# Cluster analysis of the national examination: School grouping to maintain the sustainability of high school quality

**Raoda Ismail[1]; Heri Retnawati[2]; Okky Riswandha Imawan[1]\***

[1]Universitas Cendrawasih, Indonesia
[2]Universitas Negeri Yogyakarta, Indonesia
\*Corresponding Author. E-mail: okkyriswandha.2021@student.uny.ac.id

## ARTICLE INFO

## ABSTRACT

This study aims to classify high schools in Papua Province, Indonesia, based on the 2019 National Examination scores so they can be considered in maintaining the sustainability of school quality in Papua. In this study, all senior high schools in Papua Province were grouped into three clusters: Cluster 1 (high), Cluster 2 (medium), and Cluster 3 (low clusters) using the K-means algorithm on the 2019 National Examination data. The data were obtained through the website official Center of Educational Assessment of the Ministry of Education, Culture, Research, and Technology of the Republic of Indonesia. Clarification was done by grouping data on national examination scores from each school based on the similarity of the data with data from other schools. The results of the high school clustering using the K-means algorithm show that 18 schools are in Cluster 1, 58 schools in Cluster 2, and 68 schools in Cluster 3. The results of the analysis of the K-means algorithm show an R2 value of 0.723 and a Silhouette score of 0.42.

## INTRODUCTION

The National Examination (*Ujian Nasional* or UN) has been carried out for approximately 14 years since it was first introduced in 2005 to replace the National Final Examination. National assessment, prior to the National Examination, is carried out annually by the Education Assessment Center at the elementary, junior high, high, vocational, and equivalent schools. UN aims to measure and assess the achievement of student competence which is the output of the learning process that refers to the Graduate Competency Standards. Apart from that, the results of UN are also useful for mapping the level of students' learning success in school.

The UN results can be accessed openly on the official website of the Ministry of Education and Culture, but the data displayed are not grouped based on certain achievement categories so the readers, in this case the community, do not yet have criteria or references of the achievement of a school based on the results of a national evaluation. The importance of knowing the meaning of data, presented in general regarding the achievement and quality of learning in a school, is commonplace in the current era of digital technology. School clustering aims to support improving the schools quality by reviewing clustering based on the results of the UN for the benefit of

Raoda Ismail, Heri Retnawati, & Okky Riswandha Imawan

school development. UN result data are one of the large amounts of open data that can be grouped into several groups or clusters.

Cluster analysis, a multivariate technique, aims to classify data based on their characteristics. In cluster analysis, each object that is most closely similar to another object is grouped in the same group (Denis, 2020; Rencher, 2001; Toledo, 2005). The clusters formed have high internal homogeneity and high external heterogeneity (Ediyanto et al., 2013; Tinsley & Brown, 2000). Cluster analysis has differences compared to other multivariate analysis techniques because this analysis uses a set of variables determined by the researcher himself, not to estimate the number of variables empirically. Cluster analysis focuses on the comparison of objects based on a set of variables, and therefore, the set of variables is considered by experts as an essential stage in cluster analysis (Denis, 2020). The difference between factor analysis and cluster analysis lies in the focus of the analysis. Factor analysis focuses on groups of variables while cluster analysis focuses on grouping objects (Toledo, 2005).

Clustering, the method used for unlabeled data, is one of the methods in data mining (Estivill-Castro & Yang, 2004; Primartha, 2018). Clustering is an activity of grouping data into clusters or groups based on the similarity of characteristics between one data and another; so that data in the same cluster will have a high level of similarity, and data between clusters will have a low level of similarity (Huang, 1998). The analysis used for clustering, commonly referred to as cluster analysis, aims to form groups of objects so that each object in the same cluster will be bound to one another, and have differences with objects in other clusters (Tan et al., 2019). K-means clustering is one of the algorithms of the clustering method that is used to cluster data based on a similarity in the attributes of the data (Capó et al., 2020; Mahdavi & Abolhassani, 2008; Rencher, 2001; Wu et al., 2008). There are several algorithms that can be implemented for the clustering process using K-means, namely Euclidean distance, Canberra distance, and also Manhattan distance (Faisal et al., 2020; Kapil et al., 2016). K-means, in general, is a heuristic algorithm that can cluster a data set into a number of clusters (K) by optimally reducing the number of squared distances in each cluster. The implementation of the K-means algorithm in this study uses the Euclidean distance method. The clustering process using the K-means algorithm is more optimal in terms of time and the resulting output is of higher quality even though it uses large amounts of data (Hossain et al., 2012; Mavroeidis & Marchiori, 2013; Rajabi et al., 2020).

This study aims to implement the K-Means Clustering Method to group senior high schools in Papua Province based on the 2019 National Examination scores. This can be an important input for schools in general, and for students in particular. For schools, especially those in Papua Province, the results of this study can be valuable information to find out their school's achievements in UN, and can be an evaluation for schools to be able to improve the quality of teaching-learning processes in schools, because the results of UN are the outputs of the learning process shown by students' understanding of the material being tested nationally. For students and parents, especially junior high school students, the results of this study can be valuable information in choosing high school as the school of choice in continuing their studies.

