# Implementation of Fuzzy C-Means for Clustering the Majelis Ulama Indonesia (MUI) Fatwa Documents

**Fajar Rohman Hariri[1]**
[1]Department of Informatics, UIN Maulana Malik Ibrahim Malang

## Article Info

## ABSTRACT

Since the Indonesian Ulema Council (MUI) was established in 1975 until now, this institution has produced 201 edicts covering various fields. Text mining is one of the techniques used to collect data hidden from data that form text. One method of extracting text is Clustering. The present study implements the Fuzzy C-Means Clustering method in MUI fatwa documents to classify existing fatwas based on the similarity of the issues discussed. Silhouette Coefficient is used to analyze the resulting clusters, with the best value of 0.0982 with 10 clusters grouping. Classify fatwas based on the similarity of the issues discussed can make it easier and faster in the search for an Islamic law in Indonesia.

*Corresponding Author:*
Fajar Rohman Hariri,
Department of Informatics,
UIN Maulana Malik Ibrahim Malang,
Jl. Gajayana No.50, Dinoyo – Lowokwaru - Malang,
fajar@ti.uin-malang.ac.id

## 1. INTRODUCTION (10 PT)

Fatwa is a product of Islamic Law that has existed since the era of the Prophet SAW, that later became a product of Islamic Law that is developed until now. Fatwas of Islamic scholars who are gathered in the books of fiqh and the decisions of fatwa institutions are part of the casuistic ijtihad because it is a response or answer to questions raised by the fatwa applicant. Therefore, fatwa is one of the solutions in solving problems that occur in modern times.

Majelis Ulama Indonesia (MUI), which is a forum for discussion of Muslim clerics, zu'ama, and intellectuals and is an umbrella for all Indonesian Muslims, is the most competent institution in answering and solving every social and religious problem that always arises and is faced by the community.

Since Majelis Ulama Indonesia (MUI) was established in 1975 until now, this institution has produced 201 fatwas covering various fields, such as matters of worship, ahwal al-syakhshiyah, family planning, food and beverage issues, culture, interfaith relations, and others. Information about the MUI's fatwas can be accessed on the page https://mui.or.id/fatwa/. But as shown in Figure 1, there is no grouping of fatwas on that page.



Figure 1. Webpage https://mui.or.id/fatwa

Text data is a good example of unstructured information, which is one of the simplest forms of data that can be generated in most scenarios. Unstructured text is easily processed and perceived by humans, but is significantly harder for machines to understand. Needless to say, this volume of text is an invaluable source of information and knowledge. As a result, there is a desperate need to design methods and algorithms in order to effectively process this avalanche of text in a wide variety of applications[1].

However, this is not accompanied by knowledge that can extract the information needed from these electronic documents. Therefore a method is needed to make the classification of these documents ease and simplify [2]. One of the methods that can be used is text mining [3], [4].

Text mining method is a development of data mining that can be applied to overcome this problem[3]. The algorithms in text mining are made to be able to recognize semi-structured data such as synopsis, abstracts, and the content of documents. Text mining is used for digging hidden data from text based data. One of the techniques of text mining is clustering technique. Clustering technique is a classification technique that is commonly used in data mining [5]. Clustering can be defined as the technique of grouping a collection of data objects as clusters or classes in such a way that objects within a cluster are similar to one another, but dissimilar to the other clusters' objects [6], [7].

Text mining has been widely used for document grouping. Cepy Slamet classifies the letters in the Qur'an based on verses using the k-means clustering algorithm [8]. Deka uses the k-means clustering technique to classify book titles according to their categories so that it makes it easier for librarians to group book placement and design strategies in increasing student reading interest at the Islamic University of Indonesia [9]. Eko Yulian (2018) implements Text Mining on LGBT Themes in the Tweet Archives of the Bandung City Community [10].

The present research is about the implementation of Fuzzy C-Means Clustering method towards MUI's fatwa documents to classify the fatwas according to the similarities of the issues. Nazief & Andriani algorithms are used for the documents stemming [11]. The main difference between our research and other document fuzzy clustering research is we used MUI's fatwa documents and will use several scenarios of the number of clusters in grouping MUI's fatwa documents. The results of the clusters will be evaluated using a silhouette coefficient so that the resulting clusters can be seen.

