

Predictive Analysis of Toddler Nutrition Using C5.0 Decision Tree Method

Styawati

Information Systems
Universitas Teknokrat Indonesia
Bandar Lampung, Indonesia
styawati@teknokrat.ac.id

Andi Nurkholis

Informatics
UPN "Veteran" Yogyakarta
Yogyakarta, Indonesia
andinurkholis@upnyk.ac.id

Syahirul Alim

Computer Engineering
Universitas Teknokrat Indonesia
Bandar Lampung, Indonesia
syahirul_alim@teknokrat.ac.id

S.Samsugi

Computer Engineering
Universitas Teknokrat Indonesia
Bandar Lampung, Indonesia
s.samsugi@teknokrat.ac.id

Chafidz Asyad

Computer Engineering
Universitas Teknokrat Indonesia
Bandar Lampung, Indonesia
chafidz_asyad@teknokrat.ac.id

Abstract—The nutritional status of toddlers is one of the benchmarks that can reflect their health. Malnutrition status in toddlers does not occur suddenly but begins with limited weight gain that is not enough. Changes in the toddler's weight from time to time are an early indication of changes in the nutritional status of toddlers. This study aims to determine the nutritional status of a toddler by applying the C5.0 Decision Tree Method. The dataset is divided into two categories: explanatory factors and target class. Explanatory factors are the criteria for determining nutritional status for toddlers, which include gender, age, weight and height, head circumference, and arm circumference. While the target class represents the nutritional status of toddlers, consisting of two classes, namely good and bad. Two prediction models are generated based on the 70:30 and 80:20 data partitions. The 70:30 partition model variation produces an accuracy of 84%, recall is 0.42, precision is 0.56, and f1-score is 0.48. While the 80:20 partition model partition obtained an accuracy of 89%, recall is 0.71, precision is 0.62, and f1-score is 0.67. The best prediction model is expected to be a solution for the community and related stakeholders as a follow-up to prevent malnutrition in toddlers. Comparisons to different algorithms may be made to create a more effective model.

Keywords—C5.0, Decision Tree, nutrition, prediction model, toddler

I. INTRODUCTION

Toddlers are children aged less than five years who are included in the age group at high risk of disease. Deficiency or excess of nutritional intake in toddlers can affect nutritional status and health status [1]. Based on health report data from World Health Organization (WHO), there are several nutritional status problems commonly suffered by toddlers that affect their health, such as lack of energy, protein, obesity, lack of vitamin A, disorders due to iodine deficiency, and anemia or lack of iron (Fe) [2], [3].

Nutritional status can be determined through laboratory examination or anthropometry. Anthropometry is the easiest and cheapest way to determine nutritional status. To get the right results, a benchmark is given as a guide, namely the Z-Score [4]. Z-Score is an anthropometric index used internationally to determine nutritional status and growth, expressed as a unit of population standard deviation. Z-Score is used to calculate nutritional status anthropometrically on weight for age, height for age, and

weight for height [5]. If the toddler continues to experience malnutrition, the toddler can experience child development deviations [6]. This can result in the inhibition of the growth process of toddlers so that they have different phases from other normal children. Some symptoms often experienced are delays in speaking at their age, difficulty adapting to the environment, and having a different face from other normal children. Based on these symptoms, the deviations experienced can be categorized into several types, namely, Autism, Down syndrome, and ADHD [7]. Therefore, it is essential to know the early symptoms or criteria for a toddler experiencing malnutrition or not [8] so that the best treatment can be done so that toddlers who experience malnutrition can be helped.

