

# Analyzing Consumer Shopping Interest via Social Media Ads with K-Means and C4.5 Algorithm

Jepri Banjarnahor<sup>[1]\*</sup>, Jessy Putrionom Hutagalung<sup>[2]</sup>, Ferdinand Jery Wilkinson Sitorus<sup>[3]</sup>

Department of Sciences and Technology, Universitas Prima Indonesia Medan <sup>[1],[2],[3]</sup>  
Medan - Indonesia

[jepribanjarnahor@unprimdn.ac.id](mailto:jepribanjarnahor@unprimdn.ac.id)<sup>[1]</sup>, [jessyhutagalungg@gmail.com](mailto:jessyhutagalungg@gmail.com)<sup>[2]</sup>, [ferdinandsitorus424@gmail.com](mailto:ferdinandsitorus424@gmail.com)<sup>[3]</sup>

**Abstract**— It is increasingly important to understand how advertisements affect consumers' propensity to shop as social media becomes the primary medium for advertising. This study uses the C4.5 algorithm for classification and K-Means Clustering for data segmentation to examine the level of consumer shopping interest driven by Facebook and Instagram ads. This strategy utilizes information collected from user interactions with ads on these two social media platforms to determine consumer interest trends more precisely. The research findings show that, compared to conventional methods, this combination of techniques can increase the accuracy of predicting consumer purchase intention by as much as 85%. These results not only validate the usefulness of clustering and classification methods in digital advertising data analysis, but also offer insights that companies can apply to optimize their marketing strategies. By understanding more specific consumer segments, companies can target their ads more precisely, thereby increasing conversions and the effectiveness of advertising campaigns. This research makes a significant contribution to the field of data analysis and digital marketing and opens up opportunities for further research in the integration of more sophisticated analysis methods.

**Keywords** : Social media, K-means clustering, C4.5 Algorithm, marketing, consumer shopping interest

## I. INTRODUCTION

In today's digital age, social media has become a key platform for companies to reach out to their consumers. Ads on platforms such as Facebook and Instagram not only increase brand visibility but also influence consumer shopping behavior. Through the advertising features provided by both platforms, businesses can promote their products or services to the right segmented target audience. Advertising is a form of communication used by individuals, companies, or organizations to promote products, services, or ideas to the wider community[1]. Social media is a digital platform that allows users to create, share, and interact with content and communicate with others. Through social media, users can participate in social networks, online communities, and discussion forums that bring together individuals with similar interests[2].

Previous research has explored the use of various methods to analyze the influence of ads on social media on consumer shopping interest. However, most of these studies face limitations in terms of data segmentation accuracy and

consumer interest predictions that often do not match real behavior. This research contributes by introducing a new approach that combines K-Means Clustering for more accurate data segmentation and C4.5 algorithm for classification, which significantly improves the ability to predict consumer shopping interest more accurately. The K-Means method and C4.5 algorithm have been widely used in various studies for data segmentation and classification. Digital advertising on social media platforms such as Facebook and Instagram has a significant impact on consumer purchasing decisions in Indonesia, especially through increased brand visibility and the appeal of advertising content tailored to consumer preferences[3]

To overcome these challenges, using K-means clustering and C4.5 algorithm methods is an effective solution. By using K-means clustering, which is one of the most common clustering algorithms for grouping data according to similar characteristics[4]. Companies can group the type of increase in consumer interest in buying within 1 month based on the specifications of advertisements displayed on Facebook and Instagram social media, thus allowing companies to understand consumer groups who respond to advertisements by viewing the advertisements displayed. In addition, with the C4.5 algorithm, this method is used to perform a series of classification problems in machine learning and data mining[5]. Companies can classify the level of interest in future consumer spending based on the patterns identified, allowing companies to more precisely target advertisements to potential consumers who have a high level of interest. By using these two methods together, companies can gain deeper insights into consumer responses to advertisements on the Facebook and Instagram platforms, allowing them to make smarter and more effective marketing decisions.

## II. METHODOLOGY RESEARCH

In conducting research, the method used is the R&D method. Research and development (R&D) is a systematic process involving investigation and experimentation to develop new knowledge and create innovation. R&D is carried out by companies or organizations to produce new and better products, services, or technologies or to improve existing ones. R&D activities include basic and applied research aimed at

developing new solutions and increasing competitiveness in the market [6].

