# INDONESIAN TERRITORY CLUSTERING BASED ON HARVESTED AREA AND RICE PRODUCTIVITY USING CLUSTERING ALGORITHM

Imelda Putri Kurniawati*, Hasih Pratiwi, Sugiyanto
Faculty of Mathematics and Natural Sciences, Universitas Sebelas Maret, Central Java, Indonesia
E-mail: imeldaputri408@student.ac.id*, hasihpratiwi@staff.uns.ac.id, sugiyanto61@staff.uns.ac.id

**ARTICLE INFO**

**ABSTRACT**

Rice (Latin: *Oryza sativa*) is one of the most important cultivated plants in civilization. This plant is the main commodity for almost all Indonesian people. Indonesia is in third place as the largest rice producing country in the world. However, based on data from the Statistics Indonesia, Indonesia will still import rice until 2022. The transfer of paddy fields is one of the reasons why Indonesia is still importing rice to this day. Many lands that used to be paddy fields have turned into airports, industrial land, housing, and so on. Rice production is one of the important topics to be discussed in order to develop rice production in areas that are still relatively low. The purpose of this research is to classify cities/regencies in Indonesia based on rice production data in 2021. In this study, three clustering methods were used, namely, Partitioning Around Medoid (PAM), Clustering Large Applications (CLARA) and Fuzzy C-Means (FCM). Then the three methods are compared based on their silhouette coefficient values. The best obtained method is FCM method with two clusters and a silhouette value of 0.828. Results *clustering* with the best method is used as a reference in making maps *clustering*. Areas that are still relatively low are expected to increase rice productivity. The PAM algorithm produces two clusters with a silhouette coefficient value of 0.58. The CLARA algorithm with 100 samples produces three clusters with a silhouette coefficient of 0.59.

## INTRODUCTION

The agricultural sector includes several sub-sectors, namely food crops, horticulture, plantations, animal husbandry, fisheries, and also forestry. The agricultural sector, especially paddy field farming, has a large multifunctional value in the context of increasing national food security, farmer welfare, and also protecting the environment (Wahyudi, 2018). One of the agricultural crops which is the main commodity for the people of Indonesia is rice, because rice is a staple food for the majority of the Indonesian population. Rice (Latin: Oryza sativa) is one of

the most important cultivated plants in civilization (Yusuf et al., 2020). Indonesia is in third place as the largest rice producing country in the world. According to the Food and Agriculture Organization (FAO), Indonesia produced 66 million metric tons of rice in 2011 (Statistic Indonesia, 2015). The second largest country is India with a production of 158 million metric tons. The country in first place is the People's Republic of China with a production of 201 million metric tons (Statistic Indonesia, 2015).

Data mining is a process for obtaining information from large databases (Tan et al., 2006). Han et al. (2022) define data mining as the process of extracting interesting patterns from large data. Data mining can also be interpreted as extracting new information taken from large chunks of data that helps in decision making (Turban, et al., 2007). Data Mining is the process of finding something meaningful from a new correlation, existing patterns and trends by sifting through large data stored in a repository, using pattern recognition technology and mathematical and statistical techniques (Larose & Larose, 2014). The most primarily accepted definition of data mining is to turn raw data into useful data or information (Velmurugan & Santhanam, 2011). According to Deka (2014) Clustering is a data mining technique used to obtain a set of objects that have the same characteristics with large enough data.

There are various kinds of algorithms for clustering including the algorithm *Partitioning Around Medoids* (PAM)*, Clustering Large Applications* (CLARA) and, *Fuzzy C-Means* (FCM). PAM and CLARA belong to partition methods. The partition method is the simplest and most basic method in cluster analysis, each object is grouped into several special clusters (Han et al., 2012). The PAM algorithm is based on finding k objects that are representative around other objects. This object represents a cluster around other objects in the data set. The representative object is called a medoid (Kaufman & Rousseeuw, 2009). The CLARA algorithm is a development of the PAM algorithm. CLARA has more robust properties against outliers and is used to handle large data sizes. The CLARA algorithm uses a sampling approach, then applies the PAM algorithm to obtain the optimal medoid. Fuzzy C-Means was first introduced by Jim Bezdek in 1981 (Bezdek et al., 1984). According to Yan, the Fuzzy C-means algorithm is a data clustering technique in which the existence of each data point in a cluster is determined by the degree of membership (Yan et al., 1995). FCM is a type of soft clustering where in grouping data, each data can belong to more than one cluster.