Previous works have studied many related to predicting the nutritional status of a toddler. The first study applied the Naive Bayes method to produce a predictive model of the nutritional condition of a toddler at Posyandu Melati IV Magetan District, East Java which obtained an accuracy of 60% [9]. The second study applied the Fuzzy K-Nearest Neighbor classification to predict the nutritional status of a toddler at the Kertosono Public Health Center, Nganjuk Regency, which resulted in an accuracy of 84.37% [10]. The last study applied the C4.5 algorithm for the nutritional status of a toddler at the Mranti Purworejo Health Center, which resulted in an accuracy of 88.24% [11]. The three studies succeeded in providing recommendations for predicting the nutritional status of a toddler that can be used as recommendations for follow-up handling. Based on the results of less-than-optimal accuracy in previous studies, the C5.0 algorithm (which is the development of the C4.5 algorithm) [12], [13] can be used to improve the prediction model's performance.

This study aims to produce a predictive model of nutritional status in a toddler by applying the C5.0 algorithm. The dataset used comes from Posyandu Kebun Dalam, divided into two categories: explanatory factors and target class. The explanatory factor is the criterion for determining the nutritional status class for toddlers, in the case of this study, which includes the gender of the toddler, age, weight and height, head circumference, and arm circumference. At the same time, the target class represents the nutritional status of toddlers, consisting of two classes, namely good and bad. It is hoped that the resulting research

can be used as a predictive model for the nutritional status of a toddler for relevant stakeholders (community, health workers, and the government) to anticipate an increase in malnutrition in toddlers.

II. RESEARCH METHOD

This research case study is located at the Posyandu (Integrated Services Post) in Kebun Dalam Village, Abung Tinggi District, North Lampung, Indonesia. The choice of the posyandu implies the high level of undernutrition in the area and its surroundings. Thus, the prediction model produced by this study will help find out information on the main factors causing and symptoms of malnutrition in a toddler. The stages of the research explain the sub-menu of the parts of the stages carried out, which can be seen in Figure 1.

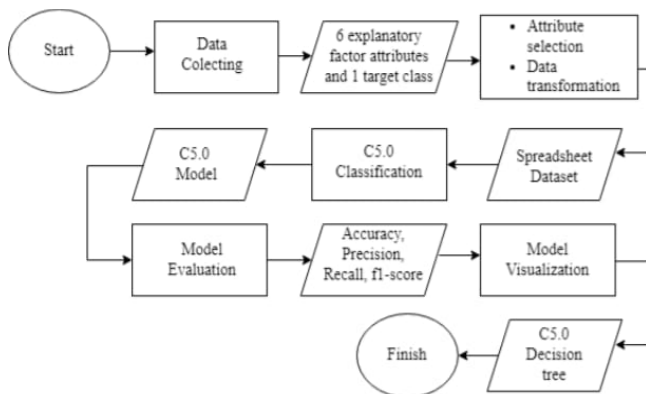


Fig. 1. Research stages

Figure 1 shows five main stages: data collecting, data preprocessing, modelling, model evaluation, and model visualization. The explanation of each stage is as follows:

A. Data Collecting

TABLE I. DATA RESEARCH

Attribute	Description
Age	Describing the age of children with a range of 0-5 years which is a toddler.
Gender	Describe the gender of the toddler, i.e. male or female
Weight (W)	Describing the amount of body weight that is owned by (kg)
Height (H)	Describing the total height of toddlers in centimeters (cm)
Head Circumference	Describe the size of the distance from the forehead to the back of the head as one of the criteria for measuring nutritional adequacy in centimeters (cm).
Arm Circumference	Describe the size of a toddler's arm as one of the benchmarks for a baby's nutritional status in centimeters (cm)
Mother's Name	Describing the mother's name from the existing toddler list
Mother's Age	Describing the sum of the ages of mothers of a toddler
Address	Describing the addresses of toddlers in Kibang and Kebun Dalam villages
Toddler Name	Describing the names of toddlers at the Kebun Dalam Posyandu
Status	Describe the nutritional status class consisting of Good or Poor Nutrition

At this stage, data collection was carried out through observations carried out by taking into account the local environment in Posyandu Kebun Dalam in the period 2017-2021. In addition, an interview process was also carried out to find out the information needed for the Posyandu Kebun Dalam. The posyandu then provides medical record data on the nutritional status of toddlers within the last 1 year in a spreadsheet format with 231 rows of data, 10 explanatory attributes and 1 label class.

