# Analyzing the Distribution of Health Workers in Semarang City Using K-Means Clustering Method

## Akhfan Setiyaji[1], Hindriyanto Dwi Purnomo[2]

[1,2]Informatics Department, Faculty of Information Technology, Satya Wacana Christian University
Diponegoro Street No. 52 – 60, Salatiga, Central Java, Indonesia 50711
E-mail: [1]672018348@student.uksw.edu, [2]hindriyanto.purnomo@uksw.edu

### Abstract

This research employed the K-Means Clustering method to examine the distribution of health workers in Semarang City, emphasizing their pivotal role in the public health infrastructure. Leveraging current data encompassing health worker locations and quantities, the clustering analysis discerned areas exhibiting similar distribution characteristics through the application of the K-Means technique. Quantitative analysis revealed distinct clusters, shedding light on the spatial patterns of health workforce dispersion within Semarang City. The study's quantitative findings furnish valuable insights crucial for formulating more efficacious health policies. By delineating the utility of the K-Means Clustering method in public health planning and providing quantitative evidence of health worker distribution, this research substantially augments geographical comprehension in the examined region.

**Keywords**: K-Means Clustering, Health, Community, Distribution, Health Workers, Spatial Analysis, Semarang City.

## 1. INTRODUCTION

Health workers have a very important role in providing services at the Community Health Center. Puskesmas as the main door for health services to the community must be able to provide optimal basic health services and in accordance with competency standards. Health workers according to Health Law no. 36 of 2009 is someone who has the knowledge, skills and permission to carry out health actions or efforts and is willing to dedicate himself to society in the health sector [1]. Meanwhile, based on Minister of Health Regulation Number 75 of 2014, it is stated that health workers who work based on staff standards at community health centers have at least 9 different types of health workers [2].

The ratio of health workers in the distribution of health workers is still significantly large, leading to frequent instances where Community Health Centers in Semarang City fail to meet standards. Thus, there is a pressing need for policy interventions regarding the planning and procurement of health workers to ensure that their distribution aligns with competency standards, particularly in Semarang City. One

301

of the classic challenges encountered by Semarang City is the unequal availability and distribution of health workers at the basic service level, particularly within Community Health Centers, hindering efforts to achieve fair and equitable health services for the community.

The occupancy of Semarang City health centers with 9 (nine) types of complete health personnel according to the Minister of Health Regulation is apparently still below 50%. There are 1,513 Community Health Centers that do not have doctors at all (15%). SISDMK data also shows that there are community health centers lacking dentists (45.53%), community health centers lacking nurses (2.4%), community health centers lacking midwives (3.4%), community health centers lacking public health personnel (29.47%), Community Health Centers lack environmental health personnel (26.9%), Community Health Centers lack medical laboratory experts - ATLM (35.01%), Community Health Centers lack nutrition personnel (24%) and Community Health Centers lack pharmaceutical personnel (22.88%).

Data analysis can be done in various ways, one of which is by using Data Mining. Data mining is a data processing method to find important patterns hidden in data. The results of data processing can become useful information for the future. Data mining has various methods or even models that can be used. For this research we used clustering techniques. Clustering is a method of grouping data, where each data will be combined into groups that have similar data characteristics to each other [3].

One of the well-known algorithms in the clustering method is K-Means. K-Means is a method for analyzing data. This algorithm determines the number and value of clusters (k) randomly. This value is the initial center of the cluster or can be called the centroid [4]. In clustering methods, the K-Means algorithm is popular because the algorithm is relatively simple and efficient to use. This algorithm is included in the unsupervised learning category where we do not need to carry out a training process or in other words there is no learning stage so it can be applied in various fields. K-Means divides data into k separate regions where k is a positive integer. The K-Means algorithm is very famous because of its simplicity and ability to classify large data and outliers very quickly [5].

Research using the K-Means algorithm with clustering techniques has often been used, one of which is for the problem of the distribution of districts and cities based on data on the number of health workers in Banten Province [6]. analyze and classify data using the K-Means algorithm to be able to group the number of health workers in Karawang Regency. So this research aims to use the K-Mens Clustering method to determine the distribution of health workers in Semarang City as a novelty in research.