Several previous studies have conducted clustering of rice production in Indonesia. Utomo (2018) compared the clustering of rice productivity in Indonesia using the K-Means and Fuzzy C-Means algorithms. Root mean square error was chosen as the method to determine the best method. The RMSE for the K-Means algorithm is 0.978112 and for the Fuzzy C-Means algorithm is 0.98203 so that for rice productivity data in Indonesia in 2015, the K-Means algorithm is better to use (Utomo, 2018). Munthe (2019) applies a clustering time series to classify the value of rice production in Indonesia. Data taken from 1968-2015. Clustering time series analysis using hierarchical and non-hierarchical methods produces the same distribution of cluster members in the three optimal clusters. The first cluster consists of 3 provinces, the second cluster consists of 15 provinces, while the third cluster consists of 8 provinces. The Silhouette Coefficient value of this study is 0.64 or in the Good Classification category. Nasution et al. (2022) grouped rice production in Indonesia during the Covid-19 pandemic using the K-Medoids method. From this study, three clusters were produced, namely high, medium, and low. There are 3 provinces with high-level groups namely West Java, Central Java and East Java. There are 2 provinces with medium level groups, namely South Sumatra and South Sulawesi and 29 other provinces are at low level (Nasution et al., 2022).

Based on this background, research was conducted on the variables of harvested area and rice productivity in every city/district in Indonesia. Clustering is carried out using the partition method, namely *Partitioning Around Medoids* (PAM) and *Clustering Large Applications* (CLARA) as an improvement from PAM method. The method is also used by considering the degree of membership and fuzzy sets as a basis for weighting, namely Fuzzy C-Means (FCM). The Silhouette Coefficient value is used to compare the three methods so that the best method is obtained in clustering harvested area and rice productivity in Indonesia.

## METHOD

The method used in this research is a quantitative method. This study used two variables, namely harvested area and rice productivity. Furthermore, clustering was carried out using the PAM, CLARA, and FCM methods. The following are the steps used in this study.

### A. Data collection

The data used in this study is secondary data for 2021 obtained from the Statistics Indonesia website for each province in Indonesia. The data used is data on harvested area and rice productivity in every city/district in Indonesia which are contained in provincial publications in figures for 2022. Harvested area data is collected every month using a sub-district area approach throughout Indonesia. Paddy productivity data was collected through direct measurements on tiled plots measuring 2.5 m × 2.5 m. Productivity data collection is carried out every four months at farmer's harvest time.

### B. Conducting data pre-processing

1. Cleaning data

    There is a missing value in the data so it is necessary to carry out the imputation process to overcome it. Missing value is filled with its mean value.

2. Data transformation

    In order to have the same measurement scale, both variables are transformed with logarithms with base 10. Data that was originally two to five digits is converted into one digit with two decimal places.

### C. PAM algorithm implementation

1. Determine the number of clusters to be formed
2. Randomly determine the initial medoid
3. Calculating the non-medoid distance to the medoid for each group
4. Places objects based on their closest distance to the medoid. Then, calculate the total distance obtained.
5. Randomly select non-medoid objects in each new medoid group.
6. Calculates the distance of each non-medoid object to the new medoid candidate
    Places objects based on their closest distance to the new medoid candidate. Calculates the total distance obtained with the new medoid candidate.
7. Calculate the difference in the total distance (S) obtained from the subtraction between the total distance in the new medoid candidate and the total distance in the initial medoid. If S < 0, then the candidate for the new medoid is the new medoid. This value indicates that the total distance between each object and the new medoid candidate is less than the total distance between each object and the old medoid, so the new medoid candidate becomes the new medoid.
8. Recalculate Steps 5 to 7 so that no medoid changes occur. The medoid does not change

Imelda Putri Kurniawati, Hasih Pratiwi, Sugiyanto

when the value of S > 0. When the total distance between each object and the new medoid candidate is greater than the total distance of each object with the old medoid, then no medoid exchange occurs. This is because the new medoid candidate is not more centered than the old medoid, so the iteration process will stop.