## 2. METHOD
## 2.1. Type of Research

This research is an experimental study. To get the best performance of the Fuzzy C-Means algorithm in clustering, several experiments were conducted with different parameters. Data obtained from MUI fatwa documents obtained from page https://mui.or.id/fatwa/ with pdf extension.

## 2.2. Metode Analisis Data

From the fatwa documents obtained, an experiment will be conducted with several test scenarios described in table 1 below

Table 1. Testing Scenario

| Dimensional Reduction ( % ) | Cluster Value |
|---|---|
| 0 ( All ) | 2 , 3, 4 , 5, 6 , 7 , 8, 9, 10 |
| 10 | 2 , 3, 4 , 5, 6 , 7 , 8, 9, 10 |
| 20 | 2 , 3, 4 , 5, 6 , 7 , 8, 9, 10 |
| 30 | 2 , 3, 4 , 5, 6 , 7 , 8, 9, 10 |
| 40 | 2 , 3, 4 , 5, 6 , 7 , 8, 9, 10 |
| 50 | 2 , 3, 4 , 5, 6 , 7 , 8, 9, 10 |
| 60 | 2 , 3, 4 , 5, 6 , 7 , 8, 9, 10 |
| 70 | 2 , 3, 4 , 5, 6 , 7 , 8, 9, 10 |
| 80 | 2 , 3, 4 , 5, 6 , 7 , 8, 9, 10 |
| 90 | 2 , 3, 4 , 5, 6 , 7 , 8, 9, 10 |

Later will be compared with the cluster that was produced by using The Silhouette Coefficient.

Generally speaking, the research will be conducted based on the flow shown at figure 2:
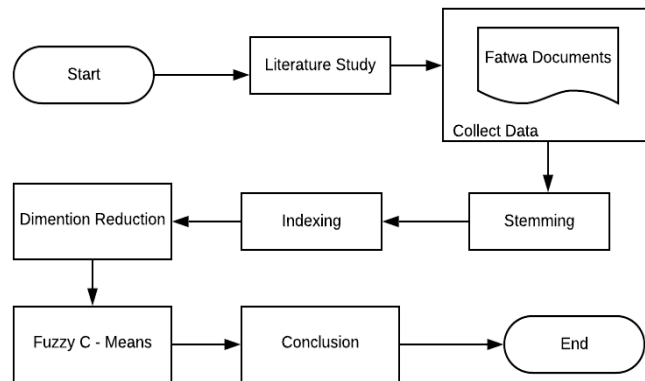


Figure 2. Research flow

These are the flow of the research:
1) Literature Study

At this stage the information collected is needed for the research. The information needed is obtained by studying the literature on document clustering problems, Fuzzy C-Means method, and stemming algorithms for Indonesian.
2) Data Collection

Looking for data related to the system, among others collecting data on MUI fatwa documents.
3) Stemming

Breaking down the words and taking the basic words from the abstract of the documents using Nazief & Andriani Algorithm.
4) Indexing

Each of the basic words is listed and the frequency of the words are observed. Then, the matrix term frequent is created.
5) Dimensional Reduction

Abstract documents that are available will produce many basic words, and there are some words that will only appear in one or a few documents, so the next step is to reduce the dimensions by removing the words.
6) Fuzzy C-Means Method Calculation

From the data after dimension reduction, it will be processed using the Fuzzy C-Means method to produce data clusters.
7) Drawing Conclusion

After observing several experimental scenarios, see how the performance of the K Means method, and see how the best conditions produce the best accuracy.

## 2.3. Fuzzy C-Means Clustering

Clustering is a technique that can divide a set of data into groups, where each group has a degree of equality. The clustering technique makes it possible to do with two approaches namely hard grouping and soft grouping. The soft grouping one is Fuzzy C Means with a degree of membership between zero and one.

