IMPUTATION

Missing data is a condition where some features are lost in the dataset. Missing data can be caused by system errors such as no response to sensors or input receiving devices. It can also be caused by human errors such as incomplete data entry in the database or respondents' misunderstanding in filling out questionnaires in large-scale surveys so that they pass through the form provided. Existing methods in data mining can only process data that has complete features so that special handling is needed for this problem.

There are 3 methods used for handling missing data, namely: Case Deletion, Parameter Estimation, and Imputation Techniques [21]. Case deletion is the easiest method, namely by deleting data that contains missing. The weakness of this method is that it is possible to delete important information when missing data is deleted. The imputation technique is a method of handling missing data that is more widely studied. Data imputation is estimating the value of missing data by getting a pattern from data that has complete features. Some popular imputation methods are: Mean, Median/Mode and clustering, k-NN Imputation, and Class Mean Imputation [22]. In K-NN Imputation, filling in missing values is done by taking into account the vector distance between attributes [23]. In Class Mean Imputation, the missing value will be replaced by the mean value of all available values in a related group or class [8]. There are 3 data missing

mechanisms, including Missing Completely at Random (MCAR), where the distribution of missing data on an attribute does not depend on the observed or missing data. This method will use a complete dataset and then generate missing data randomly based on certain proportions. The advantage of this method is that it makes it easier for researchers to estimate computationally from the proposed model. Another mechanism is Missing at Random (MAR), where the distribution of missing data on an attribute depends on the observed data but does not depend on the missing data. The last one is Not Missing at Random (NMAR), if the distribution of missing data on an attribute depends on the missing data [24].

3. RESEARCH METHODOLOGY

The stages of the research methodology are: (1) Processing the dataset and deleting data from the dataset with a predetermined ratio; (2) dataset imputation using the k-NNI method; (3) Studying the training data using the Decision Tree C5.0 method. (4) Testing the resulting model with the test set that has been prepared. The research methodology for completing the classification of the Car Evaluation dataset with missing values is shown in Figure 1. The class distribution of the research dataset can be seen in Table 1.

Table 1 -Class Distribution of Research Dataset

No	Attributes	Total Samples	Total Samples (%)
1	Unacc	1210	70.02
2	Acc	384	22.22
3	Good	69	3.99
4	Very Good	65	3.76

3.1 Training Dataset

The dataset used for testing the algorithm uses the Car Evaluation dataset from the UCI Machine Learning Repository. The Car Evaluation dataset has a multiclass type and categorical attribute characteristic values. The training set was obtained by taking 1728 data samples from the Car Evaluation dataset. The Car Evaluation Dataset contains information on several car evaluation parameters with six characteristic attributes, including Purchase Price, Maintenance Price, Number of Doors, Person Capacity, Boot Baggage Size and Safety Estimate. In addition, the dataset produces four types of classes, namely Unacceptable, Acceptable, Good and Very Good.