## 2.  METHODS

The research method used in this research is a quantitative method. Quantitative data is research that aims to examine a certain population or sample and random sampling through the use of tools to collect data, and statistical data analysis [17]. The research method explains the stages used in the research, namely by using the K-Means clustering algorithm and the Davies Bouldin Index (DBI) cluster evaluation.

### 2.1  Research Strategy and Design

This research adopts the waterfall system development method, renowned for its structured and sequential approach to analyzing systems. In this method, each step must be completed before proceeding to the next; hence, if any preceding step remains unfinished, subsequent processes cannot commence. Consequently, Process 4 cannot initiate unless Processes 1, 2, and 3 have been executed. The waterfall method encompasses distinct phases: Analysis, Design, Code, Testing, Implementation, and Maintenance. Pressman describes the Waterfall Method as a classic, systematic, and linear sequential model for software development, also known as the "Linear Sequential Model". The author's choice of the waterfall method stems from its structured nature, which ensures software quality and facilitates easy maintenance. This model is advantageous for users as it allows comprehensive planning and preparation of necessary data and process requirements from the outset, enabling smooth and problem-free execution.
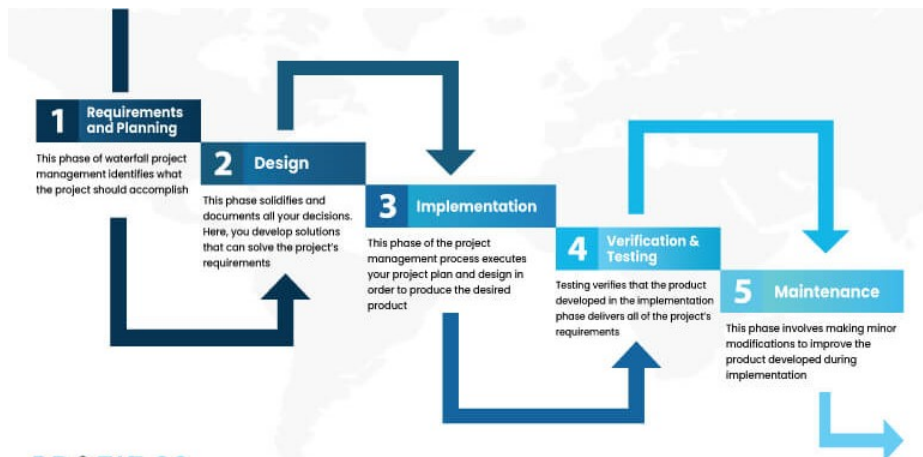


**Figure 1.** Waterfall Model

The following are the stages of the Waterfall Model as in Figure 1 according to Roger S Pressman [18]:

1) Requirements, this stage starts from collecting requirements to be applied to the device system to be developed. This stage collects the required requirements with incentives to be understood by the user, these needs are intended for the user as a future user of the system.
2) Design, this stage designs the interface design requirements for the system to be developed. Interface design needs to be done when you want to develop the system.
3) Implementation, this program creation stage must be implemented using software.
4) Verification, this system testing stage is carried out logically and functionally which is used to find out the parts of the system that have been tested previously. This testing is to minimize errors that occur in the program.
5) Maintenance, the system maintenance or maintenance phase needs to be carried out because system requirements can continue to be improved after use.

## 2.2 Data CollectionT

According to [19] the subject is one part or member of the sample. Research subjects are parties who are used as sources of information or data sources in a study. This research took data in the city of Semarang. This stage is carried out by direct observation of the object being studied. This will observe how health worker data is managed, so that a soft copy of the 2023 Integrated Data is obtained. The 2023 Health Worker data is given to Semarang City Health Workers.