In 1973, Dunn discovered a method [12]. Later on, this method was developed by Bezdek in 1981 [13]. This method centers on fuzzy logic that is similar to the K-Means method and the naming of this method is Fuzzy C Means Clustering (FCM).

FCM is a data clustering technique in which the existence of each data point in a cluster is determined by the degree of membership. FCM uses a fuzzy grouping model with a blurring index using Euclidean Distance so that data can be a member of all classes or clusters formed with different membership degrees from 0 to 1 [14].

The basic concept of FCM, first is to determine the cluster center, which will mark the average location for each cluster. In the initial condition, the cluster center is still inaccurate. Each data point has a degree of membership for each cluster that is formed. By improving the cluster center and the degree of membership of each data point repeatedly, it will be seen that the cluster center will shift to the right location. This repetition is based on the minimization of objective functions that describe the distance from

the data point given to the center of the cluster weighted by the degree of membership of that data point [15].

The work process of this method is that the data collected has been entered into the data group, the data entry into groups based on the membership value [16]. The FCM algorithm is as follows [17]:

1) Determination of the number of groups (c), maximum iteration (MaxIter), fuzzifier (m), then determining the smallest objective value, the expected objective value ($\varepsilon$). Determination of initial objective function (P0 = 0) and initial iteration (t = 1)

2) Increasing the random number $u_{ik}$, a lot of data is symbolized by i, then the number of groups is symbolized by k. The elements of membership of U are i and k.

3) Calculating the center of the to-i group with the equation:

$$P_i = \frac{\sum_{k=1}^{N} (u_{ik})^m X_k}{\sum_{k=1}^{N} (u_{ik})^m} \qquad (1)$$

4) The formula for calculating object functions on the t iteration with an equation:

$$J(P,U,X,c,m) = \sum_{i=1}^{C} \sum_{k=1}^{N} (u_{ik})^m d_{ik}^2 (X_k, P_1) \quad (2)$$

Information:

- c is the expected number of groups,

- N is a lot of research objects,

- $u_{ik}$ is the membership value of the k-specific object in the i group (part of the matrix)

- U, m are fuzzifiers, and $d_{ik}2$ $(x_k, p_i)$ is the distance between the k approach and the center of the k-I group.

5) Calculation of changes in the membership matrix with the equation:

$$u_{ik} = \frac{1}{\sum_{j=1}^{c} \left[\frac{d_{ik}^2}{d_{jk}^2}\right]^{\frac{1}{m-1}}} \qquad (3)$$

the $u_{ik}$ is the membership value of the k-th object with the center of the group, i, $d_{ik}$, $d_{ik}2$ is the distance between the k-th object and the center of the i-th group, $d_{ik}2$ is the distance between the k-object with the center of the i group j, and m is the fuzzifier.

6) Condition checking

- If |Jt−Jt−1|<$\varepsilon$ or $t>MaxIter$ then stop;

- If not :t = t + 1, repeat to step 3;

The Fuzzy C-Means algorithm is often used for grouping data that is used to estimate something, estimation with FCM is an efficient estimate and does not require many parameters [18] . Several studies have produced statements that the Fuzzy CMeans method can be used to classify data based on certain attributes [19].

## 2.4. Dimensional Reduction

Dimensional Reduction Process is done to shrink the date so it can reduce the time for computation. In Dimensional Reduction, the researcher must pay attention to the characteristics of the data because the eliminated dimension can also be the eliminated characteristics of the data [20]. Because of that, in this research dimensional reduction of the spare part is chosen. The most frequent term that mostly has the value of zero will be eliminated, because those terms represent rarely used terms so it can be ignored. The

number of the dimensional reduction is according to the parameters. To avoid the elimination of characteristics of the data, value parameters are made to limit the dimensional reduction. Therefore, even though the term value is mostly zero, the term will not be eliminated if one of the frequencies is more than the parameters.