B. Data Preprocessing

Before classification using the C5.0 algorithm, it is necessary to prepare data to obtain optimal modelling results [14][15]. The two stages carried out are attribute selection and data transformation.

- Using attributes in data mining can affect the resulting pattern. Therefore, it is necessary to select the attributes to be used. From the ten explanatory factors in Table 1, six attributes were selected: Toddler Age, Gender, Weight, Height, Head Circumference, and Arm Circumference. The six selected attributes affect growth in toddlers, especially nutritional status. The selection of these six attributes was based on the expertise of the midwife of Posyandu Kebun Dalam, who stated that four other attributes were not included in determining the nutritional status of a toddler. Attribute selection is made manually by deleting columns in spreadsheet file. The results of the attribute selection are shown in Table II.

TABLE II. ATTRIBUTE SELECTION RESULT

Age	Gender	W	H	Head Circum.	Arm Circum.	Status
3	L	11.5	87.4	46	16.5	Poor
3	L	11.4	88.4	48	14	Poor
3	P	12.5	101.7	49	16.5	Good
3	P	11.2	99	47	11.6	Good
3	L	14	87.4	52	16	Poor
3	L	13.9	89.3	46	16.5	Poor
3	L	13	101.2	49	14	Poor
2	P	12.2	80	45	16	Good
3	P	13.5	87.4	47	16.7	Good
2	P	11.5	80	45	16	Good

- The data that has been selected for attributes is then transformed by changing the data shape into a categorical to facilitate modelling. This is because the C5.0 classification process, which will be carried out using the DecisionTreeClassifier function from the Sklearn library version 1.0 contained in Python programming version 3.7.13, requires numeric data. For example, the value of the attribute Gender, male and female, is converted to numeric form 1 or 0. This process is carried out using the Label Encoder library in python programming. The results of the data transformation selection is shown in Table III.

TABLE III. DATA TRANSFORMATION RESULT

Age	Gender	W	H	Head Circum.	Arm Circum.	Status
3	0	11.5	87.4	46	16.5	1
3	0	11.4	88.4	48	14	1
3	1	12.5	101.7	49	16.5	0
3	1	11.2	99	47	11.6	0
3	0	14	87.4	52	16	1
3	0	13.9	89.3	46	16.5	1
3	0	13	101.2	49	14	1
2	1	12.2	80	45	16	0
3	1	13.5	87.4	47	16.7	0
2	1	11.5	80	45	16	0

C. C5.0 Classification

The C5.0 algorithm is an extension of the C4.5 algorithm, which has advantages, especially in large dataset. The C5.0 algorithm is better than the C4.5 algorithm in efficiency and memory [16]. In general, the flow of the decision tree process in the C5.0 and C4.5 algorithms has similarities, where both algorithms perform entropy and gain calculations. The C4.5 algorithm will only stop at the gain calculation, while the C5.0 algorithm will continue by calculating the gain ratio based on the gain and entropy values [17]. The gain ratio measure is used to select the test attribute for each node in the tree. The attribute with the highest gain ratio value will be chosen as the parent for the next node. The C5.0 algorithm breaks down the training data based on the largest gain information value attribute. In detail, the process flow of the C5.0 algorithm is shown in Figure 2 [18].

