OPEN ACCESS

RESEARCH

# K-Prototype Algorithm in Grouping Regency / City in South Sulawesi Province Based on 2020 People's Welfare

## Muhammad Refaldy*, Suwardi Annas, & Zulkifli Rais

Statistics Study Program, Faculty of Mathematics and Natural Sciences, Universitas Negeri Makassar, Indonesia

**Abstract:** Clustering is something that is used to analyze data both in machine learning, data mining, pattern engineering, image analysis and bioinformatics. To produce the information needed for a data analysis using the clustering process, this is because the data has a large variety and amount. Researchers will use the K-Prototype method where this method becomes an efficient and effective algorithm in processing mixed-type data. The K-Prototype algorithm has problems in finding the best number of clusters. So, in this paper, researchers will conduct research by finding the best number of clusters in the K-Prototype method. There are many ways to determine this, one of which is the Elbow method. The determination of this method is seen from the SSE (Sum Square Error) graph of several number of clusters. The results of the clustering formed 2 clusters which were considered optimal based on the value of k that experienced the greatest decrease. The results showed that, cluster 1 is a cluster that has characteristics of people's welfare which is better than cluster 2.

**Keywords**: K-Prototype Algorithm, Cluster, Elbow Method, People's Welfare Indicators

*Corresponding author: Muhammad Refaldy, Statistics Study Program, Faculty of Mathematics and Natural Sciences, Makassar State University, Indonesia

E-mail: refaldy15@gmail.com

## 1. Introduction

Cluster *analysis* is one of the methods in *multivariate* analysis that has the main objective of grouping objects based on the characteristics they have. Cluster analysis groups individuals or research objects, so that each object that is closest in common to other objects is in the same *cluster. Clusters* formed in one cluster have relatively the same characteristics *(homogeneous),* while between *clusters* have different characteristics *(heterogeneous).* This grouping is carried out based on the observed variables (Usman and Sobari, 2013).

According to Mattjik & Sumertajaya (2011) *clustering* algorithms are divided into two types, namely *hierarchical* and *non-hierarchical* methods. *Hierarchical clustering algorithms are* used to group objects in a structured manner based on their similarity in nature and the desired clusters are not yet known. According to Johnson & Wichern (2002) *non-hierarchical clustering* algorithms are used for grouping objects where the number of *clusters* to be formed can be determined in advance as part of the clustering procedure.

The *K-Prototype* algorithm is one of *the clustering* methods based on *partitioning.* This algorithm is the result of the development of the *K-Means* (Huang,1998) algorithm to handle *clustering* on data with numerical and categorical mixed-type attributes. The development carried out by Huang maintains the efficiency of the *K-Means* algorithm in the face of large-sized data and can be applied to numerical and categorical data. The fundamental development of the *k-Prototype* algorithm is in measuring the similarity between objects and their *centroids*

*(prototypes)*. Based on the discussion above, the researcher applied the *k-prototype* algorithm because of the presence of categorical and numerical variables in the people's welfare indicators to group regencies/cities in South Sulawesi Province.

## 2. Literature Review

### 2.1. Clustering Definition

*Clustering* is one of the data mining methods that is *unsupervised,* meaning that this method is applied without training and without a teacher and does not require target output. Where a *cluster* consists of a collection of objects that are similar to one another and different from objects contained in other clusters. The *clustering* algorithm consists of two parts, namely *hierarchically* and *non-hierarchically*. The *hierarchical* algorithm finds *clusters* sequentially where *the* clusters were predefined, while *the non-hierarchical* algorithm determines all groups at any given time (Madhulatha, 2012). *Clustering* can also be said to be a process where grouping and dividing data patterns into several numbers of datasets so that it will form similar patterns and be grouped on the same cluster and separate themselves by forming different patterns into different clusters (Huang et al., 2005). Cluster analysis classifies objects so that each object that is closest in common to other objects is in the same *cluster. The* clusters formed have high internal homogeneity and high external heterogeneity (Prasetyo, 2014).