## METHOD

The study of methods for finding patterns from data is data mining. Data mining is a set of steps used to explore previously unlabeled or undefined data or values, as well as data sourced from databases (Wu et al., 2008). Finding meaning from data is a structured process which is described in the following stages (Han & Kamber, 2011): (1) data cleaning is cleaning data from inconsistent data; (2) data integration is the process of combining data from several different sources; (3) data selection is the process of selecting data from the database according to the purpose; (4) data transformation is the process of changing the form of data into data suitable for the mining process; (5) data mining is an important process that uses certain methods to get patterns from data; (6) pattern evaluation is the process of identifying patterns; (7) knowledge pre-

sentation is the one that can represent the required information, and the process by which the information that has been obtained is then used by the data owner.

The grouping of data mining has resulted in several parts, namely description, estimation, and prediction (Nariya et al., 2017). This research began with collecting data that would be implemented by the K-means algorithm. The data in this study, which were collected through the official website of the Center for Educational and Cultural Assessment of the Ministry of Education, Culture, Research, and Technology of the Republic of Indonesia, are data from the 2019 Papua Province National Examination, which consists of data on 144 high schools from 30 cities. Furthermore, the collected data were selected for data selection to obtain a research sample. Then, the data were analyzed using the K-means algorithm, which is one of the algorithms of the clustering method that is used to cluster data based on a similarity in the attributes of the data. (Berry & Maitra, 2019; Wu et al., 2008).

K-Means, in general, is a heuristic algorithm that can cluster a data set into a number of clusters (K) by optimally reducing the number of squared distances in each cluster (Cai & Tang, 2021; Estivill-Castro & Yang, 2004; Mahdavi & Abolhassani, 2008). The clustering process using the K-Means algorithm is more optimal in terms of time and the resulting output is of higher quality even though it uses large amounts of data (Demidenko, 2018; Dorman & Maitra, 2021; Kapil et al., 2016; Rajabi et al., 2020). The stages of the K-means algorithm are as follows (Kapil et al., 2016; Tabachnik & Fidel, 2014; Toledo, 2005): (a) determining the number of k clusters to be searched; (b) selecting the centroid point sequentially or randomly from the initial data as many as k; (C) calculating the distance from each data using the number of k centroid points; (d) each centroid is recalculated based on the obtained cluster mean values; (e) grouping based on the smallest distance; and (f) steps 3-5 are repeated until the smallest distance group no longer experiences a change in pattern. These stages are briefly presented in Figure 1.
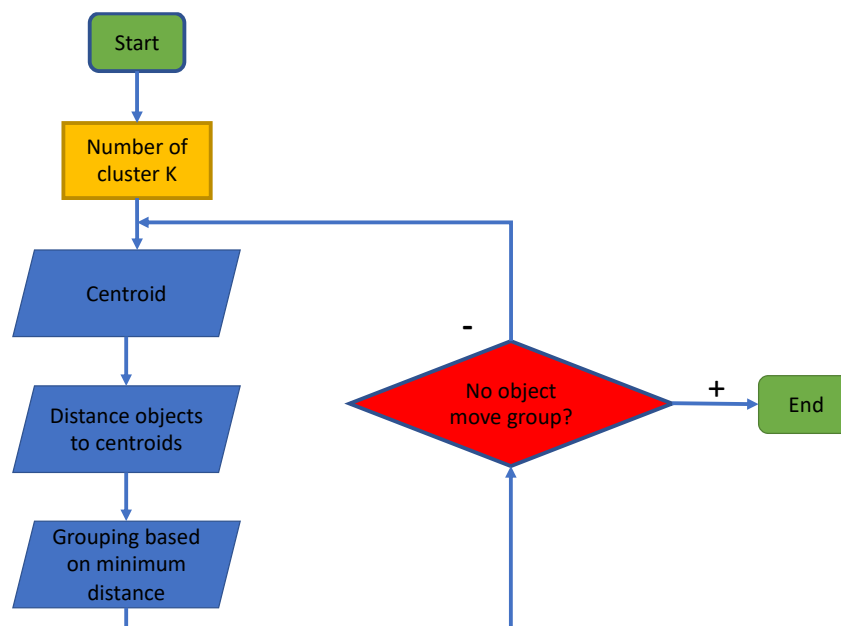


Figure 1. Stages of K-Means Clustering

The K-means clustering stage is a simple but very maximal stage in the data analysis process, so that the analysis results are valid. Figure 1 shows that the analysis stage begins with determining the value of K as the center of the centroid. The initial value of K serves as a manually defined parameter. Next, each datum is added up with all the centroids and the distance between the clusters is calculated, which will produce cluster values for each datum and will be grouped randomly. Then all stages are repeated with a new centroid in each repetition until there is no change in the centroid, and then the data in the cluster can be declared valid (Singh et al., 2013).