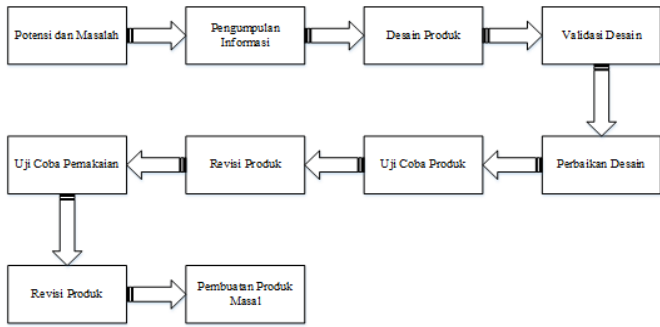


Figure 1: Stages of Research and Development (R&D)

The research and development phase of developing an analytical model to optimize social media marketing strategies involved several steps. First, we identified the opportunities and challenges associated with using Facebook and Instagram as effective advertising platforms and analyzed consumer shopping interest from ad interactions. Data from user interactions with ads was then collected. Consumer data is clustered in this study according to how they interact with advertisements on social media, especially Facebook and Instagram, using the K-Means Clustering technique. The selection of K-Means was based on its efficacy in managing large and diverse datasets. The C4.5 method was then used to categorize the resulting groups according to the level of consumer interest in shopping after clustering. The new methodology that combines these two techniques results in more accurate shopping interest predictions of up to 85%. Data on customer interactions was collected first, then the data was cleaned and prepared for analysis. After K-Means clustering, C4.5 was used to further process the data to classify the results. The cross-validation method is used to assess the resulting model, and the model is contrasted with conventional baseline data. This method works by identifying a number of centroids in the data, and then allocating each data to its closest centroid [7]. The C4.5 algorithm creates a top-down decision tree, where the topmost attribute is the root, and the bottom is called the leaf [8]. Applied the C4.5 algorithm to e-commerce data in order to predict consumer purchase intentions, demonstrating its effectiveness in classifying consumer behavior based on various product and advertisement attributes[9]. By using a sample, researchers can make inferences or conclusions about the population as a whole [10].

### III. RESULT AND DISCUSSION

This research analyzes the level of consumer shopping interest through advertisements on Facebook and Instagram social media using the K-means clustering and C4.5 algorithm methods. Practically, the results of this study can be used by advertisers to optimize marketing strategies on Facebook and Instagram [11]. Optimized marketing strategies by applying the K-Means algorithm to segment social media data, which allowed companies to better understand consumer groups and enhance the precision of their ad targeting[12]. Analyzed the

influence of Instagram advertisements on consumer loyalty using data mining techniques, highlighting the importance of personalized advertising strategies in fostering long-term customer engagement[13].

#### A. Sampel Data

By using samples, researchers can make inferences or conclusions about the population as a whole. In this study, using a sample of advertising specification data on social media Facebook and Instagram [14].

#### 1) Dataset

Table I. Dataset

Month	Year	Information	Image	Price	Promotion
January	2023	80	70	80	80
February	2023	75	80	89	76
March	2023	86	85	77	80
April	2023	75	80	92	87
May	2023	86	89	89	75
June	2023	87	75	77	79
July	2023	88	65	79	60
August	2023	70	89	80	90
September	2023	75	90	80	80
October	2023	92	78	84	80
November	2023	87	89	78	96
December	2023	70	87	87	80
January	2024	75	92	88	98
February	2024	78	80	79	67
March	2024	90	85	90	80
April	2024	80	80	92	86
May	2024	81	82	78	80

Datasets are used in various fields, including data science, statistics, machine learning, and scientific research, for analysis and decision making [15]. The datasets taken are specifications of advertisements displayed on Facebook and Instagram social media, while the datasets used in this study are as follows:

#### 2) Data Analysis

This process involves various techniques and methods applied to raw data to extract meaningful and relevant insights [16]. To determine the number of clusters, the sample data will be divided into two clusters (groups). Set two data points from the dataset as initial cluster centers.

Table II. Initial Cluster

C1	Up	75	80	92	87
C2	Down	86	85	77	80

After setting the k value and the initial cluster center, the next step is to calculate the distance between each student's data and the cluster center:

$$C1 = \sqrt{(80 - 75)^2 + (70 - 80)^2 + (80 - 92)^2 + (80 - 87)^2}$$

$$C1 = \sqrt{(5)^2 + (-10)^2 + (-12)^2 + (-7)^2}$$

$$C1 = \sqrt{25 + 100 + 144 + 49}$$

$$C1 = \sqrt{318}$$

$$C1 = 17,83255$$

$$C2 = \sqrt{(80 - 86)^2 + (70 - 85)^2 + (80 - 77)^2 + (80 - 80)^2}$$

$$C2 = \sqrt{(-6)^2 + (-15)^2 + (-3)^2 + (0)^2}$$

$$C2 = \sqrt{36 + 225 + 9 + 0}$$

$$C2 = \sqrt{270}$$

$$C2 = 16,43168$$

The closest cluster distance for each record can be determined. Determine the cluster for each record and update the cluster center point. After the calculation of all the data is complete, the results can be seen in the table below.