### D. CLARA algorithm implementation
1. Determine the number of clusters to be formed (c).
2. Dividing data randomly with several subsets of fixed size. The minimum size is (40+2c).
3. Applying the PAM algorithm to each subset in order to obtain the best medoid in each subset.
4. Calculating all the distances of objects that are not medoid to objects that are medoid using the Euclidean Distance formula.

$$d(i,j) = \sqrt{\sum_{i=1}^{n}(x_i - x_j)^2 + (y_i - y_j)^2}$$

where

$d(i,j)$: Euclidean distance between the ith observation and the jth observation

$x_i$ : object in the *i*-th observation of *x*

$x_j$ : object in the *j*-th observation of *x*

$y_i$ : object in the *i*-th observation of *y*

$y_j$: object in the *j*-th observation of *y*

*n*: many observations

5. Places objects based on their closest distance to the medoid.
6. Calculates the total distance, then compares the total distance of all subsets. The subset with the smallest total distance is selected.

### E. FCM algorithm implementation
1. Input the data to be clustered in the form of an X matrix of size n × m (n: the number of data and m: the number of variables for each data).
   Xij : the *i*-th data (i=1,2,.....,n), the *j*-th variable (j = 1,2,.....,m).

   Determine:
   a) The number of clusters to be formed (c), $1 \leq c \leq N$
   b) Rated weighting (w)
   c) Maximum iteration (MaxItr).
   d) Smallest error (ε)
   e) Initial objective function (P0=0)
   f) Initial iteration (t = 1)
2. Generating random numbers $\mu ik$ , i = 1,2,…,n; k = 1,2,…,c ; as the elements of the initial partition matrix U.
   $\mu ik$ is the degree of membership which refers to how likely a data can be a member of a cluster.
3. Calculating the sum of each column (variable):

$$Q_i = \sum_{k=1}^{c} \mu_{ik}$$
$$Q_i = \mu_{i1} + \mu_{i2} + \cdots + \mu_{ic}$$

with $i$ = 1,2,…,n

4. Calculating the $k$-th cluster center ($V_{kj}$) with $k$ = 1,2,…..,$c$, and $j$ = 1,2,…., $m$ as follows:

$$V_{kj} = \frac{\sum_{i=1}^{n}((\mu_{ik})^w X_{ij})}{\sum_{i=1}^{n}(\mu_{ik})^w}$$

where

$P_t$: objective function

$X_{ij}$: element X row i, column j

$V_{kj}$: cluster center

5. Fixed degree of membership of the partition matrix:

$$P_t = \sum_{i=1}^{n}\sum_{k=1}^{c}\left(\left[\sum_{j=1}^{m}(X_{ij} - V_{kj})^2\right](\mu_{ik})^w\right)$$

where

$i$= 1,2,……,n

$k$= 1,2,……,c

$X_{ij}$: the $i$-th data sample, the $j$-th variable

$V_{kj}$: the $k$-th cluster center for the $j$-th variable

$w$: weighting rank

6. Check stop conditions

If ($|P_t$-$P_{(t-1)}$ $|<$ ε) or (t>MaxIter) then stop.

If not, the iteration is increased $t$ = $t$+1, repeat step 4.

## F. Comparing the results of clustering with the three methods based on silhouette coefficient values

Silhouette shows how well objects are located in a cluster compared to being in other clusters (Rousseeuw, 1987). This method measures the validation of the goodness of a data, a single cluster or the entire cluster. Each method is searched for its silhouette coefficient. The method with the greatest silhouette coefficient values is the best method.