### 2.2. Elbow Method

The *Elbow* method is a method used to generate information in determining the best number of *clusters* by looking at the percentage of comparison results between the number of *clusters* that will form an elbow at a point. This method provides ideas by choosing the cluster value and then adding the cluster value to be used as a data model in determining the best *cluster*. And besides the resulting percentage of the calculation becomes a comparison between the number of *clusters* added. The results of different percentages of each *cluster* value can be shown using a graph as a source of information. If the value of the first *cluster* with the value of the second *cluster* gives an angle in the graph or the value decreases the most then the value of the *cluster* is best. To get the comparison is to calculate the SSE *(Sum of Square Error)* of each cluster value. Because the greater the number of *cluster* k, the smaller the SSE value will be. SSE formula on *k-prototype*

$$SSE = \sum_{K=l}^{K} \sum_{x_l \in S_K} \|X_i - C_k\|_2^2$$

### 2.3. Similarity Measure

The general form of the measure of similarity is expressed as follows

$$\mathrm{d}(X_i , Z_l) = X_i Z_l \sum_{j=1}^{m} \mathrm{d}( , ) x_{ij} z_{lj}$$

$z_l = [z_{i1}, z_{i2}, \dots , z_{im}]^T$ is *the prototype* for *cluster* l. The measure of similarity for numerical attributes known as *euclidean distances* is shown in the following equation

$$d(X_i , Z_l) = X_i Z_l (\sum_{j=1}^{m_r} (x_{ij}^r - z_{lj}^r)^2)^{½}$$

$x_{ij}^r$ is the value on the numeric attribute j
$z_{lj}^r$ is the average or *prototype* of the numeric attribute to *j cluster l.*
$m_r$ is the sum of numerical attributes.

While the measure of similarity for categorical data is

$$\mathrm{d}(X_i , Z_l) = X_i Z_l \gamma_l \sum_{j=l+1}^{m_c} \mathrm{d})(x_{ij}^c - z_{lj}^c$$

Where *the simple matching similarity measure* for categorical data is

$$\delta(x_{ij}^c - z_{lj}^c) = (x_{ij}^c - z_{lj}^c) \begin{cases} 0 & (x_{ij}^c = z_{lj}^c) \\ 1 & (x_{ij}^c \neq z_{lj}^c) \end{cases}$$

$\gamma_l$ is the weight for the category attribute on *cluster* l whose value is the standard deviation value for the numeric attribute on each *cluster*

$x_{ij}^c$ is the value of the categorical attribute

$z_{lj}^c$ is the mode attribute to j *cluster* l

$m_c$ is the sum of categorical attributes.

He, *modifying the simple matching similarity measure* into equation (6) to increase the similarity of objects in the *cluster* with categorical attributes so that the clustering results are better.

If

$$(d) = (x_{ij}^c - z_{lj}^c) \begin{cases} 1 - \omega(x_{ij}^c, l) & (x_{ij}^c = z_{lj}^c) \\ 1 & (x_{ij}^c \neq z_{lj}^c) \end{cases}$$

$\omega(x_{ij}^c, l)$ is the weighing value for which the value is$x_{ij}^c \omega(x_{ij}^c, l)$

$$\omega(x_{ij}^c, l) = \frac{f(x_{ij}^c | c_l)}{|c_l| f(x_{ij}^c | D)}$$

$f(x_{ij}^c | c_l)$ is the frequency of values in $x_{ij}^c$ *cluster* l

$|c_l|$ is the number of objects in *the cluster* l

$f(x_{ij}^c | D)$ is the frequency of values on the entire dataset. $x_{ij}^c$

### 2.4. K-Prototype Algorithm

The *K-Prototype* algorithm is one of the partitioning-based *clustering* methods. This algorithm is the result of the development of the *k-means* algorithm (Huang,1998) to handle *clustering* on data with numerical and categorical mixed-type attributes. The development carried out by Huang maintains the efficiency *of the k-means* algorithm in the face of large-sized data and can be applied to numerical and categorical data. The fundamental development of the *k-prototype* algorithm is in measuring the *similarity between* objects and their *centroids (prototypes)*. The *k-prototype* algorithm of numerical data is calculated by *euclidean* distance while category data is calculated using *k-modes* distance.