Table III. Iterations 1

Data to-	C1	C2	Cluster
1	17,83255	16,43168	C2
2	11,40175	17,49286	C1
3	20,4939	0	C2
4	0	20,4939	C1
5	18,84144	13,60147	C2
6	21,40093	10,0995	C2
7	35,9444	28,42534	C2
8	16,09348	19,51922	C1
9	17,11724	12,4499	C2
10	20,14944	11,57584	C2
11	22,40536	16,55295	C2
12	12,16553	18,97367	C1
13	16,76305	24,79919	C1
14	24,04163	16,18641	C2
15	17,4069	13,60147	C2
16	5,09902	17,94436	C1
17	16,88194	5,91608	C2

Table IV. Cluster Grouping

Data to-	C1	C2
1	17,83255	16,43168
2	11,40175	17,49286

3	20,4939	0
4	0	20,4939
5	18,84144	13,60147
6	21,40093	10,0995
7	35,9444	28,42534
8	16,09348	19,51922
9	17,11724	12,4499
10	20,14944	11,57584
11	22,40536	16,55295
12	12,16553	18,97367
13	16,76305	24,79919
14	24,04163	16,18641
15	17,4069	13,60147
16	5,09902	17,94436
17	16,88194	5,91608

To update the cluster center point value, the cluster center equation can be used as follows:

$$C_{1,1} = \frac{75+75+70+70+75+80}{6}$$

$$= 74,16667$$

$$C_{1,2} = \frac{80+80+89+87+92+80}{6}$$

$$= 84,66667$$

$$C_{1,3} = \frac{89+92+80+87+88+92}{6}$$

$$= 88$$

$$C_{1,4} = \frac{76+87+90+80+98+86}{6}$$

$$= 86,16667$$

$$C_{2,1} = \frac{80+86+86+87+88+75+92+87+78+90+81}{11}$$

$$= 84,54545$$

$$C_{2,2} = \frac{70+85+89+75+65+90+78+89+80+85+82}{11}$$

$$= 80,72727$$

$$C_{2,3} = \frac{80+77+89+77+79+80+84+78+79+90+78}{11}$$

$$= 81$$

$$C_{2,4} = \frac{80+80+75+79+60+80+80+96+67+80+80}{11}$$

$$= 77,90909$$

From the above calculations, the latest cluster point center results are obtained:

Table V. New Cluster

C1	Up	74,16667	84,66667	88	86,16667
C2	Down	84,54545	80,72727	81	77,90909

After the cluster grouping is obtained in the 3rd iteration, the next calculation process is carried out using the C4.5 algorithm to determine the cluster in the following month. Yes, the explanation can be seen as follows.

Table VI. K-Means Clustering Result

No	Month	Year	Cluster
1	January	2023	C2
2	February	2023	C1
3	March	2023	C2
4	April	2023	C1
5	May	2023	C2
6	June	2023	C2
7	July	2023	C2
8	August	2023	C1
9	September	2023	C1
10	October	2023	C2
11	November	2023	C1
12	December	2023	C1
13	January	2024	C1
14	February	2024	C2
15	March	2024	C2
16	April	2024	C1
17	May	2024	C2
18	June	2024	?

From the table above, what is taken to become training data for data testing, that is, only the attribute data, as for the attribute data taken, A1, A2, A3, A4. Where the training data has 17 records, it will be classified with the testing data, which has 1 record. For example, the calculation can be explained as follows:

$$\text{EntropyA}^1 (>=80) = (-\frac{2}{10} \times \text{Log}_2(\frac{2}{10})) + (\frac{8}{10} \times \text{Log}_2(\frac{8}{10})) = 0,7219$$

$$\text{EntropyA}^1 (<80) = (-\frac{6}{7} \times \text{Log}_2(\frac{6}{7})) + (\frac{1}{7} \times \text{Log}_2(\frac{1}{7})) = 0,5917$$

$$\text{GainA}^1 = 0,997 - ((\frac{10}{17} \times 0,7219) - ((\frac{7}{17} \times 0,5917)) = 0,3292$$

And so on.