1. Calculates the average distance of an object with all other objects that are in one cluster with the equation:

$$a(i) = \frac{1}{|A| - 1}\sum_{j\in A, j\neq i} d(i,j)$$

where

$j$: another object in one cluster A

$d(i, j)$: the distance between objects i and j

2. Calculate the average distance from object i to all objects in other clusters, and take the smallest value using the equation:

$$d(i,C) = \frac{1}{|C|}\sum_{j\in C} d(i,j)$$

where

$d(i, C)$: the average distance of object i to all objects in other cluster C with A≠C, A is the number of members of cluster A, C is the number of members of cluster C

Calculating the average minimum object distance, b(i), which shows the average distance of object i and objects in other clusters is determined using the following equation:

$$b(i) = \min_{C \neq A} d(i,C)$$

3. Calculating Silhouette value:

$$s(i) = \frac{b(i) - a(i)}{\max(a(i), b(i))}$$

4. Calculating *Silhouette Coefficient*:

$$SC = \frac{1}{n}\sum_{i=1}^{n} s(i)$$

where

$n$: many observations

## G. Interpret the results

Clustering results with the best method will be visualized with a map. The software used is ArcGIS. Areas with high rice harvest area and productivity are colored pink while areas classified as low are colored green.

# RESULTS AND DISCUSSION

## A. Partitioning Around Medoids (PAM)

The PAM algorithm uses the partition method to group $n$ objects into k clusters. This algorithm uses objects in a collection of objects to represent a cluster. The object chosen to represent a cluster is called a medoid. Clusters are built by calculating the closeness between medoid and non-medoid objects. The process in this method begins with determining the cluster center and placing objects into the nearest cluster center (Hair et al., 2014). By using the R Studio software, 2 clusters were produced using the PAM method. The medoid in this algorithm is shown in table 1, namely in the City of Tojo Una-Una and the City of Sabang.

**Table 1**
**PAM algorithm medoid**

|  | Harvest Area | Productivity |
|---|---|---|
| Tojo Una-Una | 3.100632 | 1.620968 |
| Sabang City | 4.358406 | 1.676511 |



**Figure 1. PAM Algorithm Cluster Plot**

In the first cluster there are 191 data and are marked in red. In the second cluster there

are 323 data and marked in blue. Silhouette in the first cluster is 0.356471 and in the second cluster is 0.716172. The silhouette in the entire dataset is 0.5825087, meaning that there is a fairly strong bond between the objects and groups that are formed.

## B. Clustering Large Application (CLARA)

CLARA is the development of PAM. This algorithm uses a sampling approach and then applies the PAM algorithm to obtain the optimal medoid. The CLARA algorithm consists of two phases, namely the build and swap phases. In the build phase, representative objects are selected repeatedly with the aim of obtaining the smallest and most similar average distance to the representative objects. The swap phase is carried out to reduce the average distance by replacing representative objects, after k representative objects are selected, each object from the data set is placed to the nearest representative object. The medoid in this algorithm is shown in table 2, namely in the cities of Pakpak Bharat, Aceh Jaya, and South Tangerang City.

**Table 2**
**CLARA algorithm medoid**

|  | Harvest Area | Productivity |
|---|---|---|
| Pakpak Bharat | 3.106976 | 1.595165 |
| Aceh Jaya | 3.973205 | 1.766487 |
| South Tangerang City | 4.358406 | 1.676511 |



**Figure 2. Cluster plot of the CLARA algorithm**

The best Silhouette Coefficient results are obtained with the best sample of 46 and the number of clusters is three. The first cluster of 167 data is marked in red. The second cluster of 139 data is marked in green and the third cluster of 208 data is marked in blue. The resulting Silhouette Coefficient is 0.59. The resulting Silhouette Coefficient value is better than the PAM method.