$$d(X_i, Z_l) = X_i Z_l \left( \sum_{j=1}^{m_r} (x_{ij}^r - z_{lj}^r)^2 + \gamma_l \sum_{j=l+1}^{m_c} \delta(x_{ij}^c - z_{lj}^c) \right)^{1/2}$$

There are four main stages to the *K-Prototypes* algorithm. Before heading to the *k-prototypes* algorithm process , first determine the number of *k clusters* to be formed, the limit is at least 2 and the maximum √n or n / 2 where n is the number of data *points*. The stages of *the K-Prototype* Algorithm are:

Step 1 : Determine k with *cluster* initials

Stage 2: Calculate the distance of all data points in the dataset to the initials of the initial cluster, allocate the data points into the *cluster* that has the closest *prototype* distance to the measured object.

Stage 3 : Calculate the center point of the new cluster after all objects are allocated. Then reallocate all the data *points* in the dataset against the new *prototype*.

Stage 4 : If the central point of *the cluster* has not changed or has converged then the algorithm process stops but if the central point is still significantly changing then the process returns to stages 2 and 3 until the maximum iteration is reached or there is no displacement of objects.

### 2.5. People's Welfare Indicators

According to the dictionary of W.J.S Poerwadarminta (1990), prosperous is defined as a state of "safe, safe, and prosperous". So that the meaning of welfare includes security, safety and prosperity. The term people (social) in the narrow sense is related to the social development sector or the development of people's welfare which aims to improve the quality of human life, especially those categorized as unlucky groups and vulnerable groups (groups that have the potential to become poor people). In this case, the development policy of people's welfare generally concerns programs or social services to overcome social problems such as, poverty, displacedness, physical and psychic dysfunction, social tuna, moral tuna, and juvenile delinquency. As a consequence, the notion of people's welfare policy is often interpreted as charitable activities or public assistance carried out by the government for poor families and their children; which social science experts relate to the condition of the "*Human Development Index*", namely: the high low level of people's lives which is seen from three main indicators: the level of life expectancy (expectation of life), the level of education *(literacy education)* and the level *of income (income).*

Although there are no firm boundaries on welfare, the level of welfare in the form of food, education, health, and is often associated with other social protections such as job opportunities, protection of old age, freedom from poverty, and so on. There are ten indicators used to determine the level of welfare, namely age, amount of income, consumption or expenditure of heaven, state of residence, living facilities, health of family members, ease of obtaining health services, ease of entering children into education and ease of obtaining facilities.

## 3. Research Method

### 3.1. Types Of Research

This research uses qualitative descriptive research methods that can explain the indicator variables used, as well as by approaching or exploring people's welfare indicator data.

### 3.2. Data Collection Techniques

The data used in this study is secondary data, namely data obtained indirectly or obtained from existing sources. The data was taken from the Central Statistics Agency of South Sulawesi Province.

### 3.3. Data Analysis Techniques

The data analysis techniques in this study are:

a. Conduct a descriptive analysis to see the characteristics of the variables of people's welfare indicators for each regency/city in South Sulawesi Province
b. Determine the optimum *cluster* using the elbow method
c. Determine the center of the initial *cluster*
d. Calculate the distance between the object and the center of the *cluster* using *the eucledian* distance on the numerical variable and *the distance of k-modes* on the category variable.
e. Update the center *cluster* after the distance between *clusters* is recalculated the difference
f. If the *cluster* center changes then repeat from step 4, otherwise it will be continued in the next process
g. Determine the *output* of the grouping results and describe each group based on its variable characteristics

## 4. Results and Discussion

### 4.1. Descriptive Analysis

Based on research variables, there were 24 samples consisting of regencies/cities in South Sulawesi Province using people's welfare indicators in 2020. The characteristics of the data are created in a descriptive form to describe each variable. The characteristics of the data can be seen in Table 1. as follows.