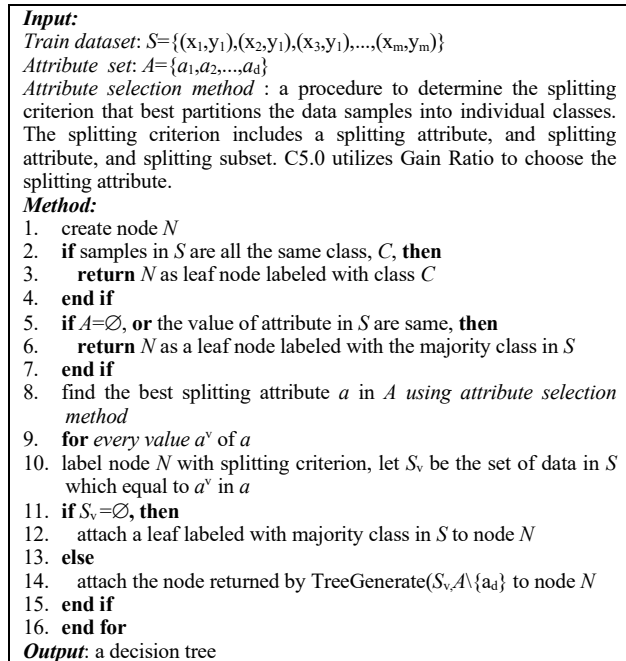


Fig. 2. C5.0 decision tree algorithm

D. Model Evaluation

To obtain the best results, an evaluation of the model is carried out by calculating the accuracy that will indicate the correctness of classifying the data to the actual class [19]. The higher the accuracy value, the lower the prediction

error of the test data, so that it reflects the model has a good performance [20]. The evaluation carried out using the confusion matrix method includes aspects of accuracy, precision, recall, and f1-score so that it will be able to determine the best model. The formulation used for evaluation is shown in Equation 1 to Equation 4.

$$Accuracy = \frac{TP+TN}{TP+FN+TN+FP} \quad (1)$$

$$Precision = \frac{TP}{TP+FP} \quad (2)$$

$$Recall = \frac{TP}{TP+FN} \quad (3)$$

$$F_1 - score = \frac{2 \times Precision \times Recall}{Precision + Recall} \quad (4)$$

Where TP (True Positive) is positive samples are predicted to be positive correctly. FN (False Positive): positive samples are predicted to negative wrongly. FP (False Negative): negative samples are predicted to be positive wrongly. TN (True Negative): negative samples are predicted to negative correctly.

E. Model Visualization

Visualization represents the modeling results that have been carried out and presented in graphical type. Visualization is used to facilitate the delivery of information and simplify decision-making. This study will display the best classification results from applying the C5.0 algorithm as a decision tree. This stage will use the Graphviz function from the Sklearn library in python programming.

III. RESULT AND DISCUSSION

The data preprocessing stage produces seven attributes consisting of six explanatory factors and one target class. The total row of data reaches 231, all of which are numeric. The dataset was then classified using the C5.0 algorithm to produce a predictive model for the nutritional status of toddlers.

A. C5.0 Decision Tree Model

The C5.0 algorithm modeling on the dataset uses the DecisionTreeClassifier function from the Sklearn library. In this study, two model variations were compared to obtain the best rule results, especially regarding accuracy. The model variation is based on the data partition ratio, where the first variation uses a ratio of 70:30, then the second variation uses a ratio of 80:20. In the 70:30 ratio model variation, the dataset is divided by a ratio of 70% of the dataset to be training data and 30% to be test data. While in the 80:20 ratio model variation, the dataset is divided by a ratio of 80% of the dataset being training data and 20% being test data.

B. Model Evaluation

Evaluation is done by testing the model on the test data based on the variation of the model, namely 70:3 and 80:20. The evaluation results of the two models cover aspects of accuracy, precision, recall, and f1-score, which are

calculated using equations 1 to 4. Details of the evaluation results of the two models can be seen in Table IV.