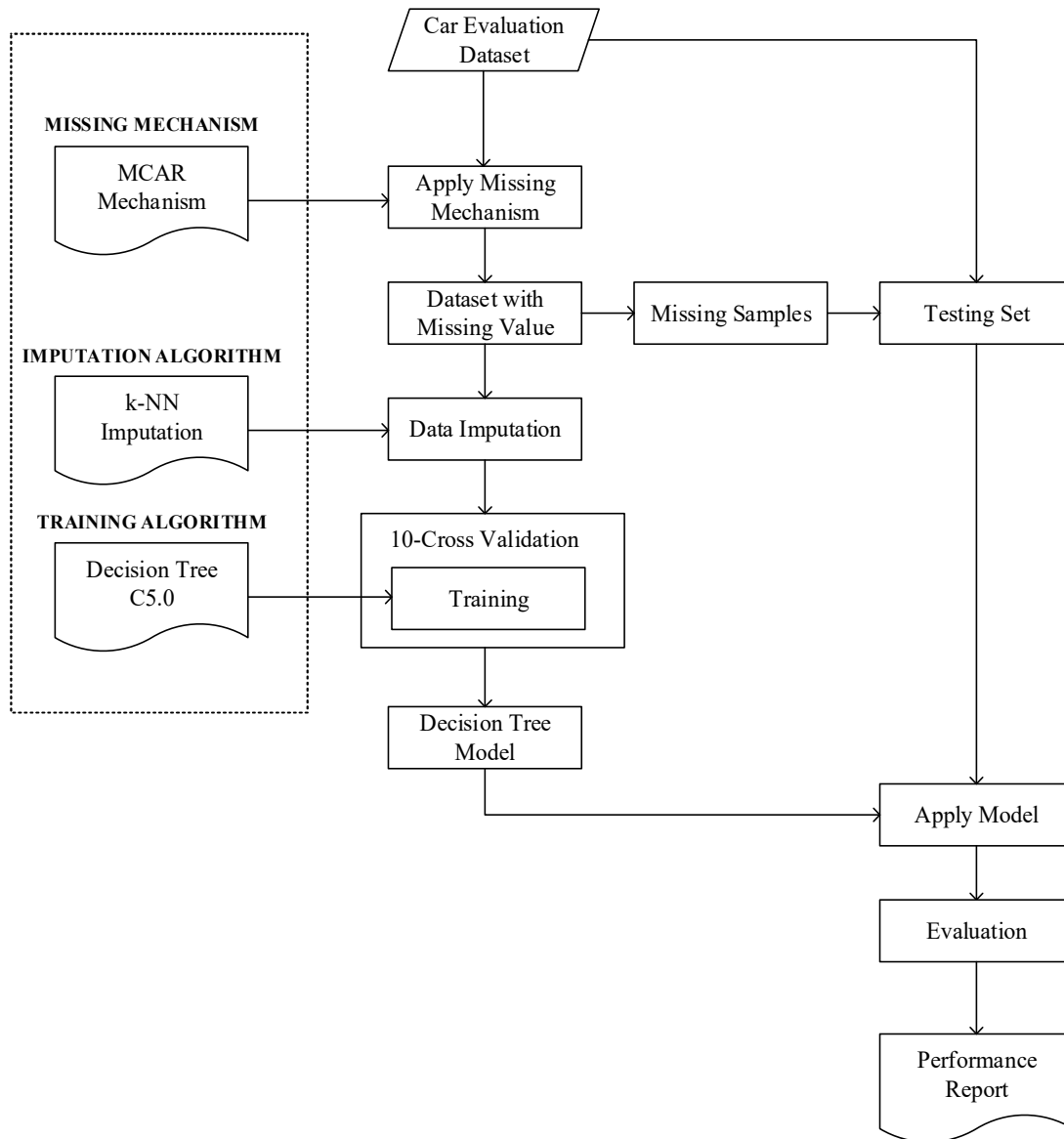


Figure 1 - Research Methodology Design

3.2 Missing Completely At Random Mechanism

Missing Completely at Random (MCAR) is a missing value mechanism where data loss occurs randomly. MCAR is most often encountered in actual cases and significantly affects the performance of the classification results. Therefore, it is used in this study. The following is the MCAR mechanism pseudocode:

Algorithm: Initialize MCAR
 Input: 'data' as data input, 'mp' as percentage of missing value
 Output: 'data' with missing value

BEGIN

```

Set x as instance numbers in data, set y as attribute
number in data
Set counter = 0, mv = x* y * mp
While counter < mv
    Data [random (0, x), random (0, y)] = null
END
    
```

3.3 K-NN Imputation

k-NN Imputation (k-NNI) is an imputation method based on the k-NN classification algorithm to impute the missing value based on several values that are close to the missing value. k-NNI can handle missing values by determining the nearest neighbor symbolized by k, then calculating the smallest distance from each neighbor that does not have a missing value. The distance between the missing

value and its neighbors can be calculated using the Euclidean distance formula.

The steps for entering missing values with the k-NNI method are as follows.

- o Determine the value of k
- o Find data with missing values in data set
- o Calculate the nan-euclidian distance from the initial observation data and other observation data, using the formula:

$$d(x_a, x_b) = \sqrt{w_{ab} * \sum_{j=1}^m (x_{aj} - x_{bj})^2 \dots \dots \dots (1)}$$

- o Choose k observational data with the smallest value
- o Select data on attributes related to missing values from selected observation data
- o Fill in the missing values with approximate values from the selected data

3.4 Decision Tree C5.0

Decision tree is a classification method by studying data into a tree-shaped pattern to produce decisions. C5.0 is the algorithm used in the decision tree method to classify and develop the previous algorithms, namely C4.5 and ID3. The C5.0 algorithm makes a decision tree model pattern based on the entropy value and information acquisition.