## 2.3 K-Means Clustering

After the data is collected, the next stage is the analysis stage. At this stage what is done is to process the data collected from the previous stages. The K-Means Clustering stages are divided into several processes, namely:
1) Data Preprocessing
   The stage before data processing in data mining is better known as the data preprocessing stage, while the benefit of this process is to improve the quality of data grouping results using the K-means algorithm. An explanation of the data preprocessing stages is as follows:
   a) Attribute Selection
      In collecting data from bpskotasemarang.com, there were 10 health worker variables in the data. Next, the selection of attributes is carried out by sorting the attributes according to research needs with the aim of reducing the scope of the research. There are 3 attributes used, namely Doctor, Nurse and Midwife.
   b) Data Cleaning
      Cleaning is carried out to remove data that is not suitable for entering the data mining process, such as noise data and missing values. In this

research, data that was not appropriate to the research was also deleted, such as data filled in by respondents who were outside the research case study.

c) Data Transformation

At this stage, data is changed from data types that initially cannot be processed mathematically into data that can be processed. Data transformation aims to avoid damaged and invalid data. In the K-Means process, alpha numeric (text) data is transformed into numeric.

2) K-Means Clustering Calculation Process

At this stage the completion of the K-Means Clustering algorithm is carried out using a predetermined formula. As explained in the Introduction Chapter, this research aims to determine the distribution of Semarang City Health Workers using K-Means Clustering. The way it works is by dividing the data into clusters through a systematic process. After obtaining the clusters, an analysis of the cluster formation patterns is carried out which can produce useful information. The entire data calculation process uses manual calculations and then continues with testing with Rapid Miner 8.1 software to reduce the impact of human error that occurs during manual calculations.

3) System Testing

This testing stage focuses on functional and non-functional testing and will later carry out User Acceptance Testing with end users to verify whether the designed system meets their needs.

## 3. RESULTS AND DISCUSSION

### 3.1 Data Preprocessing

The data sourced from bpskotasemarang.com comprised 16 datasets detailing the total count of healthcare workers in each sub-district within Semarang City. Within these 16 sub-district datasets, emphasis was placed on three key variables: Doctors, Nurses, and Midwives. These variables were identified as crucial for the community in accessing essential healthcare services.

Next, the data goes through the Data Processing stage to produce a dataset as in Table 1.

**Table 1.** Semarang City Health Worker Dataset

| Subdistrict | Doctor | Nurse | Midwife |
|---|---|---|---|
| Mijen | 9 | 29 | 18 |
| Gunung Pati | 7 | 14 | 10 |
| Banyumanik | 166 | 155 | 29 |
| Gajah Mungkur | 248 | 419 | 69 |
| Semarang Selatan | 436 | 2123 | 165 |
| Candisari | 8 | 9 | 7 |
| Tembalang | 165 | 405 | 91 |

| Subdistrict | Doctor | Nurse | Midwife |
|---|---|---|---|
| Pedurungan | 72 | 25 | 22 |
| Genuk | 142 | 547 | 45 |
| Gayamsari | 61 | 117 | 7 |
| Semarang Timur | 160 | 219 | 80 |
| Semarang Utara | 9 | 10 | 4 |
| Semarang Tengah | 221 | 663 | 51 |
| Semarang Barat | 30 | 48 | 17 |
| Tugu | 9 | 14 | 10 |
| Ngaliyan | 123 | 244 | 47 |

### 3.2  K-Means Clustering Calculation

This stage, calculations are carried out using the K-Means algorithm based on the dataset that has been obtained. In calculating the data, the data will be grouped into 3 clusters based on the small number of health workers (C1), medium cluster (C2), and large cluster (C3). In the calculation results, 3 iterations or repetitions are carried out until a value of 0 is obtained or the same value as the previous iteration. Meanwhile, the determination of the cluster center is as follows.
1) Centroid 1, taken from Candisari District data with Doctor, Nurse, and Midwife variables (8, 9, 7).
2) Centroid 2, taken from Genuk District data with Doctor, Nurse, and Midwife variables (142, 547, 45).
3) Centroid 3, taken from South Semarang District data with Doctor, Nurse, and Midwife variables (436, 2123, 165).

The distance between each data and the initial cluster center that has been determined is calculated. Then clusters 1, 2, and 3 are grouped based on the minimum distance of each cluster.