There are several methods that can be implemented for the clustering process using K-means, namely Euclidean distance, Canberra distance, and Manhattan distance (Crawford et al., 2021; Faisal et al., 2020; Kapil et al., 2016). The implementation of the K-Means algorithm in this study used the Euclidean distance method. This method was used to measure the distance between two different data (Denis, 2020; Tinsley & Brown, 2000). The Euclidean distance formula was presented in Formula (1), in which d(*i,j*) = Euclidean distance, $X_i$ = point value 1, and $X_j$ = point value 2.

$$d(i,j) = \sqrt{|X_{i1} - X_{j1}|^2 + |X_{i2} - X_{j2}|^2 + \cdots + |X_{ip} - X_{jp}|^2}$$ ………………….. (1)

To find out the weaknesses of the K-means algorithm and get the number of clusters, then in determining the number of clusters needed based on the data used in the analysis, a validity index is needed, which functions as a method to find the results of the clustering algorithm. (Capó et al., 2020; Khairati et al., 2019; Lithio & Maitra, 2018). Therefore, in this study, the validity index, namely Silhouette, was used. The Silhouette validity index can measure the average value of each point in the data set. The measurement was by calculating the difference between the values of separations and compactness divided by the maximum value between the two. The Silhouette value that was getting closer to 1 indicates the best number of clusters. The interpretation of the Silhouette interval value can be seen in Table 1. In addition, the formula for identifying multicollinearity by calculating the Variance Inflation Factor (VIF) is shown in Formula (2).

Table 1. Interpretation of Silhouette Values

| Interval | Interpretation |
| --- | --- |
| < 0.25 | No substantial structure found |
| 0.26 – 0.50 | Weak structure |
| 0.51 – 0.70 | Reasonable structure |
| 0.71 – 1.00 | Strong structure |

$$VIF = \frac{1}{1 - R^2}$$ ………………………. (2)

In Formula (2), $R^2$ is the value of the coefficient of determination of the dependent variable with the independent variable. Multicollinearity is indicated if the VIF value is more than 10. If there is multicollinearity, the variables that are correlated with the model will be excluded to see the ability of the cluster to distinguish the existing data according to the variables or characteristics of the subjects used for clustering. Checks were carried out using the value of R-Square ($R^2$) (Toledo, 2005) as seen in Formula (3), in which $k$ = number of clusters, $n_j$ = number of data in cluster j, $X$ = grand mean, and $X_j$ = average value of each j cluster.

$$R^2 = \frac{SS_b}{SS_T}$$

$$SS_b = \sum_{j=1}^{k} n_j (\bar{X}_j - \bar{\bar{X}})^2$$

$$SS_T = \sum_{j=1}^{k} \sum_{m=1}^{n_j} (X_{jm} - \bar{\bar{X}})^2$$ ……………………… (3)

To check the validity of the cluster solution, it is possible to test the hypothesis that there is no cluster in the population from which the sample is drawn. For example, the hypothesis could be that the population represents a single unimodal distribution such as a multivariate normal, or that observations arise from a uniform distribution (Rencher, 2001). The cross-validation approach can also be used to check the validity or stability of clustering results. The data are randomly divided into subsets A and B, and cluster analysis is performed separately on A and B, respectively. The results should be similar if the clusters are valid (Bilodeau & Brenner, 2000; Huberty & Elejnik, 2007; Rencher, 2001). The analysis results must be similar if the cluster is valid.

## FINDINGS AND DISCUSSION

Schools that are included in the high cluster have total UN scores in the interval 161-217, those in the medium cluster have total UN scores in the 121-160 interval, and those in the low cluster has total UN scores in the 90-120 interval. The interval used is determined by the number of UN scores in Indonesian, Mathematics, and English subjects. UN scores in Physics, Chemistry, and Biology subjects were not included because there were 63 schools that did not have UN scores in one, two, or all of these subjects, so that UN scores that can be included in the cluster analysis are only the scores in the subjects that have UN scores for all high schools in Papua Province, namely Indonesian Language, Mathematics, and English subjects.

Table 2. Change in Cluster Centers

| Iteration | 1 | 2 | 3 |
|---|---|---|---|
| 1 | 7.631 | 11.764 | 14.660 |
| 2 | 1.027 | 1.819 | 7.061 |
| 3 | 1.071 | 0.454 | 3.892 |
| 4 | 0.448 | 0.154 | 0.674 |
| 5 | 0.173 | 0.151 | 0.000 |
| 6 | 0.000 | 0.000 | 0.000 |