Entropy (S) is a value that expresses the uncertainty or impurity in a random data set from a data set expressed in bits. The entropy value is needed to calculate the information gain. Information gain is a measure of the effectiveness of an attribute in classifying data. Information gain is used to determine the order of attributes used to form a classification model pattern.

The following is the formula used to get the entropy and information gain values:

$$Entropy(S) = - \sum_{i=1}^n P_i (P_i) \dots \dots \dots (2)$$

$$Gain(S,T) = Entropy(S) - \sum_{j=1}^n (p_j \times Entropy(p_j)) \dots \dots \dots (3)$$

The decision concept of the C5.0 decision tree method is as follows:

- First, calculate the total entropy value of the dataset using equation 2.
- Calculate the entropy value and information gain for each attribute criterion using equations 2 and 3.
- Finally, determine the root node based on the largest information gain value using equation 3.

Define an internal node to generate a leaf node based on the entropy value and information gain. The process stops when all attributes have been used. Formation of rules based on the formed classification model pattern.

3.5 Validation

In this study, 10-fold cross-validation was used for the validation process. Cross-validation (CV) is an analytical method that can be used to evaluate the performance of a classifier, where the dataset is divided into two subsets, namely learning data and test data. The selection of the type of CV is based on the size of the dataset. The way k-fold cross-validation works is by dividing the dataset into two groups, namely training data and test data, then the testing process is carried out with k repetitions. The test results are then averaged to produce an accuracy value. A value of k that is too small can affect the accuracy value to be low. This is because the value of k is small so it is easily affected by noise. While the value of k that is too large will result in an ineffective process. This is because a large value of k takes a long time to test. The 10-fold cross-validation method has become the standard method for learning and testing data [25]. The following is an illustrative example of 10-fold cross validation.

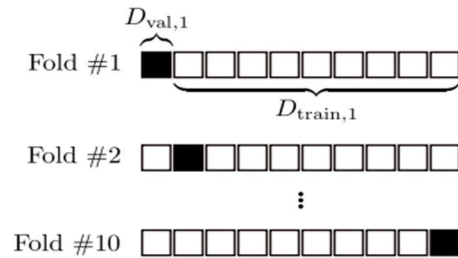


Figure 2 - Illustration of 10-cross validation method

3.6 Testing Dataset

The model generated from the validation process is undergoing a testing process using the test set that has been prepared. The test set is obtained by taking 100 data samples from the complete dataset before loss of data is generated. We selected 100 samples from the complete dataset based on the samples having missing values in the missing dataset. Files provided for testing are predicted in one of the predefined labels - Unacceptable, Acceptable, Good, and Very Good.

4. RESULTS AND DISCUSSION

This section presents the empirical results and discussion of applying Decision Tree C5.0 Algorithm and k-NN Imputation on Car Evaluation dataset with MCAR Mechanism. Therefore, we analyze the impact of different missing rate on the dataset to the predictive accuracy. Comparison of results by other classification algorithms such as Decision Tree C4.5 and k-NN, combined with CMI

and k-NNI was also presented in order to determine which algorithm produces the best performance on the case study. The scenario testing design can be seen in Figure 3.

There are four test scenarios:

1. Performance on C5.0, C4.5, and k-NN at 0% missing rate
2. Performance on C5.0, C4.5, and k-NN at 5-50% missing rate
3. Performance on C5.0 + k-NNI, C4.5 + k-NNI, and k-NN + k-NNI at 5-50% missing rate
4. Performance on C5.0 + CMI, C4.5 + CMI, and k-NN + CMI at 5-50% missing rate