Table VII. Node

Node	Attributes	Value	Case	C1	C2	Entropy	Gain
1	Total		17	8	9	0,997502546	
	A <sup>1</sup>	>=80	10	2	8	0,721928095	0,329209
		<80	7	6	1	0,591672779	
	A <sup>2</sup>	>=80	13	8	5	0,961236605	0,262439
		<80	4	0	4	0	
	A <sup>3</sup>	>=80	11	7	4	0,945660305	0,156185
		<80	6	1	5	0,650022422	
	A <sup>4</sup>	>=80	12	7	5	0,979868757	0,093499
		<80	5	1	4	0,721928095	

From the node table above, which has gone through several stages in the formation of nodes, a decision tree is formed, which can be seen below:

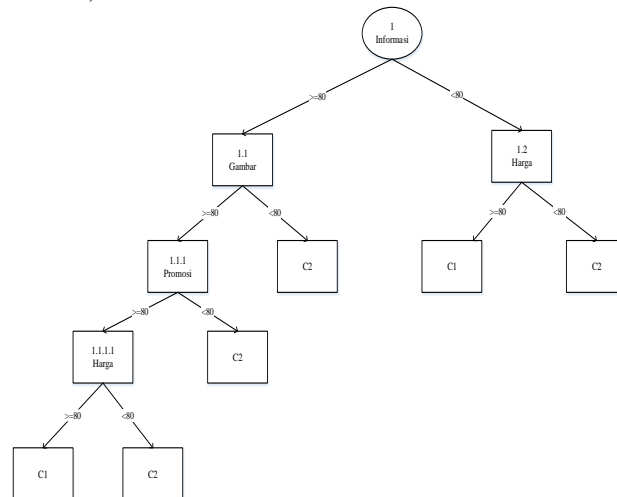


Figure 2 Decision Tree

Table VIII. C4.5 Classification

No	Month	Year	Information	Image	Price	Promotion	Interest
18	June	2024	85	87	70	89	C2

So it can be explained in the 18th data the results of the classification of the C4.5 algorithm, then for data 18 enter the cluster "C2".

### 3) System Implementation

System implementation refers to applying or implementing a new system or changes to an existing system in a production or operational environment. This involves hardware setup, software installation, network configuration, and user testing and training [17]. The explanation can be seen as follows:

```

Test mode: evaluate on training data
--- Clustering model (full training set) ---

kMeans
=====
Number of iterations: 4
Within cluster sum of squared errors: 3.7643443046637324

Initial starting points (random):
Cluster 0: 75,80,92,67
Cluster 1: 70,89,80,90

Missing values globally replaced with mean/mode

Final cluster centroids:
Attribute Full Data Cluster 0 Cluster 1
(17.0) (7.0) (10.0)
-----
A1 80.8824 78.7143 82.4
A2 81.8824 84.7143 79.9
A3 83.4706 89.5714 79.2
A4 80.8289 83.4289 79.2

Time taken to build model (full training data) : 0 seconds

--- Model and evaluation on training set ---
Clustering Instances
0 7 (41%)
1 10 (59%)
    
```

Figure 3: Weka K-Means Clustering Classification Results

The Weka analysis above uses the K-means algorithm to group the data into two clusters. The model went through four iterations with a sum of squared errors value of

3.7644334066637334. The starting point of clustering was initialized randomly. The final cluster shows that:

1. Cluster 0 has centroids A1: 78.7143, A2: 71.4286, A3: 65.5714, and A4: 78.4286.
2. Cluster 1 has centroids A1: 83.1429, A2: 79.2857, A3: 74.7143, and A4: 82.4286.

A total of 41% of the data is included in Cluster 0 and 59% in Cluster 1. This model takes 0 seconds.

No.	1: Instance_number	2: A1	3: A2	4: A3	5: A4	6: Cluster
	Numeric	Numeric	Numeric	Numeric	Numeric	Nominal
1	0.0	80.0	70.0	80.0	80.0	cluster1
2	1.0	75.0	80.0	89.0	76.0	cluster0
3	2.0	86.0	85.0	77.0	80.0	cluster1
4	3.0	75.0	80.0	92.0	87.0	cluster0
5	4.0	86.0	89.0	89.0	75.0	cluster0
6	5.0	87.0	75.0	77.0	79.0	cluster1
7	6.0	88.0	65.0	79.0	60.0	cluster1
8	7.0	70.0	89.0	80.0	90.0	cluster1
9	8.0	75.0	90.0	80.0	80.0	cluster1
10	9.0	92.0	74.0	84.0	80.0	cluster1
11	10.0	87.0	89.0	78.0	96.0	cluster1
12	11.0	70.0	87.0	87.0	80.0	cluster0
13	12.0	75.0	92.0	88.0	98.0	cluster0
14	13.0	78.0	80.0	79.0	67.0	cluster1
15	14.0	90.0	85.0	90.0	80.0	cluster0
16	15.0	80.0	80.0	92.0	86.0	cluster0
17	16.0	81.0	82.0	78.0	80.0	cluster1