**Table 1.** Data Characteristics

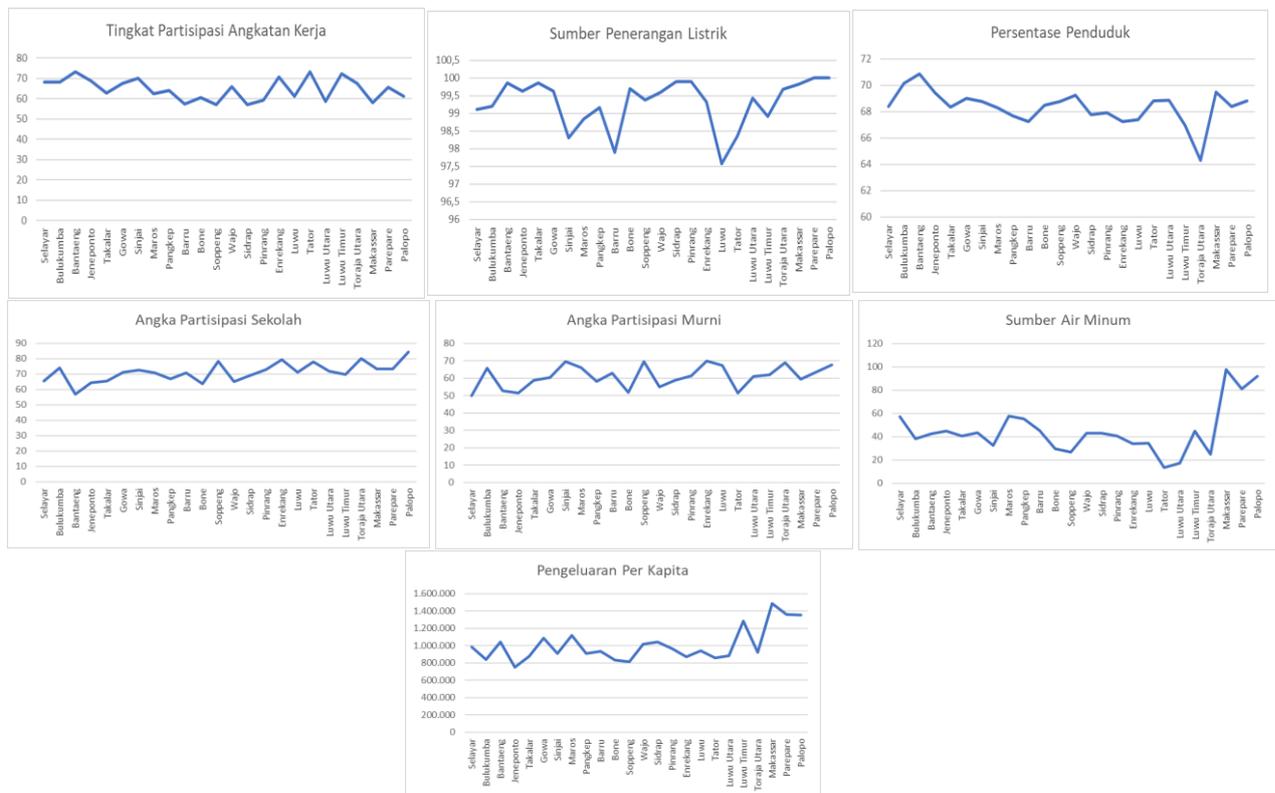| Variable | Minimum | Maximum | Average | Standard Deviation |
|---|---|---|---|---|
| Labor Force Participation Rate (%) | 56,92 | 73,25 | 64,59 | 5,38 |
| Source of Electric Lighting (%) | 97,57 | 100,00 | 99,29 | 0,67 |
| Percentage of Population (%) | 64,32 | 70,86 | 68,36 | 1,26 |
| School Participation Rate (%) | 56,85 | 84,31 | 71,14 | 6,16 |
| Pure Participation Rate (%) | 49,87 | 69,80 | 60,94 | 6,39 |
| Source of Drinking Water (%) | 13,55 | 97,71 | 44,99 | 20,78 |
| Per Capita Expenditure (Rupiah) | 752.620 | 1.489.084 | 1.003.736 | 191.877 |