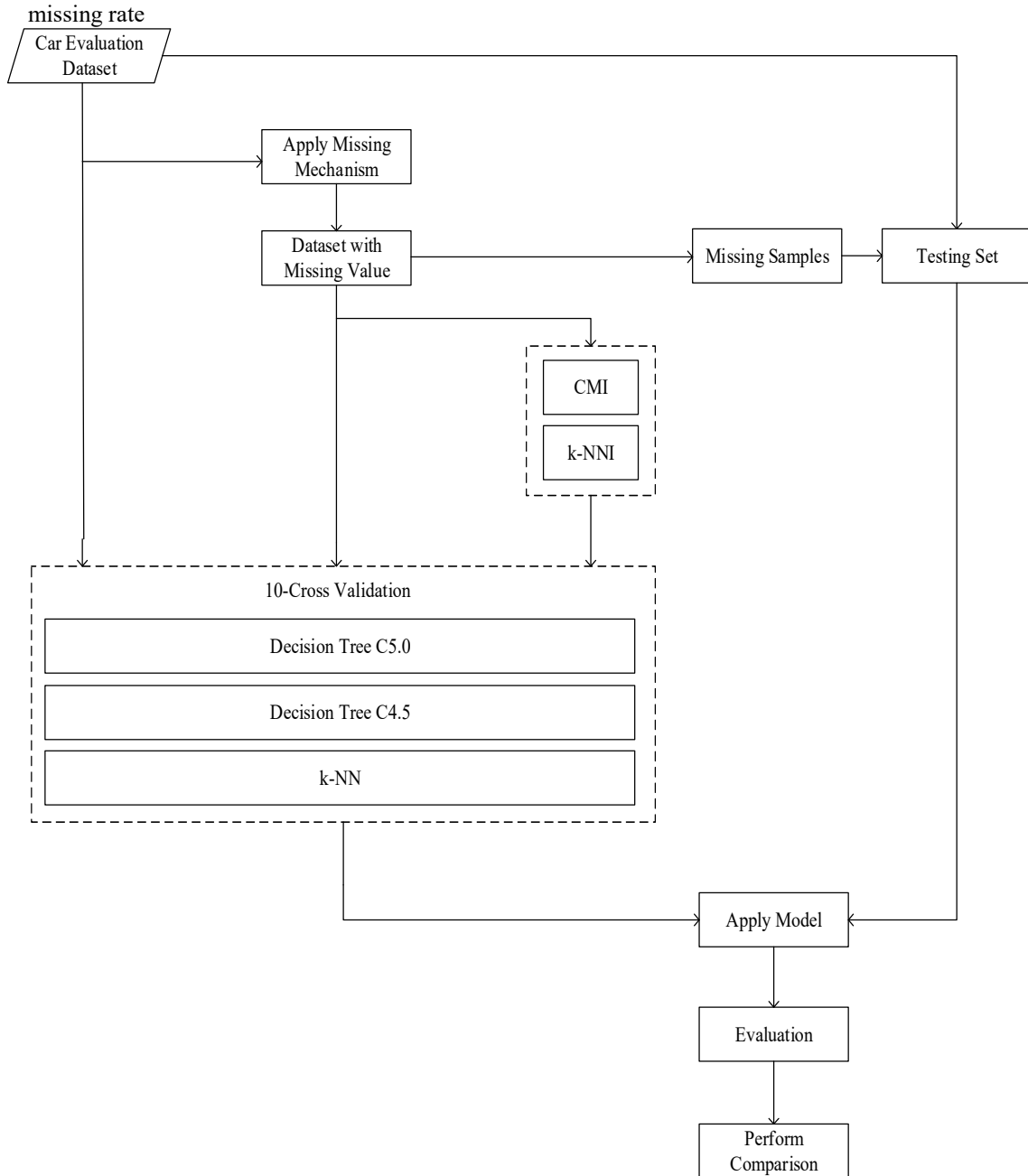


Figure 3 – Scenario Testing Design

In this study, the Rstudio tool is used to implement the missing mechanism and imputation process, and the RapidMiner tool is used to implement the classification algorithm. RStudio developed The

VIM package to explore and analyze the structure of missing values in a dataset, relate missing values to several imputation methods, and verify the imputation process using visualization tools [26]. In

this study, we use Rstudio to generate the MCAR mechanism in the dataset and implement the k-NNI algorithm using the k-NNI method provided by the VIM package. The results of the analysis of the performance of the C5.0 algorithms on the Car Evaluation dataset during validation process with or without the addition of the k-NNI method are shown in Figure 4.