Figure 3 Final Result of Weka Clustering

The figure above explains the results of the comparison between manual calculations and calculations using the Weka application, there are 'Y' errors totaling '4', while valid 'T' totals '13', more valid results than error results, so the results of this calculation can be used to help companies adjust their ad targeting more precisely.

```

Size of the tree : 5

Time taken to build model: 0 seconds

=== Stratified cross-validation ===
=== Summary ===

Correctly Classified Instances      15      83.3333 %
Incorrectly Classified Instances    3      16.6667 %
Kappa statistic                    0.6582
Mean absolute error                 0.1667
Root mean squared error             0.383
Relative absolute error             33.5294 %
Root relative squared error        76.61 %
Total Number of Instances          18

=== Detailed Accuracy By Class ===

          TP Rate  FP Rate  Precision  Recall  F-Measure  MCC      ROC Area  PRC Area  Class
0.900  0.250  0.818  0.900  0.857  0.663  0.913  0.919  C2
0.750  0.100  0.857  0.750  0.800  0.663  0.913  0.825  C1
Weighted Avg.  0.833  0.183  0.835  0.833  0.832  0.663  0.913  0.877

=== Confusion Matrix ===

a b  <-- Classified as
9 1 | a = C2
2 6 | b = C1
    
```

Figure 4: Classify the C4.5 algorithm Weka Result

The Weka analysis results above show the performance of the classification model using the stratified

cross-validation method. Out of 18 instances, the model classified 83.333% (15 instances) correctly and 16.6667% (3 instances) incorrectly. The kappa statistic obtained was 0.6582. The resulting error values are as follows: mean absolute error (0.1667), root mean squared error (0.383), relative absolute error (33.5294%), and root relative squared error (76.61%). The accuracy details show that for class C2, the precision value is 0.818, recall 0.900, and F-measure 0.857. For class C1, the precision value is 0.857, recall is 0.750, and F-measure is 0.800. The area under the ROC curve is 0.913 for both classes. The confusion matrix shows that 9 instances of class C2 and 6 instances of class C1 are correctly classified, while 2 instances of class C1 are classified as C2.

No.	1: A1	2: A2	3: A3	4: A4	5: prediction margin	6: predicted Cluster	7: Cluster
	Numeric	Numeric	Numeric	Numeric	Numeric	Nominal	Nominal
1	70.0	89.0	80.0	90.0	1.0	C1	C1
2	90.0	85.0	90.0	80.0	1.0	C2	C2
3	75.0	92.0	88.0	98.0	1.0	C1	C1
4	86.0	85.0	77.0	80.0	1.0	C2	C2
5	75.0	80.0	89.0	76.0	1.0	C1	C1
6	81.0	82.0	78.0	80.0	1.0	C2	C2
7	75.0	90.0	80.0	80.0	1.0	C1	C1
8	85.0	87.0	70.0	89.0	-1.0	C1	C2
9	70.0	87.0	87.0	80.0	1.0	C1	C1
10	92.0	74.0	84.0	80.0	1.0	C2	C2
11	75.0	80.0	92.0	87.0	1.0	C1	C1
12	88.0	65.0	79.0	60.0	1.0	C2	C2
13	87.0	89.0	78.0	96.0	-0.8	C2	C1
14	87.0	75.0	77.0	79.0	0.8	C2	C2
15	80.0	80.0	92.0	86.0	-0.8	C2	C1
16	80.0	70.0	80.0	80.0	0.8	C2	C2
17	78.0	80.0	79.0	67.0	1.0	C2	C2
18	86.0	89.0	89.0	75.0	1.0	C2	C2

Figure 5: Final Result Classify Weka

The figure above explains the results of the comparison between manual calculation and calculation using the Weka application. In the 18th data, the results of manual classification and Weka have similarities for the results, namely in the 18th data into the "C2" cluster, so that the results of this calculation can be used in helping companies adjust their ad targeting more precisely.