TABLE IV. MODEL EVALUATION RESULT

Aspect	Model variation result	
	70:30 partition data	80:20 partition data
Accuracy	0.84	0.89
Sensitivity/ Recall	0.42	0.71
Precision	0.56	0.62
F1-Score	0.48	0.67

Table 4 shows that the 80:20 model produces better accuracy than the 70:30 model, with a difference of 0.05. This can be caused by the ratio of training data used in the 80:20 model larger, the resulting model can represent more data. Thus, the number of rules produced by the 80:20 model is more than the 70:20 model, which is 21 versus 19. In addition, the test data used by the 80:20 model is also less, so the model can predict it better. In other aspects, namely recall, precision, and f1-score, the 80:20 model also produces better scores than the 70:30 model. This confirms that the 80:20 data partition can create a better model than the 70:30 partition. Based on the previous discussion, it can be stated that the best model in this study is the 80:20 variation. Even so, the 80:20 partition itself sometimes has drawbacks in a case, namely if it is indicated that it is overfitting. The 80:20 model produced in this study does not indicate overfitting, as evidenced by the fact that there are still errors when testing, which is 11% of the test data. The details of the 80:20 confusion matrix model can be seen in Table V.

TABLE V. CONFUSION MATRIX OF MODEL 80:20

Total = 47 (20% of 231)	Predicted: Good	Predicted: Poor
Actual: Good	5	2
Actual: Poor	3	37

C. Model Visualization

The best model results are 21 rules, where all attributes are involved in the rules. This shows that the attribute selection process was successful and implies that these attributes affect the nutritional status of toddler. Here are some of the best model rules:

1. IF weight \leq 15.1 kg AND age \leq 2.5 years AND weight \leq 8.45 kg THEN nutrition status is Poor
2. IF weight $>$ 15.1 kg THEN nutrition status is Good
3. IF weight \leq 15.1 kg AND age \leq 2.5 years AND weight $>$ 8.45 kg AND weight $>$ 8.95 kg THEN nutrition status is Good
4. IF weight \leq 15.1 kg AND age \leq 2.5 years AND weight $>$ 8.45 kg AND weight $>$ 8.95 kg AND age \leq 1.5 THEN nutrition status is Good
5. IF weight \leq 15.1 kg AND age \leq 2.5 years AND weight $>$ 8.45 kg AND weight \leq 8.95 kg AND age $>$ 1.5 THEN nutrition status is Poor
6. IF weight \leq 15.1 kg, age $>$ 2.5 years AND weight $>$ 13.95 kg AND arm circumference $>$ 16.875 THEN nutrition status is Poor

7. IF weight \leq 15.1 kg, age $>$ 2.5 years AND weight \leq 13.95 kg AND arm circumference $>$ 16.6 THEN nutrition status is Poor
8. IF weight \leq 15.1 kg, age $>$ 2.5 years AND weight \leq 13.95 kg AND arm circumference \leq 16.6 AND arm circumference \leq 14.125 AND head circumference \leq 48.25 THEN nutrition status is Poor
9. IF weight \leq 15.1 kg, age $>$ 2.5 years AND weight \leq 13.95 kg AND arm circumference \leq 16.6 AND arm circumference \leq 14.125 AND head circumference $>$ 48.25 AND height \leq 89.0 THEN nutrition status is Poor
10. IF weight \leq 15.1 kg, age $>$ 2.5 years AND weight \leq 13.95 kg AND arm circumference \leq 16.6 AND arm circumference $>$ 14.125 AND height \leq 88.65 THEN nutrition status is Poor

The result of the best model rules is then visualized into a decision tree using the Graphviz function from the Sklearn library. Thus, a decision tree result is saved into a pdf document. The following is a partial decision tree of the best model shown in Figure 3.