TF data that is included in the database has different values. Terms that have mostly zero TF can be ignored, as long as no TF value that is more than the parameters. In short, the process of dimensional reduction is a process of eliminating the feature or term that the all or almost all the values of the frequencies are zero. If some of the values is not zero, then those TF have to be more than certain parameters. This is needed so the data is not losing the characteristics, so the process can run well [21].
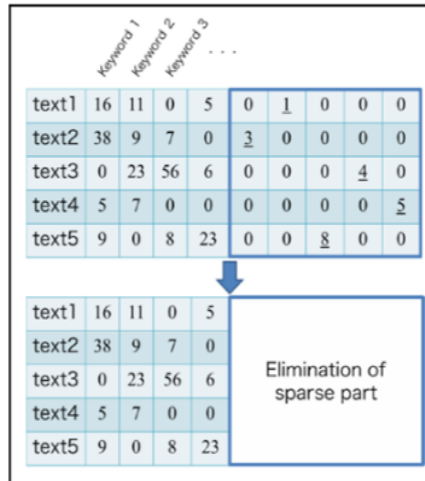


Figure 3. Dimensional Reduction

In figure 3, shown that the parameter is 40%. Then, the dimensional reduction is done by eliminating the columns that have the value more than 40% from the data. The 5th until 9th column will be eliminated because those don't qualify the criteria.

## 2.5. Silhouette Coefficient

Silhouette Coefficient is one of the methods used to test the quality and strength of a cluster. The silhouette coefficient method is a combination of the cohesion method and the separation method. The cohesion method itself is a method used to measure how close relationships are between objects in the same cluster. While the separation method is used to measure how far a cluster is separated from other clusters.

Silhouette has three stages in its calculation, the following are silhouette coefficient calculation stage [22]:

a. Calculate the average distance of objects with all documents in one cluster by using equation ( 4 )

$$a(i) = \frac{1}{[A]-1}\sum \quad j \in A, j \neq 1 \ d(i,j) \ .......... (4)$$

b. Then calculate the distance of the object with all documents between clusters by using equation (5)

$$d(i,C) = \frac{1}{[A]}\sum \quad j \in C \ d(i,j) .......... (5)$$

c. Then calculate the silhouette value by using equation(6)

$$s(i) = \frac{b(i)-a(i)}{max \ (a(i),b(i))} .......... (6)$$

## 3. RESULTS AND DISCUSSION
### 3.1 Preprocessing
Preprocessing stage:
a. Tokenization

For the tokenization process, in reading the contents of pdf files in this study using the itextsharp library. This library is quite good in changing the contents of pdf files into text.
b. Stopword and Stoplist Removal

Terms that are generated in the previous process are various, among the existing terms, there are some terms that are considered useless or worthless, such words are for example conjunctions, prepositions, etc. There are a total of 694 words included in the word list. Which can then be removed to do the next process.

c. Stemming

The stemming method used is the Nazief & Adriani Algorithm. From the fatwa document data, after the tokenization process, stopword removal, then for each word, the basic words are taken with the algorithm by applying regular expressions.

d. Indexing

The next process after preprocessing is indexing. From the 114 existing MUI fatwas documents, they were preprocessed, and changed them to the most frequent matrix terms, resulting in 1531 words (terms).

e. Dimensional Reduction

From 114 Fatwa document data after processing it produces 1531 terms (words). Of the 1531 words available, there are some words that only appear in one or two documents, so that the word can be ignored. Therefore, to eliminate these words, and also to speed up the calculation process, a dimensional reduction process is carried out.

The results of the number of terms generated after the process of dimension reduction are in the following table 2:

Table 2. Dimensional Reduction

| Dimensional Reduction Value | Number of Terms |
|---|---|
| All | 1531 |
| 10 % | 142 |
| 20 % | 85 |
| 30 % | 53 |
| 40 % | 36 |
| 50 % | 26 |
| 60 % | 17 |
| 70 % | 13 |
| 80 % | 12 |
| 90 % | 5 |

The percentage above means that the words taken are words that appear in more than how many % of documents. The fewer parameters, the fewer terms will be produced.