Table 2 shows a convergence achieved due to no or small change in cluster centers. The maximum absolute coordinate change for any center is .000. The current iteration is 6. The minimum distance between initial centers is 40.335.
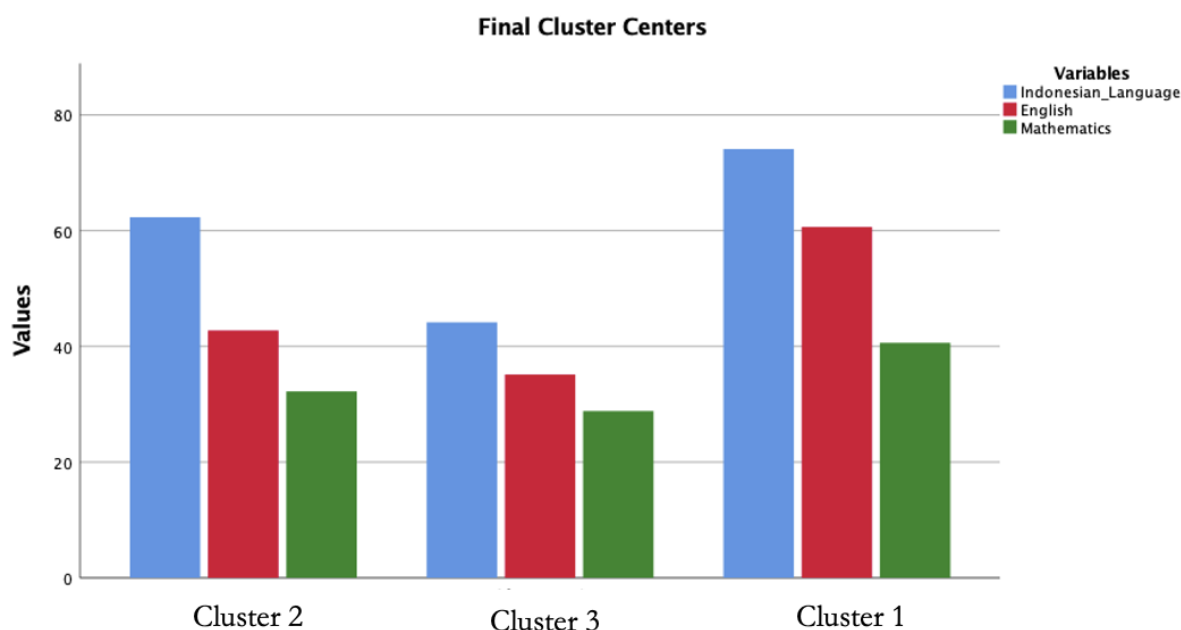


Figure 2. Final Cluster Centers

Raoda Ismail, Heri Retnawati, & Okky Riswandha Imawan

Figure 2 shows that the average UN scores of high schools in Cluster 1 are higher than the average UN scores of high schools in Clusters 2 and 3 on the three variables of UN scores, namely Indonesian Language, English, and Mathematics. To still get an idea of how data aggregate in a high-dimensional feature space, van der Maaten and Hinton developed a t-distributed stochastic neighbor embedding algorithm. This algorithm is used to generate the "t-SNE cluster plot", which can be found under "Plots" (van der Maaten & Hinton, 2008), and the results are shown in Figure 3.
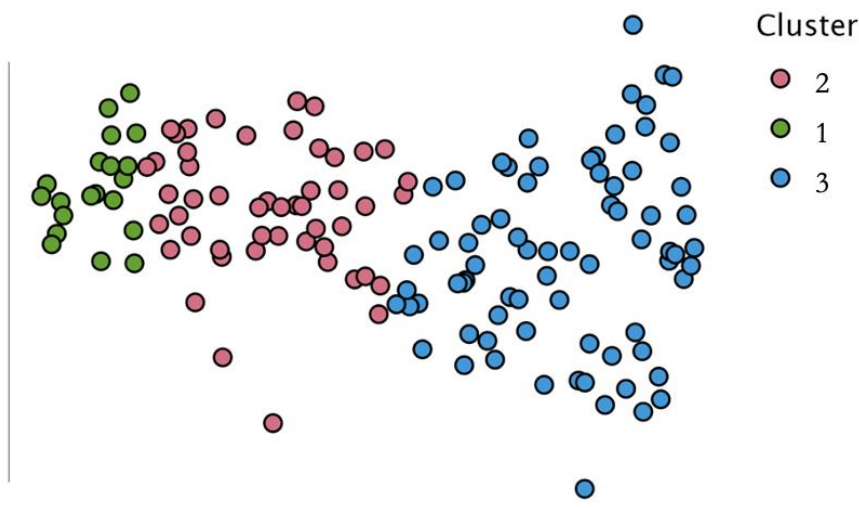
Figure 3. Clustering Result

The t-SNE results in Figure 3, with high probability, project similar data points in the higher-dimensional feature space to points that are adjacent to each other in the plane, and dissimilar data points in the higher-dimensional feature space to points that are not similar to each other in the higher-dimensional feature space points that are far apart from each other. Meanwhile, Table 3 shows that 144 senior high schools in Papua Province are grouped into three clusters. The value of $R^2$ is the ratio between the sum of squares and the total number of squares, which is also usually reported in ANOVA or regression models. A model with $R^2$ whose value is close to the upper limit of one is considered a suitable model, while $R^2$ whose value is close to the lower limit of zero indicates that the model is not suitable. But the $R^2$ makes no difference between a fit model and an overfit model. The value of $R^2$ indicates the magnitude of the coefficient of determination is 0.723, which means that the effective contribution of the variable is 72.3%. The Silhouette score in Table 2 is 0.42, which represents the mean internal consistency of the clustering by assessing how similar each case is with respect to its own cluster compared to other clusters. For Silhouette scores, the general rule is that the closer the grouping to the upper limit of 1, the more consistent it is, while Silhouette scores closer to the lower limit of -1 indicate poor fit.