**Figure 1.** Line Chart of Each Variable

The average labor force participation rate from regencies/cities in South Sulawesi Province in 2020 was 64.59% with a diversity of 5.38%. The highest labor force participation rate was in Tanah Toraja Regency, which was 73.25% while the lowest in Sidrap Regency was 56.92%. The average source of electric lighting from regencies/cities in South Sulawesi Province in 2020 was 99.29% with a diversity of 0.67%. The highest source of electric lighting is in Palopo and Parepare Cities at 100% while the lowest in Luwu Regency is 97.57%. The average percentage of residents from regencies/cities in South Sulawesi Province in 2020 was 68.36% with a diversity of 1.26%. The highest percentage of the population according to the age of 15-64 is in Bantaeng Regency, which is 70.86% while the lowest in North Toraja is 64.32%. The average school participation rate from regencies/cities in South Sulawesi Province in 2020 was 71.14% with a diversity of 6.16%. The highest school participation rate was in Palopo City at 84.31% while the lowest in Bantaeng Regency was 56.85%. The average pure participation rate from regencies/cities in South Sulawesi Province in 2020 was 60.94% with a diversity of 6.39%. The highest pure participation rate was in Enrekang Regency at 69.8% while the lowest in Selayar Regency was 49.87%. The average source of drinking water from regencies/cities in South Sulawesi Province in 2020 was 44.99% with a diversity of 20.78%. The highest source of drinking water is in Makassar City, which is 97.71% while the lowest

in Tanah Toraja Regency is 13.55%. The average per capita expenditure from regencies/cities in South Sulawesi Province in 2020 was Rp. 1,004,736 with a diversity of Rp. 191,877. The highest per capita expenditure is in Makassar City, which is Rp. 1,489,084 while the lowest in Jeneponto Regency is Rp. 752,620.

### 4.2. Determining the number of clusters formed

The *k-prototype* algorithm is a technique for performing *non-hierarchical* groupings on an object. Before the algorithm process is run, it is necessary to determine the number of *clusters* first. Specifying a different number of *clusters* will result in different conclusions or *descriptions* of clusters. The determination of the number of *clusters* is obtained from evaluating *the elbow method* calculations on each *cluster* formed so that *homogeneous clustering* results can be obtained in one *cluster* and *heterogeneous* between clusters. The selection of the number of *clusters* is carried out in stages starting from the number of *clusters* as many as 1 to 8 *clusters*. The results of determining the optimal cluster using the *elbow* method can be seen in Figure 1 as follows.
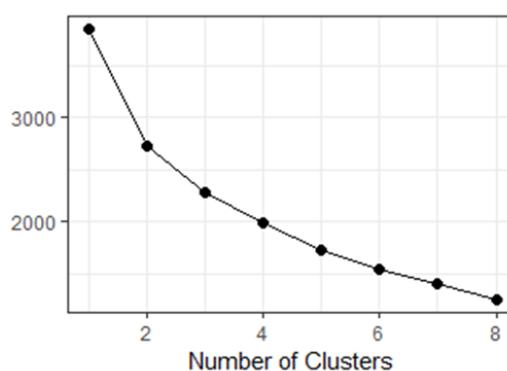


**Figure 2.** Plot the optimal number of clusters

Figure 2 shows that there are several k values that have decreased the most and subsequently the result of the k value will decrease slowly until the result of the k value stabilizes. For example, *the cluster* value k=1 to k=2, then from k=2 to k=3, there is a drastic decrease in forming an elbow at the point k=2 then the ideal *cluster* k value is k=2. The following are the results of district/city membership in each *cluster* presented in Table 2.