## C. Fuzzy C-Means (FCM)

This method works with a fuzzy model where data is grouped based on its membership value. This algorithm begins by determining the desired number of clusters and initializing the membership value which contains all the data which will then be grouped based on the cluster. Cluster centers are calculated from the shortest distance to the points that have the greater

membership value. In other words, these membership values will act as temporary weight values in a cluster. The output of Fuzzy C-Means itself is a row of cluster centers and also several degrees of membership for each of these data points (Bezdek, 1980). This method works with a Fuzzy model that allows all data from all group members to be formed with different degrees of membership between 0 and 1 (Bora et al., 2014; Sanmorino, 2012). The following are the degrees of membership in the top and bottom five rows in the data.

**Table 3**
**Degree of membership of the FCM algorithm**

|  | Clusters 1 | Clusters 2 |
|---|---|---|
| Simeulue | 0.76736900 | 0.232630997 |
| Aceh Singkil | 0.04270866 | 0.957291341 |
| South Aceh | 0.94790789 | 0.052092107 |
| Southeast Aceh | 0.95229164 | 0.047708358 |
| East Aceh | 0.99361697 | 0.006383033 |
| ..... | ..... | ..... |
| Peak | 0.9996126 | 0.000387404 |
| Dogiyai | 0.9996126 | 0.000387404 |
| Jaya Intan | 0.9996126 | 0.000387404 |
| Deiyai | 0.9996126 | 0.000387404 |
| Jayapura City | 0.0541255 | 0.945874495 |

The silhouette value obtained by this method is 0.828 and the number of clusters is two. There are 369 data in the first cluster (red) and 145 data in the second cluster (blue). The resulting plot is as follows.



**Figure 3. Cluster plot of the FCM algorithm**

## D. Visualization of clustering results

The results used for visualization are the results of the Fuzzy C-Means method because this method has the largest silhouette coefficient value compared to the other two methods. This visualization uses ArcGis software. Areas included in the first cluster are colored red, namely areas with high rice productivity. Areas included in the second cluster are colored green, namely areas with low rice productivity.
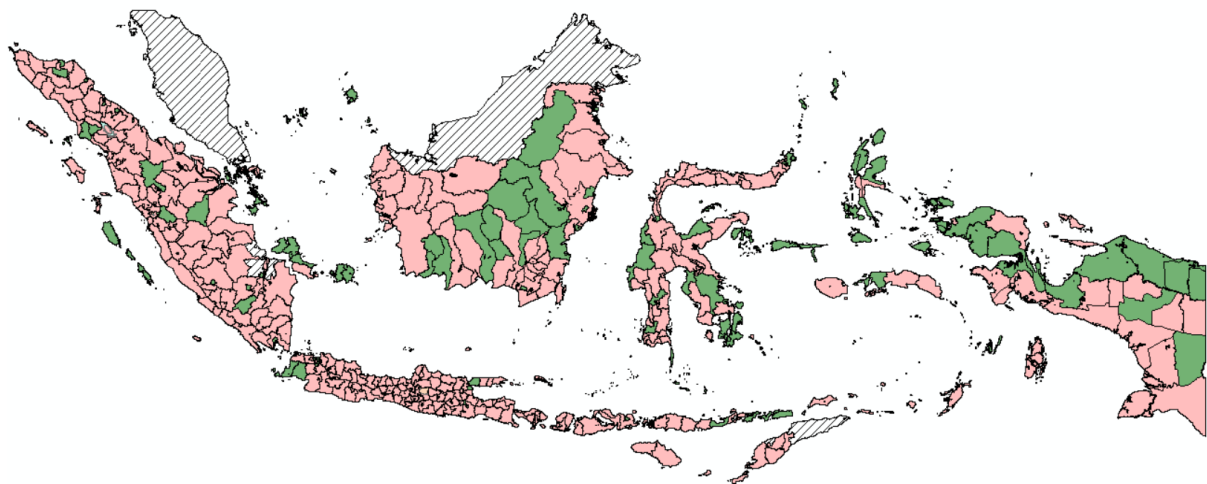
**Figure 4. Visualization of the map of Indonesia based on clustering results using the FCM method**