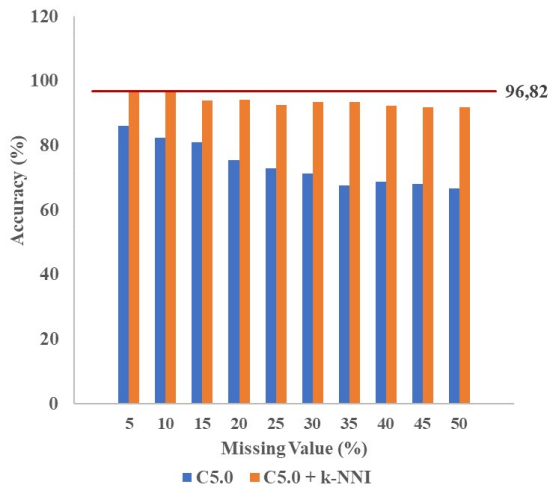


Figure 4 - Algorithm performance result on validation with and without k-NNI method.

Based on the results of the accuracy values obtained in the validation process, it was found that the number of missing values contained in the dataset can affect the predictive ability of the Decision Tree C5.0 method, which is characterized by a decrease in accuracy obtained when there are more missing values. Decision Tree C5.0 with k-NNI only experienced a slight decrease in the accuracy value of 0.18-4.89% from the accuracy value before the missing value was given, which was 96.82%. This method is much better than just using the Decision Tree C5.0 method to classify datasets with missing values, with a decrease in accuracy of 10.65-30.15% on datasets with missing values of 5-50%. The results of the accuracy value prove that the addition of the k-NNI method to Decision Tree C5.0 can improve the performance of the Decision Tree C5.0 method in classifying Car Evaluation datasets with missing values, with an additional accuracy of 10-25%.

Improved performance can also be seen in the testing process, by testing the decision tree model that has been formed into 100 testing sets, to find out the number of correct predictions from the model that has been formed. The results of the performance of the C5.0 on the Car Evaluation dataset during the testing process with or without the addition of the k-NNI method are shown in Table 2.

Table 2 C5.0 with k-NNI performance on testing

Missing Rate (%)	C5.0	C5.0 + k-NNI
5	95	99
10	92	99
15	94	97
20	92	96
25	86	90
30	86	94
35	83	92
40	77	88
45	82	88
50	82	88

From the test results above, it can be seen that during testing process the use of the Decision Tree C5.0 method with k-NNI produces a good accuracy value, with the lowest value of 88% in the dataset with a missing value of 50% and the highest value of 99% in the dataset with a missing value of 5% and 10%. This value is much better than just using the Decision Tree C5.0 method which produces an accuracy value of 82-95%.

Next, we present a performance comparison with other classification and imputation algorithms by implementing scenario testing design in Figure 3 to determine which algorithm yields the best performance in the case study. The results of the performance comparison during validation process are shown in Figure 5.

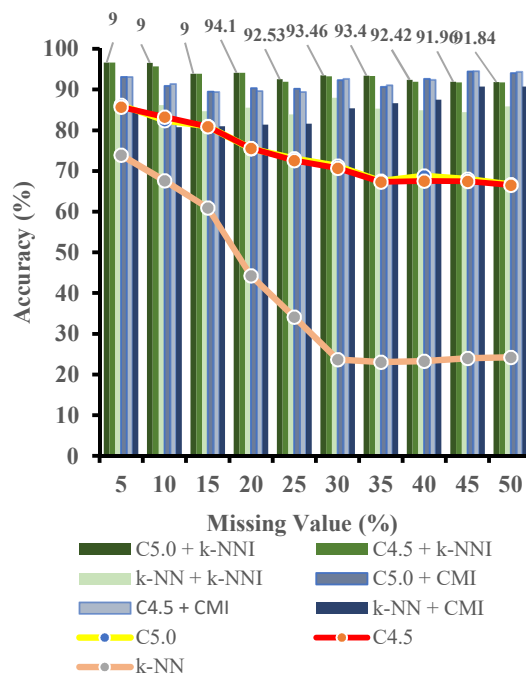


Figure 5 – Performance comparison result on validation with and without imputation method