3.2 Experiment Result

After the dimensional reduction process, the clustering process is carried out using K-Means, with different dimension reduction scenarios and different number of clusters, resulting in the value of *Silhoutte Coefficient* as shown in table 3 below:

Tabel 3. Silhouette Coefficient

| Dimensional Reduction Value | Number of Clusters | Silhouette Coefficient |
|---|---|---|
| All 0% | 1 | 0,1636 |
| | 2 | 0,1215 |
| | 3 | 0,1296 |
| | 4 | 0,1199 |
| | 5 | 0,1237 |
| | 6 | 0,1232 |
| | 7 | 0,1138 |
| | 8 | 0,1193 |
| | 9 | 0,0982 |
| | 10 | 0,1636 |
| 10 % | 1 | 0,2131 |
| | 2 | 0,1800 |
| | 3 | 0,1616 |
| | 4 | 0,1480 |
| | 5 | 0,1369 |
| | 6 | 0,1187 |
| | 7 | 0,1394 |
| | 8 | 0,1357 |
| | 9 | 0,1297 |
| | 10 | 0,2131 |

| Dimensional Reduction Value | Number of Clusters | Silhouette Coefficient |
|---|---|---|
| 20 % | 1 | 0,2717 |
| | 2 | 0,2084 |
| | 3 | 0,1943 |
| | 4 | 0,1736 |
| | 5 | 0,1330 |
| | 6 | 0,1303 |
| | 7 | 0,1402 |
| | 8 | 0,1379 |
| | 9 | 0,1339 |
| | 10 | 0,2717 |
| 30 % | 1 | 0,3247 |
| | 2 | 0,2520 |
| | 3 | 0,2098 |
| | 4 | 0,1620 |
| | 5 | 0,1479 |
| | 6 | 0,1540 |
| | 7 | 0,1651 |
| | 8 | 0,1502 |
| | 9 | 0,1534 |
| | 10 | 0,3247 |
| 40 % | 1 | 0,3283 |
| | 2 | 0,2502 |
| | 3 | 0,1888 |
| | 4 | 0,1916 |
| | 5 | 0,1398 |
| | 6 | 0,1737 |
| | 7 | 0,1468 |
| | 8 | 0,1607 |
| | 9 | 0,1508 |
| | 10 | 0,3283 |
| 50 % | 1 | 0,3247 |
| | 2 | 0,2520 |
| | 3 | 0,2098 |
| | 4 | 0,1620 |
| | 5 | 0,1479 |
| | 6 | 0,1540 |
| | 7 | 0,1651 |
| | 8 | 0,1502 |
| | 9 | 0,1534 |
| | 10 | 0,3247 |
| 60 % | 1 | 0,1799 |
| | 2 | 0,2293 |
| | 3 | 0,2018 |
| | 4 | 0,2052 |
| | 5 | 0,2134 |
| | 6 | 0,1895 |
| | 7 | 0,2099 |
| | 8 | 0,1909 |
| | 9 | 0,2029 |
| | 10 | 0,1799 |
| 70 % | 1 | 0,2288 |
| | 2 | 0,2769 |
| | 3 | 0,2452 |
| | 4 | 0,2408 |
| | 5 | 0,2320 |
| | 6 | 0,2432 |
| | 7 | 0,2268 |
| | 8 | 0,2291 |
| | 9 | 0,2099 |
| | 10 | 0,2288 |
| 80 % | 1 | 0,2037 |
| | 2 | 0,2537 |
| | 3 | 0,2410 |
| | 4 | 0,2564 |
| | 5 | 0,2573 |
| | 6 | 0,2572 |
| | 7 | 0,2629 |
| | 8 | 0,2543 |
| | 9 | 0,2437 |
| | 10 | 0,2037 |
| | 1 | 0,3879 |

| Dimensional Reduction Value | Number of Clusters | Silhouette Coefficient |
|---|---|---|
| 90 % | 2 | 0,3804 |
|  | 3 | 0,3673 |
|  | 4 | 0,3670 |
|  | 5 | 0,3795 |
|  | 6 | 0,3646 |
|  | 7 | 0,3875 |
|  | 8 | 0,3985 |
|  | 9 | 0,4157 |
|  | 10 | 0,3879 |