Table 3. K-Means Clustering

| Clusters | N | R² | AIC | BIC | Silhouette |
|---|---|---|---|---|---|
| 3 | 144 | 0.723 | 136.740 | 163.470 | 0.420 |

The results of the K-means clustering analysis with three clusters show 18 schools in Cluster 1 have high UN scores, 58 schools in Cluster 2 have moderate UN scores, and 68 schools in Cluster 3 have low UN scores. Table 3 presents the results of the classification of senior high schools in Papua Province using K-means clustering with three clusters.

Table 4. Clustering Results

| Cluster | Member | Total | Criteria |
|---|---|---|---|
| Cluster 1 | SMA Negeri 3 Jayapura; SMA Negeri 4 Jayapura; SMA Negeri 5 Jayapura; SMA YPPK Teruna Bakti; MA DDI Entrop; SMA Kristen Kalam Kudus; SMA Wahana Cita Jayapura; MA Persiapan Negeri Koya Barat; MA YPKP Sentani; SMA Lentera Harapan Sentani; SMA Negeri 1 Merauke; SMA YPPK Adhi Luhur Nabire; SMA Kristen Anak Panah Nabire; SMA Negeri 1 Mimika; SMA Advent Timika; SMA YPPK Tiga Raja Timika; SMA Kristen Shining Stars; MA Negeri Keerom. | 18 | High |
| Cluster 2 | SMA Negeri 1 Jayapura; SMA Muhammadiyah Jayapura ; SMA Pembangunan V Yapis Waena; SMA YPPK Taruna Dharma; SMA Negeri 2 Jayapura; SMA Mandala Trikora; SMA Hikmah Yapis Jayapura; MA Daarul Maarif Numbay; SMA Negeri 6 Skouw Jayapura; MA Al Muttaqin Buper Jayapura; SMA Negeri 1 Sentani; SMA Advent Doyo Baru; SMA YPPK Asisi Sentani; SMA Al-Fatah YPKP Sentani; SMA Yapis Nimbokrang; MA Negeri Jayapura; SMA Negeri Kaureh; SMA Persiapan Bumi Sahaja; MA Nurul Anwar; SMA Negeri 1 Biak; SMA Katholik Yos Sudarso; SMA YPK 2 Biak; SMA Yapis Biak; SMA Negeri 3 Biak; SMA Sup Byaki Fyadi; SMA Negeri 1 Serui; SMA Negeri Unggulan Yapen Waropen; SMA Negeri 1 Mearuke; SMA Negeri 3 Merauke; SMA YPPK Yoanes 23; SMA Negeri 1 Kurik; SMA Negeri 4 Merauke; SMA Negeri 1 Muting; SMA KPG Khas Papua Merauke; MA An-Najah Yamra; MA Al-Munawwaroh; MA Al-Hikmah; SMA Plus Muhammadiyah Merauke; SMA Negeri 1 Tanah Miring; SMA Negeri 1 Wamena; SMA Negeri Kurulu; SMA Negeri 1 Nabire; SMA Negeri 2 Nabire; SMA Negeri 3 Nabire; SMA Negeri 6 Nabire; SMA YPK Tabernakel Nabire; SMA Yapis Nabire; SMA Muhammadiyah Nabire; SMA Almadina Nabire; SMA Negeri 2 Mimika; SMA Negeri 4 Mimika; SMA Integral Hidayatullah; SMA Negeri 6 Mimika; SMA Negeri 1 Arso; SMA Negeri 2 Skanto; SMA Pembangunan 6 Yapis Keerom; SMA Santo Arnoldus Janssen; SMA Negeri 1 Sarmi. | 58 | Medium |
| Cluster 3 | SMA PGRI Jayapura; SMA Satria Tasangkapura; SMA 45 Entrop; SMA El-Shaddai Jayapura; SMA YPK Diaspora Kotaraja; SMA Gabungan Jayapura; MA Baiturrahim Jayapura; SMA Negeri Khusus Olah Raga; SMA Negeri 3 Sentani; SMA Negeri 2 Senatani; SMA Kristen Sentani; SMA Negeri Demta; SMA Negeri 1 Nimboran; SMA Negeri Kemtuk Gresi; SMA YPK Sentani; SMA Negeri Yokiwa; SMA Santo Antonius Padua Sentani; SMA Negeri 2 Biak; SMA YPK 1 Biak; SMA Negeri Samber; SMA YPK Immanuel Agung Samofa; SMA Negeri 2 Serui; SMA PGRI Serui; SMA Onate Serui; SMA Yasuka Ansus; MA Darussalam Serui; SMA YPK Penabur Serui; SMA YPPK Yos Sudarso; SMA YPK Merauke; SMA Negeri Plus Urumb; MA DDI Lampu Satu; SMA Negeri Plus Satu Atap 1 Merauke; SMA PGRI Wamena; SMA YPPGI Wamena; SMA YPPK St. Thomas Wamena; SMA Kristen Wamena; SMA YPPGI Nabire; SMA Bhakti Mandala Nabire; SMA Negeri 1 Plus KPG Nabire; SMA Negeri 1 Paniai; SMA Negeri 5 Sentra Pendidikan Mimika; SMA Katolik Santa Maria; SMA Taruna Dharma Timika; SMA YPPK Taruna Tegasa; MA AL-Muhtadin Arso VI; SMA Negeri 4 Arso; SMA YPK Ebenhaezer Sarmi; SMA Negeri 2 Sarmi; SMA Negeri 3 Sarmi; SMA Negeri 3 Yenggarbun; SMA Negeri 5 Warke; SMA Negeri 6 Sowek; SMA Negeri 7 Urmboridodi; SMA Negeri Waren; SMA Negeri Urei Fasei; MA Maarif NU Waropen; SMA Advent Urei Faisei; SMA YPK Fx Mote; SMA Negeri 1 Tanah Merah; SMA Negeri 1 Agats; SMA Negeri 1 Obaa; SMA YPPK Yohanes Paulus II; SMA Negeri 2 Obaa; SMA Negeri 1 Oksibil; SMA YPPK Bintang Timur; SMA Negeri 2 Tigi. | 68 | Low |