Experiments show varying results for different imputation and classification methods. We found that C5.0 + k-NNI gave better prediction accuracy than other test methods in the validation and testing process, as shown in Figure 5 and Figure 6. C5.0 + k-NNI resulted in an average prediction accuracy at over 94% in the dataset under validation and a mean 95% prediction accuracy in the test for the missing rate between 5-35%. However, when the data set has 40% and more missing values, C5.0 + CMI and C4.5 + CMI perform better than C5.0 + k-NNI. The resulting accuracy for C5.0 + CMI and C4.5 CMI is around 92.42 - 94.5%, while C5.0 + k-NNI only produces a prediction accuracy of around 91.84 - 92.42% when 40 - 50 % missing added to the dataset. Probably due to high missing value presence, there are more flaws in the data set. There are fewer neighbors for the k-NNI method to generate the imputed value, and the CMI method can produce the imputed value better because it is based on the distribution of known values of the same class, thus having a sample.

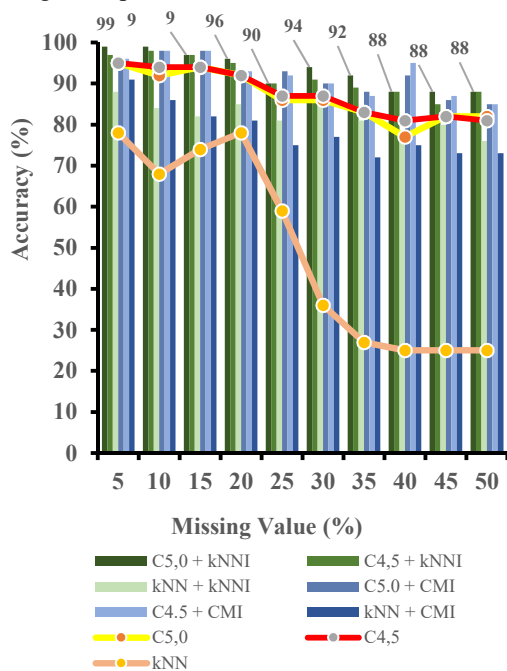


Figure 6 - Algorithm performance result on testing with and without imputation method.

Based on the results in Figure 4 and Figure 5, C4.5 without k-NNI and CMI models produces slightly better accuracy than C5.0 and k-NN in the test. Therefore, the C5.0 model with k-NNI produces higher accuracy than competing models in both validation and testing. Decision tree is the most suitable method for car evaluation databases than k-NN, Random Forest, Naïve Bayes, and Rule

induction whose decision trees have an accuracy of 91.1% [3]. When there are missing values presence inside car evaluation dataset, Decision Tree C5.0 with k-NNI provides a better prediction accuracy than the other tested classification algorithms although C4.5 still slightly outperforms C5.0 and k-NN before k-NNI was implemented.

The Group or Class Means imputation fare better than Mean Imputation in terms on accuracy when dealing with high categorical variables [8]. However, based on the overall test on Car Evaluation which has categorical type variables, C5.0 with k-NNI provides better prediction accuracy than other tested classification algorithms combined with CMI. Although in some cases, C5.0 with CMI prediction results outperforms C5.0 with k-NNI when there are high missing values presence inside dataset.

Changes in the accuracy of Decision Tree C5.0 with k-NNI form a straight line equation with R2: 0.89 - 0.91. This shows that the greater the missing value will increase the accuracy linearly. To determine the effect of missing values on the processing time of an algorithm, using an Intel Core i7 CPU running at 2.8 GHz and 16 GB RAM, the average data processing time for each missing value was obtained as shown in Figure 7.

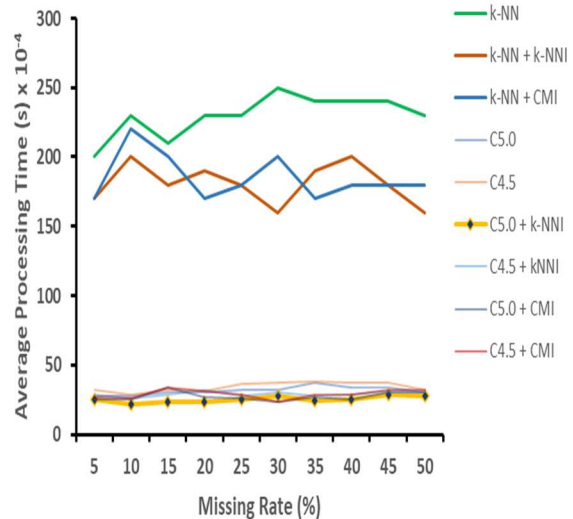


Figure 7-Average Computational Time with and without Imputation Method.