The results of the reduction of dimensions with the Silhouette Coefficient are shown in Figure 4 below:
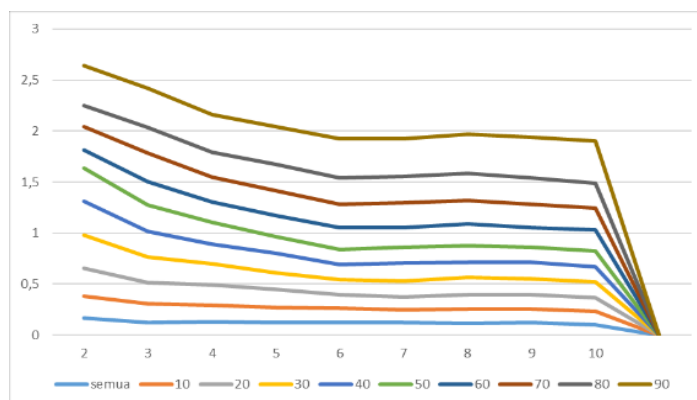


Figure 4. Silhouette Coefficient

Highest and lowest Silhouette Coefficient value from each Dimensional Reduction in table 4 below:

Tabel 4. Highest and Lowest Silhouette Coefficient

| Dimensional Reduction Value | Highest | Number of Clusters | Lowest | Number of Clusters |
|---|---|---|---|---|
| All | 0,1636 | 1 | 0,0982 | 9 |
| 10 % | 0,2131 | 1 | 0,1187 | 6 |
| 20 % | 0,2717 | 1 | 0,1303 | 6 |
| 30 % | 0,3247 | 1 | 0,1479 | 5 |
| 40 % | 0,3283 | 1 | 0,1398 | 5 |
| 50 % | 0,3247 | 1 | 0,1479 | 5 |
| 60 % | 0,2293 | 2 | 0,1799 | 10 |
| 70 % | 0,2769 | 2 | 0,2099 | 9 |
| 80 % | 0,2629 | 7 | 0,2037 | 10 |
| 90 % | 0,4157 | 9 | 0,3646 | 6 |

From the results above it is known that the reduction of dimensions and number of clusters affect the quality of the clusters produce. As shown at table 3,table 4 and figure 4, The highest value of Silhouette Coefficient 0,4157 is when 90% Dimensional Reduction and Number of clusters is 9. On the other hand, the lowest value 0,0982 is when without using Dimensional Reduction and the number of clusters is 9. Silhouette Coefficient show how close are between objects in the same cluster, so lower is better.

## 4.   CONCLUSION (10 PT)

From the results of the study, it is found that Dimensional Reduction and the number of the clusters affect the quality of the resulting cluster. The Fuzzy C-Means Clustering method is quite good in classifying MUI's Fatwa documents with the lowest value of Silhouette Coefficient 0,0982 when without dimensional reduction and the number of clusters is 9. The word in the fatwa document is very specific, because it contains quite different details of Islamic law, so that when the dimensions are reduced it results in a higher silhoutte coefficient. Future research is expected to use of several other methods or Algorithm's.

## 5. REFERENCES (10 PT)

[1] M. Allahyari, E. D. Trippe, and J. B. Gutierrez, "A Brief Survey of Text Mining : Classification , Clustering and Extraction Techniques," 2017.

[2] W. B. Zulfikar, M. Irfan, C. N. Alam, and M. Indra, "The comparation of text mining with Naive Bayes classifier, nearest neighbor, and decision tree to detect Indonesian swear words on Twitter," *5th Int. Conf. Cyber IT Serv. Manag.*, vol. pp. 1-5, d, 2017.

[3] D. S. Maylawati and G. A. P. Saptawati, "Set of Frequent Word Item sets as Feature Representation for Text with Indonesian Slang," *Int. Conf. Comput. Appl. Informatics 2016*, 2017, doi: 10.1088/1742-6596/755/1/011001.