In Figure 4 it can be seen that on the island of Papua which has a low level of productivity and harvested area, namely Wondama Bay, Bintuni Bay, South Sorong, Sorong, Raja Ampat, Manokwari, Jayawijaya, Jayapura, Nabire, Boven Digoel, Sarmi, Keerom, Waropen and Jayapura city. In addition, districts/cities in Papua Island have high productivity. On the island of Sulawesi which has a low level of productivity and harvested area, namely Sangihe Islands, Talaud Islands, North Minahasa, Southeast Minahasa, Bolaang Mongondow, East Bolaang Mongondow, Bitung City, Tomohon City, Banggai Islands, Parigi Moutong, Tojo Una-Una, Kota pAlu, Selayar Islands, Makassar City, Pare-Pare City, Buton, Muna, North Kolaka, North Buton, North Konawe, Konawe Islands, West Muna, South Buton, Kendari City, Bau-Bau City, Gorontalo City, Majene, and North Mamuju. On the island of Kalimantan which has low productivity and rice harvest area, namely Pontianak City, West Waringin City, South Barito, North Barito, Sukamara, Lamandau, Katingan, Gunung Mas, Murung Raya, Palangkaraya City, Balangan, Banjarmasin City, Banjarbaru City, Kutai West, Mahakam Ulu, Balikpapan City, Samarinda City, Bontang City, Malinau, Tana Tidung, and Tarakan City.

All cities/regencies on the island of Bali have a high level of productivity and rice harvest area. In West Nusa Tenggara Province, which has a low level of productivity and rice harvest area, namely the City of Mataram and the City of Bima. The provinces of East Nusa Tenggara which have low productivity levels and rice harvest areas are Alor, Lembata, Sikka, Sabu Raijua, and Kupang City. Regencies/cities on Java Island that have low levels of productivity and rice harvest area are East Jakarta, West Jakarta, North Jakarta, Bogor City, Sukabumi City, Bandung City, Cirebon City, Bekasi City, Depok City, Cimahi City, Magelang City, Surakarta City, Salatiga City, Pekalongan City, Tegal City, Yogyakarta City, Blitar City, Malang, Probolinggo, Pasuruan, Mojokerto, Madiun, Suarabaya, Batu, Pandeglang, Lebak, Serang City, and South Tangerang City. Regencies/cities on the island of Sumatra that have low levels of productivity and rice harvest area are Aceh Singkil, Bener Meriah, Banda Aceh, Lnagsa City, Lhokseumawe, Subulussalam, Pakpak Bharat, South Labuhanbatu, Tanjung Balai, Pematangsiantar, Tebing Tinggi, Medan, Binjai , Gunungsitoli, Mentawai Islands, Solok City, Sawahlunto, Padang Panjang, Bukittinggi, Pariaman, Indragiri Hulu, Rokan Hulu, Bengkalis, Dumai, Jambi, Ogan Komering Ulu, Palembang City, Prabumulih City, Lubuklinggau, Central

Bengkulu, Bengkulu City, Bandar Lampung City, Bangka, Belitung, West Bangka, Central Bangka and East Belitung. Regencies/ cities in the Riau Archipelago have high levels of productivity and rice harvest area except Batam City and Tanjung Pinang City.

## CONCLUSION

The PAM algorithm produces two clusters with a silhouette coefficient value of 0.58. The CLARA algorithm with 100 samples produces three clusters with a silhouette coefficient of 0.59. The FCM algorithm produces two clusters with a silhouette coefficient value of 0.83. The best of the three methods is Fuzzy C-Means (FCM) and is used as a reference in making maps of Indonesian districts/cities. There are 368 regencies/cities in Indonesia that are included in cluster 1, namely areas with high rice harvest area and productivity. There are 146 regencies/cities in Indonesia which are included in cluster two, namely areas with low rice harvest area and productivity.