Figure 7 shows that in terms of processing time, the C5.0 + k-NNI offers fast processing time compared with other methods. k-NNI and CMI can reduce the processing time of C5.0 and C4.5. Those imputation methods also can reduce the processing time of k-NN classifiers. However, k-NN classifier processing time is still the slowest compared to other methods both in complete and incomplete data.

5. CONCLUSIONS

The results of this study show that the Decision Tree C5.0 with k-NNI is suitable for handling the classification of Car Evaluation dataset when there are missing values. This study also presents comparison with other classification algorithm such as Decision Tree C4.5 and k-NN, and other imputation method such as CMI in terms of predictive accuracy and processing time to find the best method to solve the missing value classification problem.

Based on the overall result of testing the algorithm on the Car Evaluation dataset at 5%, 10%, 15%, 20%, 25%, 30%, 35%, 40%, 45%, and 50% missing rate, C5.0 with k-NNI provides a better prediction accuracy than the other tested classification algorithms even though at some cases C4.5 still slightly outperforms C5.0 and k-NN before k-NNI is applied. In terms of processing time, the C5.0 with k-NNI has better performance than the C4.5 and k-NN algorithms, both with the imputation method and without the imputation method.

This study has some limitations in the mechanism of data loss is using MCAR (Missing Completely at Random) with of 0-50% ratio of missing values, and using the Rstudio and RapidMiner data science tools. Future research should consider the use of other data loss mechanisms and other data science tools.

ACKNOWLEDGEMENT

The authors would like to thank the Department of Electrical Engineering, Brawijaya University, who have supported this research project.

REFERENCES

- [1] [1] J. Awwalu, A. Ghazvini, and A. Abu Bakar, "Performance Comparison of Data Mining Algorithms: A Case Study on Car Evaluation Dataset," *International Journal of Computer Trends and Technology*, vol. 13, no. 2, 2014, [Online]. Available: <http://www.ijctjournal.org>
- [2] S. T. Luo, C. T. Su, and W. C. Lee, "Constructing intelligent model for acceptability evaluation of a product," *Expert Systems with Applications*, vol. 38, no. 11, pp. 13702–13710, Oct. 2011, doi: 10.1016/j.eswa.2011.04.162.
- [3] P. Jain and S. Kr Vishwakarma, "A Case Study on Car Evaluation and Prediction: Comparative Analysis using Data Mining Models," 2017.
- [4] A. Baykal, "PERFORMANCE ANALYSIS OF CLASSIFICATION ALGORITHMS OF SEVERAL DATA MINING SOFTWARES," *Middle East Journal of Science*, vol. 4, no. 2, pp. 104–112, Dec. 2018, doi: 10.23884/mejs.2018.4.2.06.
- [5] O. M. Al-Mubayyed, B. S. Abu-Nasser, and S. S. Abu-Naser, "Predicting Overall Car Performance Using Artificial Neural Network," 2019. [Online]. Available: <http://www.ijeais.org/ijaar>
- [6] R. R. S. B. and V. S. A. Sofia, "Data Mining Issues and Challenges: A Review," *IJARCCCE*, vol. 7, no. 11, pp. 118–121, Nov. 2018, doi: 10.17148/ijarccce.2018.71125.
- [7] S. Gavankar and S. Sawarkar, "Decision tree: Review of techniques for missing values at training, testing and compatibility," in *Proceedings - AIMS 2015, 3rd International Conference on Artificial Intelligence, Modelling and Simulation*, Oct. 2016, pp. 122–126. doi: 10.1109/AIMS.2015.29.
- [8] F. U. F. Khan, K. U. Z. Khan, and S. K. Singh, "Is Group Means Imputation Any Better Than Mean Imputation: A Study Using C5.0 Classifier," in *Journal of Physics: Conference Series*, Jul. 2018, vol. 1060, no. 1. doi: 10.1088/1742-6596/1060/1/012014.
- [9] Q. Song, M. Shepperd, X. Chen, and J. Liu, "Can k-NN imputation improve the performance of C4.5 with small software project data sets? A comparative evaluation," *Journal of Systems and Software*, vol. 81, no. 12, pp. 2361–2370, Dec. 2008, doi: 10.1016/j.jss.2008.05.008.
- [10] M. Rajan and V. Gimpy, "Missing Value Imputation in Multi Attribute Data Set." [Online]. Available: www.ijcsit.com
- [11] Dua, D. and Graff, C. (2019). UCI Machine Learning Repository [<http://archive.ics.uci.edu/ml>]. Irvine, CA: University of California, School of Information and Computer Science.
- [12] U. R. Yelipe, S. Porika, and M. Golla, "An efficient approach for imputation and classification of medical data values using class-based clustering of medical records," *Computers and Electrical Engineering*, vol. 66, pp. 487–504, Feb. 2018, doi: 10.1016/j.compeleceng.2017.11.030.
- [13] Y. Ariyanto, B. Hariyanto, and A. N. Asri, "Analyzing Student's Learning Interests in the Implementation of Blended Learning Using Data Mining," *International journal of online and biomedical engineering*, vol. 16, no. 11, pp. 153–160, 2020, doi: 10.3991/ijoe.v16i11.16453.