[4] M. I. N. Saputra, D. Fauzy, R. A. Hakim, P. Dauni, M. D. Firdaus, and I. Taufik, "Implementation of Fuzzy C-Means algorithm to classifying research topics in informatics department, UIN Sunan Gunung," *J. Phys. Conf. Ser.*, vol. 1402, no. 2, 2019, doi: 10.1088/1742-6596/1402/2/022091.

[5] B. Santosa, *Data mining teknik pemanfaatan data untuk keperluan bisnis*. Yogyakarta: Graha Ilmu, 2007.

[6] and N. I. G. O. M. E. Zein, L. M. E. Bakrawy, "A robust 3D mesh watermarking algorithm utilizing fuzzy CMeans clustering," *Futur. Comput. Informatics*, vol. 2, pp. 148–156, 2017.

[7] N. P. and M. P. J. Iran, "Clustering Techniques and the Similarity used in Clustering: A survey," *Int. J. Comput. Appl.*, vol. 134, no. 7, 2016.

[8] C. Slamet, "Clustering the Verses of the Holy Qur'an using K-Means Algorithm," *Asian J. Inf. Technol.*, vol. 15, no. 24, 2016.

[9] D. D. C. Nugraha, "Klasterisasi Judul Buku dengan Menggunakan Metode K – Means," *Semin. Nas. Apl. Teknol. Inf.*, 2014.

[10] E. Yulian, "Text Mining dengan K-Means Clustering pada Tema LGBT dalam Arsip Tweet Masyarakat Kota Bandung," *J. Mat. "MANTIK,"* vol. 4, no. 1, 2018.

[11] L. Agusta, U. Kristen, and S. Wacana, "Perbandingan Algoritma Stemming Porter Dengan Algoritma Nazief & Adriani Untuk Stemming Dokumen Teks Bahasa Indonesia," pp. 196–201, 2009.

[12] J. C. Dun, "A fuzzy relative of the ISODATA process and its use in detecting compact well-separated clusters," *J. Cybernet*, vol. 3, no. 3, pp. 32–57, 1973.

[13] J. C. Bezdek, *Pattern recognition with fuzzy objective function algorithms*. USA: Kluwer Academic Publishers Norwell, 1981.

[14] E. T. Luthfi, "FUZZY C-MEANS UNTUK CLUSTERING DATA ( STUDI KASUS : DATA PERFORMANCE MENGAJAR DOSEN )," vol. 2007, no. November, pp. 1–7, 2007.

[15] S. Kusumadewi and H. Purnomo, *Aplikasi Logika Fuzzy untuk pendukung keputusan*. Yogyakarta: Graha Ilmu, 2010.

[16] M. Zhang, L., & Luo, *iverse fuzzy c-means for image clustering. Pattern Recognition Letters*, 1st ed. 2018.

[17] L. Zhang, M. Luo, J. Liu, Z. Li, and Q. Zheng, "Diverse Fuzzy c -Means for Image Clustering," *Pattern Recognit. Lett.*, 2018.

[18] J. Stetco, A., Zeng, X., & Keane, "Expert Systems with Applications Fuzzy C-means ++ : Fuzzy C-means with effective seeding initialization," *Expert Syst. Appl.*, pp. 7541–7548, 2015.

[19] Muhardi and Nisar, "PENENTUAN PENERIMA BEASISWA DENGAN ALGORITMA FUZZY C-MEANS DI UNIVERSITAS MEGOW PAK TULANG BAWANG," *J. TIM Darmajaya*, vol. 01, no. 02, pp. 158–174, 2015.

[20] J. O. Yang, Y., Pedersen, "Comparative Study on Feature Selection in Text Categorization," *Proc. Fourteenth Int. Conf. Mach. Learn.*, 1997.

[21] H. Toyota, T., Nobuhara, "Visualization of the Internet News Based on Efficient Self-Organizing Map Using Restricted Region Search and Dimensionality Reduction," *J. Adv. Comput. Intell. Intell. Informatics*, vol. 6, 2012.

[22] M. Jujjuri, Ramadevi & Rao, "Evaluation of enhanced subspace clustering validity using silhouette coefficient internal measure," *J. Adv. Res. Dyn. Control Syst.*, pp. 321–328, 2019.