Table 4 shows that there are 18 high schools in Cluster 1, 58 high schools in Cluster 2, and 68 high schools in Cluster 3. The difference in the number of schools in each cluster is affected by the distance between UN score data from each school. Schools with high UN scores are in the same cluster, as well as schools with medium and low UN scores are in their respective clusters.

Schools that are included in the high cluster have a total UN score of 161-217, in those in the middle cluster have a total UN score of 121-160, and those in the low cluster have a total UN score of 90-120. The interval used is determined by the number of UN scores in Indonesian, Mathematics, and English subjects. UN scores in Physics, Chemistry, and Biology subjects were not included because there were 63 schools that did not have UN scores in one, two, or all of these subjects, so that UN scores that can be included in the cluster analysis are only scores in the subjects that have UN scores for all high schools in Papua Province, namely Indonesian, Mathematics, and English subjects.

Table 4 also shows that the majority of senior high schools in Papua are still in the lower cluster. This is a special concern for the Papuan Provincial Government to further improve the quality of teaching and learning so that students' understanding can increase as evidenced in the national assessment which will take effect in 2021. The number of high schools in high clusters also shows that the quality of meaningful learning in Papua Province is still very low and small, so that only 18 out of 144 or 12.5% of schools can enter the high cluster. This is a separate mandate for the Papua Provincial government, especially the local Education Office, to further improve the quality of teaching and learning so that more students can achieve a meaningful learning understanding so that they can achieve satisfactory learning outcomes.

This is an input for schools in general, and for students in particular. For schools, especially those in Papua Province, the results of this study can be valuable information to determine the achievement of their schools in UN, and can be an evaluation for schools to be able to improve the quality of teaching and learning in schools, because the results of UN are the output of the learning process shown by students' understanding of the material being tested nationally. For students, especially junior high school students, and their parents, the results of this study can be valuable information in choosing high school as the school of choice in continuing their studies.

The study concludes that cluster formation can have a positive impact on school leadership and local education offices, although the study specifically explores the common characteristics of each school (Kurniadi & Sugiyono, 2020; Lock, 2011; Sutriyani et al., 2018). The findings can be transferred to any group of school leaders who are in the same cluster so that they can work together to improve the quality of teaching and learning which will have an impact on increasing student understanding (Aditya et al., 2020; Chikoko, 2007; Lock, 2011). Thus, fellow schools in the same cluster can get a place to support each other, share expertise, share educational vision, and participate in teacher training to improve teacher competence (Imawan & Ismail, 2020; Ismail & Imawan, 2021a, 2021b). Apart from that, schools in high clusters can be an example or role model for schools in medium and low clusters to remain enthusiastic in improving the quality of teaching and learning and understanding of students so that the quality of education in Papua Indonesia can be maintained.

Schools may come together in a cluster to address special needs they may have in common, such as how to make teaching in their schools more meaningful for students. School clusters can include a variety of activities that involve collaboration between schools. This can be administrative, material, pedagogical, or extracurricular. Resource centers can be located within a school cluster to provide professional and pedagogical support to each school within the same cluster. The result of this research is a decision support system for grouping academic data to see the achievement of each school.

## CONCLUSION

Implementation of the K-Means Algorithm can be applied in the clustering of UN scores, which consists of stages of data cleaning, data integration, data selection, data transformation, data valuing, and evaluation. The results of high school clustering using the K-means algorithm obtained the number of senior high schools. In Cluster 1, there are 18 schools with UN scores in the high category. In Cluster 2, there are 58 schools with UN scores in the medium category. In Cluster 3, there are 68 schools with UN scores in the low category. The results of the analysis of

the K-means algorithm show an $R^2$ value of 0.723 and a Silhouette score of 0.42. It is expected that further research can implement other evaluation approaches or different algorithms to obtain better results and conclusions. Further research is suggested to use a different algorithm and more than three clusters.