- [14] F. Siraj and M. Ali, "Mining Enrollment Data Using Descriptive and Predictive Approaches," in *Knowledge-Oriented Applications in Data Mining*, InTech, 2011. doi: 10.5772/14210.
- [15] A. P. Wibawa *et al.*, "Naïve Bayes Classifier for Journal Quartile Classification," *International Journal of Recent Contributions from Engineering, Science & IT (iJES)*, vol. 7, no. 2, p. 91, Jun. 2019, doi: 10.3991/ijes.v7i2.10659.
- [16] C. Shah and A. G. Jivani, "Comparison of data mining classification algorithms for breast cancer prediction," 2013. doi: 10.1109/ICCCNT.2013.6726477.
- [17] A. Elen and E. Avuclu, "A comparison of classification methods for diagnosis of Parkinson's," *International Journal of Intelligent Systems and Applications in Engineering*, vol. 8, no. 4, pp. 164–170, 2020, doi: 10.18201/ijisae.2020466309.
- [18] Y. Muliono and F. Tanzil, "A Comparison of Text Classification Methods k-NN, Naïve Bayes, and Support Vector Machine for News Classification," *Jl. Kh. Syahdan*, vol. 03, no. 02, 2018.
- [19] S. Kumar and H. Sharma, "A Survey on Decision Tree Algorithms of Classification in Data Mining," 2016. [Online]. Available: www.ijsr.net
- [20] P. Patidar and A. Tiwari, "Handling Missing Value in Decision Tree Algorithm," 2013.
- [21] H. Kang, "The prevention and handling of the missing data," *Korean Journal of Anesthesiology*, vol. 64, no. 5, pp. 402–406, May 2013. doi: 10.4097/kjae.2013.64.5.402.
- [22] Z. Zhang, "Missing data imputation: Focusing on single imputation," *Annals of Translational Medicine*, vol. 4, no. 1, Jan. 2016, doi: 10.3978/j.issn.2305-5839.2015.12.38.
- [23] L. Muflikhah, N. Hidayat, and D. J. Hariyanto, "Prediction of hypertension drug therapy response using K-NN imputation and SVM algorithm," *Indonesian Journal of Electrical Engineering and Computer Science*, vol. 15, no. 1, pp. 460–467, Jul. 2019, doi: 10.11591/ijeecs.v15.i1.pp460-467.
- [24] Y. Dong and C.-Y. J. Peng, "Principled missing data methods for researchers," 2013. [Online]. Available: <http://www.springerplus.com/content/2/1/222>
- [25] D. Berrar, "Cross-validation," in *Encyclopedia of Bioinformatics and Computational Biology: ABC of Bioinformatics*, vol. 1–3, Elsevier, 2018, pp. 542–545. doi: 10.1016/B978-0-12-809633-8.20349-X.
- [26] A. Kowarik and M. Templ, "Imputation with the R package VIM," *Journal of Statistical Software*, vol. 74, 2016, doi: 10.18637/jss.v074.i07.