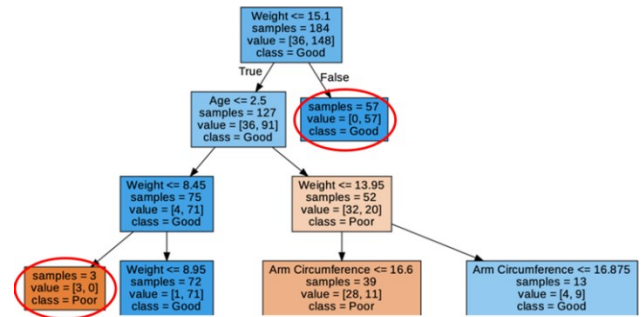


Fig. 3. C5.0 decision tree of 80:20 model

Based on Figure 3, the best model defines weight as the root node, which means that the main factor determining the nutritional status of toddler. The weight attribute chosen as the root node implies a high gain ratio value, which means it has the lowest diversity among other attributes. All the attributes used in the modelling are involved in preparing the decision tree. It can be assumed that all the attributes influence the nutritional status of toddlers. Relevant stakeholders (community, health staff, government) can use the resulting rules to become symptom information in knowing the nutritional status of toddlers. On the other hand, regulations can also be used as factors used in efforts to improve the nutritional status of a toddler. For example, if a toddler wants a good nutritional status, the weight must exceed 15.1 kg. Meanwhile, if a toddler weighs less than 8.45 kg and is under 2.5 years of age and weighs less, it has the potential for poor nutritional status. As an implication, the resulting prediction model can provide knowledge for the community to prevent nutritional deficiencies in a toddler.

IV. CONCLUSION

This study produces two decision trees for predicting the nutritional status of a toddler using the C5.0 algorithm. The dataset consists of 231 rows and is divided into two categories: six explanatory factors (Age, Gender, Weight, Height, Head Circumference, Arm Circumference and one target class (Good or Poor). The model using the 80:20 data

partition obtained an accuracy of 0.89, while the model using the 70:30 data partition only obtained an accuracy of 0.84. In addition, the 80:20 partition model also dominates all aspects of the evaluation consisting of recall 0.71, precision 0.62, and f1-score 0.67, while the 70:30 model produces a recall of 0.42, an accuracy of 0.56, and an f1-score of 0.48. The best model involves all attributes in the resulting rule, where weight is used as a root node. With the accuracy obtained, the results of the rules can help decision-making to predict the nutritional status of toddlers very well. As a development, comparisons can be made to other algorithms to produce a more optimal model.

ACKNOWLEDGMENT

This research was funded by Hibah Publikasi Terindeks Scopis (HIPUTS) Universitas Teknokrat Indonesia 2024.