#### ACKNOWLEDGMENT

Thank you to the Prima Indonesia University campus for providing the opportunity to complete this research and create this paper.

#### REFERENCES

- [1] Armanto, R., & Gunarto, M. (2022). Analysis of the impact of social media on housing sales: An empirical study of Facebook and Instagram ad usage. *Journal of Business, Management and Economics*, 3(1), e-ISSN 2745-7281.
- [2] Hartawan, E., Liu, D., Handoko, M. R., Evan, G., & Widjojo, H. (2021). The effect of advertising on Instagram social media on public buying interest in e-commerce. *Sam Ratulangi University Scientific Journal of Business Management and Innovation*, 8(1), 217-228.
- [3] Santoso, D., & Harsono, S. (2023). Analysis of the Effect of Digital Advertising on Social Media on Consumer Purchasing Decisions in Indonesia. *Journal of Digital Economics and Business*, 4(2), 102-115.
- [4] Sembiring, C. S. D., Hanum, L., & Tamba, S. P. (2022). Application of data mining using the K-Means algorithm to determine thesis titles and research journals (Case study of FTIK UNPRI). *JUSIKOM PRIMA (Journal of Information Systems and Computer Science Prima)*, 5(2).
- [5] Anestiviya, V., & Pasaribu, A. F. O. (2021). Pattern Analysis Using the

- C4.5 algorithm Method for Student Major Specialization Based on Curriculum (CASE STUDY: SMAN 1 NATAR). *Journal of Information Technology and Systems (JTSI)* Vol. 2, No. 1, March 2021, 80 - 85.
- [6] Fransisca, S., & Putri, R. N. (2019). Utilization of RFID technology for school inventory management with (R&D) method (Case study: SMK Global Pekanbaru). *Journal of Computer and Information Technology Application Students*, 1(1), 72-75.
- [7] Mawarni, Q. I., & Budi, E. S. (2022). Implementation of K-Means clustering algorithm in student discipline assessment. *Journal of Computer and Information Systems (JSON)*, 3(4), 522-528.
- [8] Anestiviya, V., & Pasaribu, A. F. O. (2021). Pattern Analysis Using the C4.5 algorithm Method for Student Major Specialization Based on Curriculum (CASE STUDY: SMAN 1 NATAR). *Journal of Information Technology and Systems (JTSI)* Vol. 2, No. 1, March 2021, 80 - 85.
- [9] Nurhadi, A. P., & Kurniawati, R. (2023). Application of C4.5 Algorithm on E-commerce Data to Predict Consumer Purchase Interest. *Journal of Information Technology and Systems*, 7(3), 87-98.
- [10] Susanto, P. C., Arini, D. U., Yuntina, L., Soehaditama, J. P., & Nuraeni. (2024). Quantitative research concepts: Population, sample, and data analysis (A literature review). *Journal of Computer Science*, 3(1).
- [11] Utomo, D. P., & Mesran. (2020). Comparative analysis of data mining classification methods and attribute reduction on heart disease datasets. *JOURNAL OF BUDIDARMA INFORMATICS MEDIA*, 4(2), 437-444.
- [12] Mulyadi, T., & Saputra, H. (2022). Marketing Strategy Optimization with Data Clustering on Social Media Using K-Means Algorithm. *Journal of Computer Science and Information Systems*, 6(1), 45-58.
- [13] Wijaya, S., & Dewi, L. (2023). Analysis of the Effect of Instagram Ads on Consumer Loyalty with a Data Mining Approach. *Indonesian Journal of Business and Management*, 9(2), 112-128.
- [14] Wiratama, M. A., & Pradnya, W. M. (2022). Optimization of data mining algorithms using backward elimination for diabetes disease classification. *National Journal of Informatics Engineering Education: JANAPATI*, 11(1).
- [15] Pramudito, D. K. (2022). Data mining implementation on Java North Coast weather forecast dataset using C4.5 algorithm. *Journal of Technology Pelita Bangsa*, 13(3).
- [16] Pambudi, A., Abidin, Z., & Permata. (2023). Application of CRISP-DM using MLR K-Fold on PT Telkom Indonesia (Persero) Tbk (TLKM) stock data (Case study: Indonesia Stock Exchange 2015-2022). *JDMISI*, 4(1), 1-14.
- [17] Rahmadi, M., Kaurie, F., & Susanti, T. (2020). Test the accuracy of postoperative patient dataset using Naïve Bayes algorithm using Weka tools. *JURIKOM (Journal of Computer Research)*, 7(1).