## ACKNOWLEDGMENT

## REFERENCES

Aditya, A., Jovian, I., & Sari, B. N. (2020). Implementasi K-means clustering Ujian Nasional sekolah menengah pertama di Indonesia tahun 2018/2019. *Jurnal Media Informatika Budidarma*, *4*(1), 51–58. https://doi.org/10.30865/mib.v4i1.1784.

Berry, N. S., & Maitra, R. (2019). TiK-means: Transformation-infused K-means clustering for skewed groups. *Statistical Analysis and Data Mining: The ASA Data Science Journal*, *12*(3), 223–233. https://doi.org/10.1002/SAM.11416.

Bilodeau, M., & Brenner, D. (2000). *Theory of multivariate statistics*. Springer. https://doi.org/10.5860/choice.37-3391.

Cai, Y., & Tang, C. (2021). Privacy of outsourced two-party K-means clustering. *Concurrency and Computation: Practice and Experience*, *33*(8), e5473. https://doi.org/10.1002/CPE.5473.

Capó, M., Pérez, A., & Lozano, J. A. (2020). An efficient K-means clustering algorithm for tall data. *Data Mining and Knowledge Discovery*, *34*(3), 776–811. https://doi.org/10.1007/S10618-020-00678-9.

Chikoko, V. (2007). The school cluster system as an innovation: Perceptions of Zimbabwean teachers and school heads. *Africa Education Review*, *4*(1), 42–57. https://doi.org/10.1080/18146620701412142.

Crawford, A. M., Berry, N. S., & Carriquiry, A. L. (2021). A clustering method for graphical handwriting components and statistical writership analysis. *Statistical Analysis and Data Mining: The ASA Data Science Journal*, *14*(1), 41–60. https://doi.org/10.1002/SAM.11488.

Demidenko, E. (2018). The next-generation K-means algorithm. *Statistical Analysis and Data Mining: The ASA Data Science Journal*, *11*(4), 153–166. https://doi.org/10.1002/SAM.11379.

Denis, D. J. (2020). *Univariate, bivariate, and multivariate statistics using R: Quantitative Tools for data analysis and data science*. John Wiley & Sons. https://doi.org/10.1002/9781119549963.

Dorman, K. S., & Maitra, R. (2021). An efficient k-modes algorithm for clustering categorical datasets. *Statistical Analysis and Data Mining: The ASA Data Science Journal*. https://doi.org/10.1002/SAM.11546.

Ediyanto, E., Mara, M. N., & Satyahadewi, N. (2013). Pengklasifikasian karakteristik dengan metode K-Means cluster analysis. *Buletin Ilmiah Mat. Stat. Dan Terapannya (Bimaster)*, *02*(2), 133–136. https://jurnal.untan.ac.id/index.php/jbmstr/article/view/3033.

Estivill-Castro, V., & Yang, J. (2004). Fast and robust general purpose clustering algorithms. *Data Mining and Knowledge Discovery*, *8*(2), 127–150. https://doi.org/10.1023/B:DAMI.0000015869.08323.B3.

Faisal, M., Zamzami, E. M., & Sutarman, S. (2020). Comparative analysis of inter-centroid K-means performance using Euclidean distance, Canberra distance and Manhattan distance.

*Journal of Physics: Conference Series*, *1566*. https://doi.org/10.1088/1742-6596/1566/1/012112.

Han, J., & Kamber, M. (2011). *Data mining: Concepts and techniques* (3rd ed.). Morgan Kaufmann.

Hossain, M. S., Ramakrishnan, N., Davidson, I., & Watson, L. T. (2012). How to "alternatize" a clustering algorithm. *Data Mining and Knowledge Discovery*, *27*(2), 193–224. https://doi.org/10.1007/S10618-012-0288-4.

Huang, Z. (1998). Extensions to the K-means algorithm for clustering large data sets with categorical values. *Data Mining and Knowledge Discovery*, *2*(3), 283–304. https://doi.org/10.1023/A:1009769707641.

Huberty, C., & Elejnik, S. (2007). Applied MANOVA and discriminant analysis. *Journal of the American Statistical Association,* *102*(479), 1075-1076. https://doi.org/10.1198/jasa.2007.s203.

Imawan, O. R., & Ismail, R. (2020). Meningkatkan kompetensi guru Matematika dalam mengembangkan media pembelajaran 4.0 melalui pelatihan aplikasi Geogebra. *Jurnal Masyarakat Mandiri (JMM)*, *4*(6), 1231–1239. https://journal.ummat.ac.id/index.php/jmm/article/view/3102.

Ismail, R., & Imawan, O. R. (2021a). Meningkatkan penguasaan TPACK guru di Papua melalui pelatihan pembuatan video pembelajaran pada masa pandemi Covid-19. *Jurnal Masyarakat Mandiri (JMM)*, *5*(1), 277-288. https://journal.ummat.ac.id/index.php/jmm/article/view/3862.

Ismail, R., & Imawan, O. R. (2021b). Optimalisasi kompetensi calon guru Matematika di Papua melalui pembuatan video pembelajaran di masa pandemi Covid-19. *Jurnal Masyarakat Mandiri (JMM)*, *5*(2), 734–745. http://journal.ummat.ac.id/index.php/jmm/article/view/4158.