Based on the results that have been tested in this study, there are several suggestions that can be made for further research, namely using the latest data after the pandemic so that the results are more accurate. In addition, you can use other methods so you can find a better method.

## REFERENCES

Bezdek, J. C. (1980). A convergence theorem for the fuzzy ISODATA clustering algorithms. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, *1*, 1–8. Google Scholar

Bezdek, J. C., Ehrlich, R., & Full, W. (1984). FCM: The fuzzy c-means clustering algorithm. *Computers & Geosciences*, *10*(2–3), 191–203. Elsevier

Bora, D. J., Gupta, D., & Kumar, A. (2014). A comparative study between fuzzy clustering algorithm and hard clustering algorithm. *ArXiv Preprint ArXiv:1404.6059*. Google Scholar

Hair, J. F., Black, W. C., Babin, B. J., & Anderson, R. E. (2014). Multivariate data analysis (MVDA). *Pharmaceutical Quality by Design: A Practical Approach. Https://Doi. Org/10.1002/9781118895238. Ch8*. Google Scholar

Han, J., Pei, J., & Tong, H. (2022). *Data mining: concepts and techniques*. Morgan kaufmann. Google Scholar

Kaufman, L., & Rousseeuw, P. J. (2009). *Finding groups in data: an introduction to cluster analysis*. John Wiley & Sons. Google Scholar

Larose, D. T., & Larose, C. D. (2014). *Discovering knowledge in data: an introduction to data mining* (Vol. 4). John Wiley & Sons. Google Scholar

Munthe, A. D. (2019). Penerapan clustering time series untuk menggerombolkan provinsi di Indonesia berdasarkan nilai produksi padi. *Jurnal Litbang Sukowati: Media Penelitian Dan Pengembangan*, *2*(2), 11. Google Scholar

Nasution, D., Solikhun, S., & Nasution, D. (2022). Penerapan K-Medoids Dalam Mengelompokkan Produksi Padi Di Indonesia Pada Masa Pandemi Covid-19. *Prosiding Seminar Nasional Riset Information Science (SENARIS)*, *4*(2), 26–35. Google Scholar

Rousseeuw, P. J. (1987). Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. *Journal of Computational and Applied Mathematics*, *20*, 53–65. Elsevier

Sanmorino, A. (2012). Clustering batik images using fuzzy C-means algorithm based on log-average luminance. *Computer Engineering and Applications Journal*, *1*(1), 25–31. Google Scholar

Statistic Indonesia. (2015). Impor Beras menurut Negara Asal Utama. *BPS. Jakarta*.

Tan, P.-N., Steinbach, M., & Kumar, V. (2006). Introduction to data mining. In *Inc., New Delhi*. Pearson Education. Google Scholar

Turban, E., J. E., Aronson, Liang, T., Asoke, K., & Ghosh. (2007). *Decision Support System and Intelligent System*.

Utomo, R. A. (2018). *Perbandingan Clustering Produktivitas Padi Di Indonesia Menggunakan algoritme K-Means dan Fuzzy C-Means*. Google Scholar

Velmurugan, T., & Santhanam, T. (2011). A survey of partition based clustering algorithms in data mining: An experimental approach. *Information Technology Journal*, *10*(3), 478–484. Google Scholar

Wahyudi, K. D. (2018). Kebijakan strategis usaha pertanian dalam rangka peningkatan produksi dan pengentasan kemiskinan. *Majalah Ilmiah Dian Ilmu*, *11*(2). Google Scholar

Yan, J., Ryan, M., & Power, J. (1995). *Using fuzzy logic: Towards intelligent systems*. Prentice-Hall, Inc. Google Scholar

Yusuf, M., Haeruddin, H., & Kusmiah, N. (2020). Pengaruh Faktor Sosial Ekonomi terhadap Pendapatan Usahatani Padi Sawah (Oryza Sativa). *Journal Peqguruang*, *2*(1), 349–352. Google Scholar