**Table 2.** First Iteration Cluster Center

| Cluster center | Doctor | Nurse | Midwife |
|---|---|---|---|
| C1 | 8 | 9 | 7 |
| C2 | 142 | 547 | 45 |
| C3 | 436 | 2123 | 165 |

**Table 3.** First Iteration Results

| Subdistrict | Doctor | Nurse | Midwife | C1 | C2 | C3 | MIN |
|---|---|---|---|---|---|---|---|
| Mijen | 9 | 29 | 18 | 22.84732 | 535,483 | 2142,142 | 22.84732 |
| Gunung Pati | 7 | 14 | 10 | 5.91608 | 550.9437 | 2157,764 | 5.91608 |
| Banyumanik | 166 | 155 | 29 | 216.2499 | 393.0598 | 1991,085 | 216.2499 |
| Gajah Mungkur | 248 | 419 | 69 | 479.1075 | 167.9166 | 1717,025 | 167.9166 |
| Semarang Selatan | 436 | 2123 | 165 | 2162,671 | 1607,673 | 0 | 0 |
| Candisari | 8 | 9 | 7 | 0 | 555.7373 | 2162,671 | 0 |
| Tembalang | 165 | 405 | 91 | 434.1901 | 151.0265 | 1740,816 | 151.0265 |
| Pedurungan | 72 | 25 | 22 | 67.65353 | 527.1745 | 2134,139 | 67.65353 |

| Subdistrict | Doctor | Nurse | Midwife | C1 | C2 | C3 | MIN |
|---|---|---|---|---|---|---|---|
| Genuk | 142 | 547 | 45 | 555.7373 | 0 | 1607,673 | 0 |
| Gayamsari | 61 | 117 | 7 | 120.3038 | 439.2095 | 2046,857 | 120.3038 |
| Semarang Timur | 160 | 219 | 80 | 269.3195 | 330.3528 | 1925,777 | 269.3195 |
| Semarang Utara | 9 | 10 | 4 | 3.316625 | 554.7423 | 2161,717 | 3.316625 |
| Semarang Tengah | 221 | 663 | 51 | 689.2177 | 140.4742 | 1480,142 | 140.4742 |
| Semarang Barat | 30 | 48 | 17 | 45.88028 | 512.1806 | 2119.52 | 45.88028 |
| Tugu | 9 | 14 | 10 | 5.91608 | 550.4571 | 2157,368 | 5.91608 |
| Ngaliyan | 123 | 244 | 47 | 264.6696 | 303.6017 | 1908,542 | 264.6696 |

The data is grouped into a cluster based on the data that has the smallest value from the cluster center. This aims to group clusters based on the smallest value, so that these results will later be used to determine new clusters. The results of cluster grouping can be seen in Table 4.

**Table 3.** First Iteration Smallest Value Grouping

| Subdistrict | C1 | C2 | C3 |
|---|---|---|---|
| Mijen | X | | |
| Gunung Pati | X | | |
| Banyumanik | X | | |
| Gajah Mungkur | | X | |
| Semarang Selatan | | | X |
| Candisari | X | | |
| Tembalang | | X | |
| Pedurungan | X | | |
| Genuk | | X | |
| Gayamsari | X | | |
| Semarang Timur | X | | |
| Semarang Utara | X | | |
| Semarang Tengah | | X | |
| Semarang Barat | X | | |
| Tugu | X | | |
| Ngaliyan | X | | |

Calculating a new cluster based on the smallest grouping value of the cluster in the first iteration, which can be seen in Table 5.

Table 4. Second Iteration Cluster Center

| | Doctor | Nurse | Midwife |
|---|---|---|---|
| **C1** | 59.45455 | 80.36364 | 22.81818 |
| **C2** | 194 | 508.5 | 64 |
| **C3** | 436 | 2123 | 165 |

The distance between each data and the cluster center that has been determined is calculated. Then clusters 1, 2, and 3 were grouped in the second iteration based

on the minimum distance of each cluster. The calculation results can be seen in Table 6.