Kapil, S., Chawla, M., & Ansari, M. D. (2016). On K-means data clustering algorithm with genetic algorithm. In *the 4th International Conference on Parallel, Distributed and Grid Computing*, 202–206. https://doi.org/10.1109/PDGC.2016.7913145.

Khairati, A. F., Adlina, A. A., Hertono, G. F., & Handari, B. D. (2019). Kajian indeks validitas pada algoritma K-means enhanced dan K-means MMCA. In *Prosiding Seminar Nasional Matematika*, *2,* 161–170. https://journal.unnes.ac.id/sju/index.php/prisma/article/view/28906.

Kurniadi, D., & Sugiyono, A. (2020). Pengelompokkan data akademik menggunakan algoritma K-means pada data akademik Unissula. *Jurnal Transformatika*, *18*(1), 93–101. https://doi.org/10.26623/transformatika.v18i1.2277.

Lithio, A., & Maitra, R. (2018). An efficient k-means-type algorithm for clustering datasets with incomplete records. *Statistical Analysis and Data Mining: The ASA Data Science Journal*, *11*(6), 296–311. https://doi.org/10.1002/SAM.11392.

Lock, A. (2011). *Clustering together to advance school improvement: Working together in peer support with an external colleague*. National College for Leadership of Schools and Children's Services.

Mahdavi, M., & Abolhassani, H. (2008). Harmony K-means algorithm for document clustering. *Data Mining and Knowledge Discovery*, *18*(3), 370–391. https://doi.org/10.1007/S10618-008-0123-0.

Mavroeidis, D., & Marchiori, E. (2013). Feature selection for k-means clustering stability: Theoretical analysis and an algorithm. *Data Mining and Knowledge Discovery*, *28*(4), 918–960. https://doi.org/10.1007/S10618-013-0320-3.

Nariya, M., Kim, J. H., Xiong, J., Kleindl, P. A., Hewarathna, A., Fisher, A. C., Joshi, S. B., Schöneich, C., Forrest, M. L., Middaugh, C. R., Volkin, D. B., & Deeds, E. J. (2017). Comparative characterization of crofelemer samples using data mining and machine learning approaches with analytical stability data set. *Journal of Pharmaceutical Sciences*, *106*(11), 3270–3279. https://doi.org/10.1016/j.xphs.2017.07.013.

Primartha, R. (2018). *Buku belajar maschine learning: Teori dan praktek*. Informatika.

Rajabi, A., Eskandari, M., Ghadi, M. J., Li, L., Zhang, J., & Siano, P. (2020). A comparative study of clustering techniques for electrical load pattern segmentation. *Renewable and Sustainable Energy Reviews*, *120*. 109628. https://doi.org/10.1016/J.RSER.2019.109628.

Rencher, A. (2001). *Methods of multivariate analysis* (2nd ed.). John Wiley & Sons.

Singh, A., Yadav, A., & Rana, A. (2013). K-means with three different distance metrics. *International Journal of Computer Applications*, *67*(10), 13–17. https://doi.org/10.5120/11430-6785.

Sutriyani, T. P., Siregar, A. M., & Kusumaningrum, D. S. (2018). Implementasi algoritma K-means terhadap Pengelompokan nilai ujian nasional tingkat SMP di Provinsi Jawa Barat. *Techno Xplore : Jurnal Ilmu Komputer Dan Teknologi Informasi*, *3*(1), 30–36. https://doi.org/10.36805/technoxplore.v3i1.797.

Tabachnik, B., & Fidel, L. (2014). *Using multivariate statistics* (6th ed.). Pearson Education.

Tan, P.-N., Steinbach, M., Karpatne, A., & Kumar, V. (2019). *Introduction to data mining* (2nd ed.). Pearson. https://www.pearson.com/us/higher-education/program/Tan-Introduction-to-Data-Mining-2nd-Edition/PGM214749.html.

Tinsley, H. E., & Brown, S. (2000). Handbook of applied multivariate statistics and mathematical modeling. In *Handbook of applied multivariate statistics and mathematical modeling*. Elsevier Science & Technology Books. https://doi.org/10.1016/b978-0-12-691360-6.x5000-9.

Toledo, M. D. G. (2005). *A comparison in cluster validation techniques*. University of Puerto Rico.

van der Maaten, L., & Hinton, G. (2008). Visualizing Data Using t-SNE. *Journal of Machine Learning Research*, *9*(1), 2579–2605. https://www.jmlr.org/papers/v9/vandermaaten08a.html.

Wu, X., Kumar, V., Quinlan, J. R., Ghosh, J., Yang, Q., Motoda, H., McLachlan, G. J., Ng, A., Liu, B., Yu, P. S., Zhou, Z.-H., Steinbach, M., Hand, D. J., & Steinberg, D. (2008). Top 10 algorithms in data mining. *Knowledge and Information Systems*, *14*, 1–37. https://doi.org/10.1007/s10115-007-0114-2.