REFERENCES

- [1] N. Lumongga, E. Sudaryati, and D. Theresia, "The relationship of visits to posyandu with the nutrition status of toddlers in amplas health center," *Budapest International Research and Critics Institute(BIRCI-Journal): Humanities and Social Sciences*, vol. 3, no. 3, pp. 2165–2173, 2020.
- [2] Z. Y. Amare, M. E. Ahmed, and A. B. Mehari, "Determinants of nutritional status among children under age 5 in Ethiopia: further analysis of the 2016 Ethiopia demographic and health survey," *Global Health*, vol. 15, no. 1, pp. 1–11, 2019.
- [3] M. Mkhize and M. Sibanda, "A review of selected studies on the factors associated with the nutrition status of children under the age of five years in South Africa," *Int J Environ Res Public Health*, vol. 17, no. 21, p. 7973, 2020.
- [4] K. Stephens *et al.*, "Evaluating mid-upper arm circumference z-score as a determinant of nutrition status," *Nutrition in Clinical Practice*, vol. 33, no. 1, pp. 124–132, 2018.
- [5] T. E. Putri, R. T. Subagio, and P. Sobiki, "Classification system of toddler nutrition status using naïve Bayes classifier based on Z-score value and anthropometry index," in *Journal of Physics: Conference Series*, 2020, vol. 1641, no. 1, p. 012005.
- [6] R. D. Hayuningtyas, S. F. N. Laila, and N. Nurwijayanti, "Analysis of factors affecting the development of children of toddler ages assessed from history of infection diseases, nutritional status and psychosocial stimulation in Ponorogo regency," *Journal for Quality in Public Health*, vol. 3, no. 2, pp. 341–347, 2020.
- [7] A. A. Sariza, Z. Maristka, L. Hayati, R. Inggarsih, and S. Purnamasari, "Dermatoglyphics findings in intellectual disability children with down syndrome, autism spectrum disorder and attention-deficit hyperactivity disorder: A descriptive cross-sectional study," *Advances in Human Biology*, vol. 11, no. 4, p. 34, 2021.
- [8] M. Wong Vega, S. Beer, M. Juarez, and P. R. Srivaths, "Malnutrition risk in hospitalized children: A descriptive study of malnutrition-related characteristics and development of a pilot pediatric risk-assessment tool," *Nutrition in Clinical Practice*, vol. 34, no. 3, pp. 406–413, 2019.
- [9] N. Rahmawati, Y. Novianto, and J. Jasmir, "Classification of nutritional conditions of toddlers using the naïve Bayes method (case study of Posyandu Melati IV)," *Scientific Journal of Informatics Engineering*, vol. 2, no. 3, pp. 257–268, 2020.
- [10] S. D. Nugraha, R. R. M. Putri, and R. C. Wihandika, "Application of fuzzy k-nearest neighbor (FK-NN) in determining the nutritional status of toddlers," *Journal of Information Technology and Computer Science Development*, vol. 2548, p. 964X, 2017.
- [11] C. Agustina, "Comparison of C. 45 algorithm and backpropagation for classification of toddler nutritional status based on anthropometric index Bb/U and BW/PB," *Speed-Center for Engineering Research and Education*, vol. 9, no. 3, 2017.
- [12] A. Z. Abdullah, B. Winarno, and D. R. S. Saputro, "The decision tree classification with C4. 5 and C5. 0 algorithm based on R to detect case fatality rate of dengue hemorrhagic fever in Indonesia," in *Journal of Physics: Conference Series*, 2021, vol. 1776, no. 1, p. 012040.
- [13] M. A. Febriantono, S. H. Pramono, R. Rahmadwati, and G. Naghdy, "Classification of multiclass imbalanced data using cost-sensitive decision tree C5. 0," *IAES International Journal of Artificial Intelligence*, vol. 9, no. 1, p. 65, 2020.
- [14] S. Styawati, A. Nurkholis, A. A. Aldino, S. Samsugi, E. Suryati, and R. P. Cahyono, "Sentiment analysis on online transportation reviews using Word2Vec text embedding model feature extraction and support vector machine (SVM) algorithm," in *2021 International Seminar on Machine Learning, Optimization, and Data Science (ISMODE)*, 2022.
- [15] S. Styawati, A. Nurkholis, F. A., Ans, S. Alim, L. Andraini, and R. A. Prasetyo, "Web scraping for summarization of freelance job website using vector space model," in *Proceeding - IEEE 9th Information Technology International Seminar*, 2023.
- [16] A. Nurkholis and S. Styawati, "Prediction model for soybean land suitability using C5. 0 algorithm," *Jurnal Online Informatika*, vol. 6, no. 2, pp. 163–171, 2021.
- [17] A. Nurkholis, Styawati, D. Alita, A. Sucipto, M. Chanafy, and Z. Amalia, "Hotspot classification for forest fire prediction using C5.0 algorithm," in *International Conference on Intelligent Cybernetics Technology & Applications (ICICyTA)*, 2021.
- [18] Y. Zhang, G. Chi, and Z. Zhang, "Decision tree for credit scoring and discovery of significant features: an empirical analysis based on Chinese microfinance for farmers," *Filomat*, vol. 32, no. 5, pp. 1513–1521, 2018.
- [19] A. Nurkholis, I. S. Sitanggang, Annisa, and Sobir, "Spatial decision tree model for garlic land suitability evaluation," *International Journal of Artificial intelligence (IJ-AI)*, vol. 10, no. 3, pp. 666–675, 2021, doi: 10.11591/ijai.v10.i3.pp666-675.
- [20] A. Nurkholis and I. S. Sitanggang, "A spatial analysis of soybean land suitability using spatial decision tree algorithm," in *Sixth International Symposium on LAPAN-IPB Satellite*, Dec. 2019, p. 65. doi: 10.1117/12.2541555.