**Table 5.** Second Iteration Results

| Subdistrict | Doctor | Nurse | Midwife | C1 | C2 | C3 | MIN |
|---|---|---|---|---|---|---|---|
| Mijen | 9 | 29 | 18 | 72.16023 | 516.0051 | 2142,142 | 72.16023 |
| Gunung Pati | 7 | 14 | 10 | 85.55652 | 531.4276 | 2157,764 | 85.55652 |
| Banyumanik | 166 | 155 | 29 | 130.2334 | 356.3303 | 1991,085 | 130.2334 |
| Gajah Mungkur | 248 | 419 | 69 | 390,329 | 104.6482 | 1717,025 | 104.6482 |
| Semarang Selatan | 436 | 2123 | 165 | 2081,914 | 1635,657 | 0 | 0 |
| Candisari | 8 | 9 | 7 | 89.3899 | 536.0459 | 2162,671 | 89.3899 |
| Tembalang | 165 | 405 | 91 | 348.1054 | 110.8253 | 1740,816 | 110.8253 |
| Pedurungan | 72 | 25 | 22 | 56.77315 | 500.4201 | 2134,139 | 56.77315 |
| Genuk | 142 | 547 | 45 | 474.3999 | 67.4333 | 1607,673 | 67.4333 |
| Gayamsari | 61 | 117 | 7 | 39.93528 | 417,385 | 2046,857 | 39.93528 |
| Semarang Timur | 160 | 219 | 80 | 180.5525 | 291.9285 | 1925,777 | 180.5525 |
| Semarang Utara | 9 | 10 | 4 | 88.60489 | 535.0956 | 2161,717 | 88.60489 |
| Semarang Tengah | 221 | 663 | 51 | 605.2737 | 157.3793 | 1480,142 | 157.3793 |
| Semarang Barat | 30 | 48 | 17 | 44.14551 | 491.0858 | 2119.52 | 44.14551 |
| Tugu | 9 | 14 | 10 | 84.34512 | 530.7271 | 2157,368 | 84.34512 |
| Ngaliyan | 123 | 244 | 47 | 177,1994 | 274.3907 | 1908,542 | 177,1994 |

The data is grouped into a cluster based on the data that has the smallest value from the cluster center according to the previous explanation. The grouping results can be seen in Table 7.

**Table 6.** Second Iteration Smallest Value Grouping

| Subdistrict | C1 | C2 | C3 |
|---|---|---|---|
| Mijen | X | | |
| Gunung Pati | X | | |
| Banyumanik | X | | |
| Gajah Mungkur | | X | |
| Semarang Selatan | | | X |
| Candisari | X | | |
| Tembalang | | X | |
| Pedurungan | X | | |
| Genuk | | X | |
| Gayamsari | X | | |
| Semarang Timur | X | | |
| Semarang Utara | X | | |
| Semarang Tengah | | X | |
| Semarang Barat | X | | |
| Tugu | X | | |
| Ngaliyan | X | | |

This stage will continue to be carried out until we get an iteration that matches the previous iteration. In the second iteration of the calculation, if seen based on the

grouping of data that has the smallest value, it has the same results as the previous iteration. Therefore, the dataset processing was stopped until the second iteration.

### 3.3 System Testing

Figure 2 and Figure 3 show the testing process using Rapidminer by looking for the cluster (k) valuebest, by using PerformanceVector to find the smallest DBI value. Meanwhile, Figure 5 and Figure 6 show the test results dataset which shows the numbers according to the results of manual testing that was carried out in the previous stage.



**Figure 2.** Testing with RapidMiner



**Figure 3.** RapidMiner Dataset Result

Then from the PerformaceVector test results in Figure 4 and Figure 5, the DBI values for clusters (k=3): -0.326; cluster (k=4): -0.366; and for cluster (k=5): -0.398. From the test results, it can be concluded that the cluster (k=3) is the best cluster used in this research because the resulting DBI value is the smallest among the other clusters, so it is considered good.
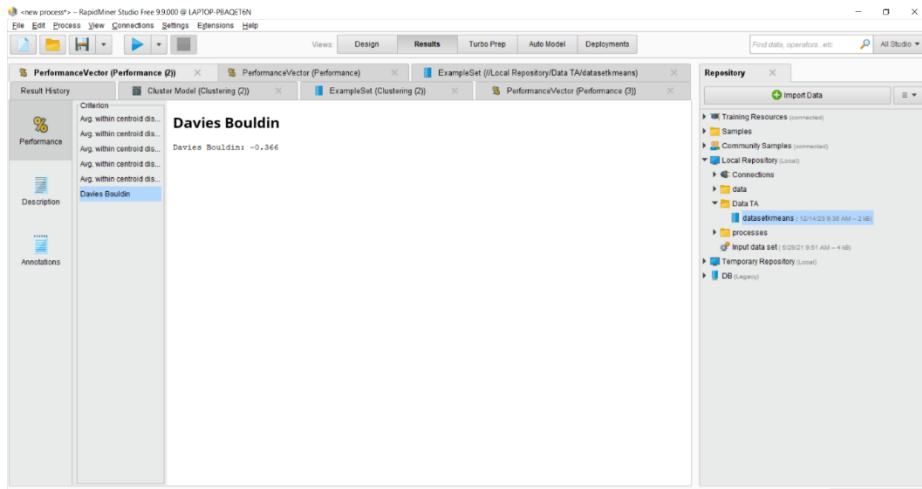


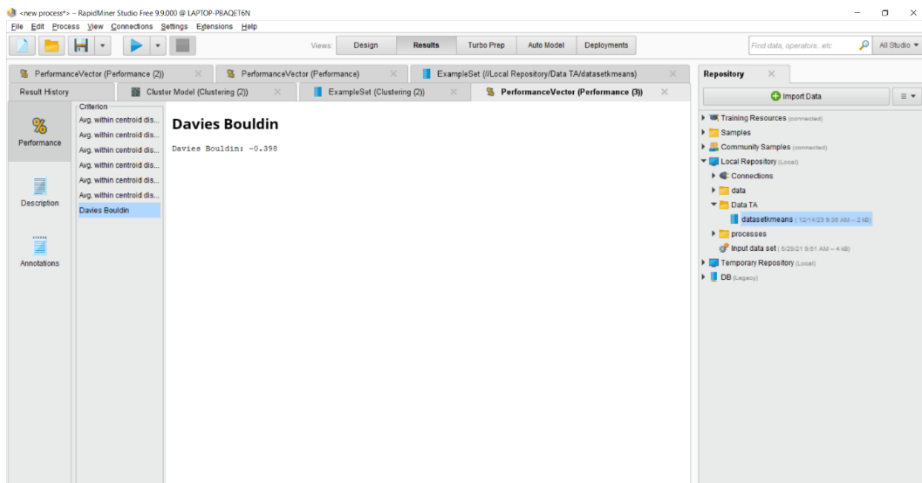**Figure 4.** DBI Cluster Value Results (K=3)



**Figure 5.** DBI Cluster Value Results (K=4)

Based on the results of testing with RapidMiner that has been carried out, it can be concluded that the amount of Health Workers in Semarang City has a large distribution in only 1 sub-district, namely South Semarang, the medium category in 4 sub-districts, namely, Gajah Mungkur, Tembalang, Genuk, and Central

Semarang. and with a small category in 11 sub-districts, namely Mijen, Gunung Pati, Banyumanik, Candisari, Pedurungan, Gayamsari, East Semarang, North Semarang, West Semarang, Tugu, and Ngaliyan.

## 4.   CONCLUSION

Based on the results of the research and trials that have been carried out, it can be concluded that the application of the K-Means Clustering algorithm in determining the distribution of the number of health workers in each sub-district in Semarang City resulted in cluster 1 with 11 sub-districts still having a shortage of health workers, including Mijen, Gunung Pati, Banyumanik, Candisari, Pedurungan, Gayamsari, East Semarang, North Semarang, West Semarang, Tugu, and Ngaliyan with a DBI value of -0.326. Then there are 4 sub-districts with a medium number of health workers, namely Gajah Mungkur, Tembalang, Genuk, and Central Semarang. Meanwhile, the number of health workers is only found in 1 sub-district, namely South Semarang. In conclusion, this research underscores the importance for the Semarang City Government and relevant agencies to prioritize the evaluation of healthcare staffing levels in accordance with standards. With the population steadily increasing and unpredictable weather conditions compromising immune systems, there is a heightened susceptibility to infectious diseases.

## REFERENCES

[1]   World Health Organization, "Health for the World's Adolescents: A Second Chance in the Second Decade," *Geneva: World Health Organization Department of Noncommunicable Disease Surveillance*, 2017.

[2]   D. Ratna Sari et al., "Application of the Naive Bayes Method in Predicting Student Satisfaction with Method Lecturer Teaching," *Proceedings of the National Information Science Research Seminar (SENARIS)*, September 2019, p. 287.

[3]   B. Serasi Ginting and M. Simanjuntak, "Grouping Diseases in Patients Based on Age Using the K-Means Clustering Method (Case Study: Bahorok Health Center)," *ALGORITHM: Journal of Computer Science and Informatics*, vol. 6341, November 2021.

[4]   N. Nugroho and F.D. Adhinata, "Use of K-Means and K-Means++ Methods for Clustering Covid-19 Data on the Island of Java," *Teknika*, vol. 11, no. 3, pp. 170–179, 2022, doi: 10.34148/teknika.v11i3.502.

[5]   I. Nisa et al., "K-Means Cluster Analysis of Health Workers in Banten Province," *J. Science and Math. Unpam*, vol. 5, no. 2, pp. 63–71, 2022.

[6]   D.K. Sitinjak, B.A. Pangestu, and B.N. Sari, "Clustering Health Workers Based on Districts in Karawang Regency Using the K-Means Algorithm," *J. Appl. Informatics Comput.*, vol. 6, no. 1, pp. 47–54, 2022, doi: 10.30871/jaic.v6i1.3855.

[7] B.L. Pailan et al., "Analysis of Health Personnel Needs Using," *Sakti - Science, Appl. Computing and Technology. Inf.*, vol. 3, no. 1, pp. 1–9, 2021.

[8] M. Zulfadhilah, Mambang, and S. Eka Prastya, "Implementation of the K-Means Clustering Method to Improve Student Networking," *Thematic*, vol. 9, no. 2, pp. 152–160, 2022, doi: 10.38204/tematic.v9i2.1053.

[9] R. Bayu Prasetyo, Y.A. Pranoto, and R.P. Prasetya, "Implementation of Data Mining Using the K-Means Clustering Algorithm for Outpatient Diseases at Dr. Clinic. Atirah, Sioyong Village, Central Sulawesi," *ITN*, 2023.

[10] R. Muliono and Z. Sembiring, "Data Mining Clustering Using the K-Means Algorithm for Lecturer Teaching Tridharma Level Clustering," *CESS (Journal of Computer Engineering, Systems and Science)*, vol. 4, no. 2, 2502–2714, 2019.

[11] G.P. Abdillah, "Application of Customer Water Use Data Mining to Determine the Potential Classification of New Customer Water Use at PDAM Tirta Raharja Using the K-Means Algorithm," *Sentika*, vol. I, p. 498, 2019.

[12] S. Hendrian, "Data Mining Classification Algorithm to Predict Students in Receiving Educational Financial Assistance," *Exacta Factors*, vol. 11, no. 3, pp. 266–274, 2018.

[13] C.H. Cheng and Y.S. Chen, "Classifying the segmentation of customer value via RFM model and RS theory," *Expert Systems with Applications*, vol. 36, no. 3, pp. 4176–4184, 2019.

[14] K. Khalili-Damghani, F. Abdi, and S. Abolmakarem, "Solving customer insurance coverage recommendation problem using a two-stage clustering-classification model," International Journal of Management Science and Engineering Management, vol. 14, no. 1, 2019, doi: 10.1080/17509653.2018.1467801.

[15] Sugiyono, Combination Research Methods (Mix Methods), Bandung: Alphabeta, 2018.

[16] R.S. Pressman, Software Engineering: A Practitioner's Approach, Yogyakarta: Andi