**Table 2.** Cluster members

| Cluster | Cluster Members |
|---|---|
| 1 | Bulukumba, Sinjai, Maros, Barru, Soppeng, Sidrap, Pinrang, Enrekang, Luwu, North Luwu, North Toraja, Makassar, Parepare, Palopo |
| 2 | Selayar, Bantaeng, Jeneponto, Takalar, Gowa, Pangkep, Bone, Wajo, Tator, East Luwu |

### 4.3. Interpretation of cluster results

The results of variable processing using *the k-prototype* algorithm with the number of *clusters* as many as 2 and many members of each *cluster*, namely 14 observations in the first *cluster* and 10 observations in the second *cluster*. The labor force participation rate variable shows that the average *cluster* 1 is 62.39% smaller than *cluster* 2 which has an average labor force participation rate of 67.66%. The electric lighting source variable shows that the average *of cluster* 1 is 99.23% smaller than *cluster* 2 which has an average source of electric lighting of 99.38%. The population percentage variable shows that the average *cluster* 1 is 68.10% smaller than *cluster* 2 which has an average population percentage of 68.72%. The variable school participation rate shows that the average *of cluster* 1 is 74.38% greater than *cluster* 2 which has an average school participation rate of 66.61%. The pure participation rate variable shows that the average *of cluster* 1 is 65.06% greater than *cluster* 2 which has an average pure participation rate of 55.19%. The drinking water source variable shows that the average *of cluster* 1 is 47.51% greater than *cluster* 2 which has an average drinking water source of 41.46%. The variable per capita expenditure shows that the average *cluster* 1 of

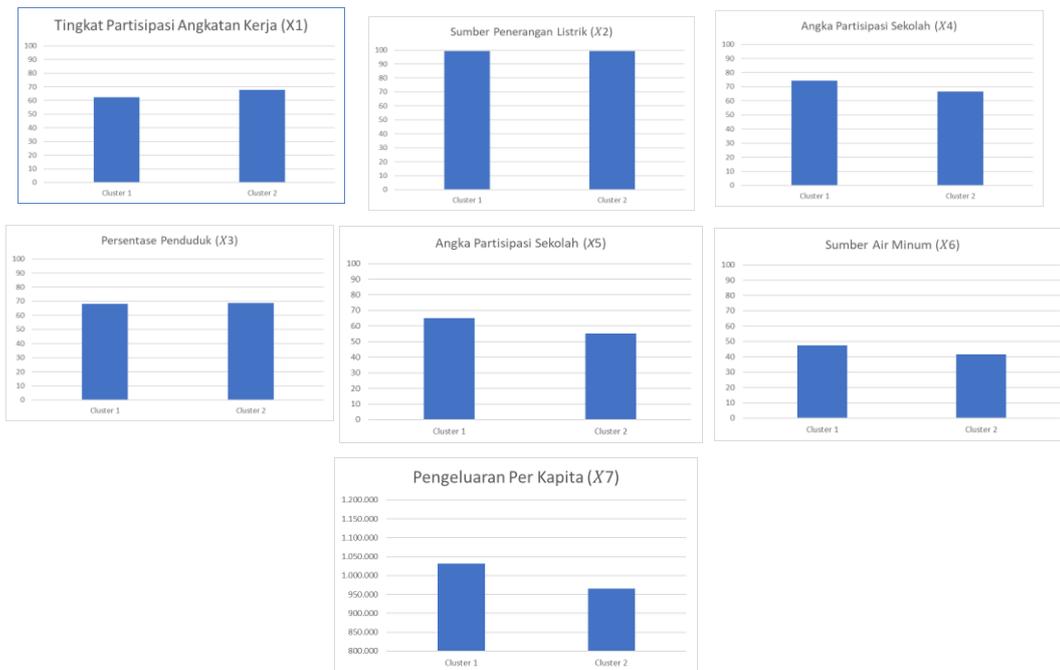Rp. 1,031,400 is greater than *cluster* 2 which has an average per capita expenditure of Rp. 965,007.



**Figure 3.** Comparison Diagram

Based on the bar chart, each variable average value for all indicators of people's welfare in *cluster* 1 consisting of Bulukumba Regency, Sinjai Regency, Maros, Barru Regency, Soppeng Regency, Sidrap Regency, Pinrang Regency, Enrekang Regency, Luwu Regency, Regency North Luwu, North Toraja Regency, Makassar City, Parepare City and Palopo City have 4 variables whose average scores are higher than *cluster* 2, namely in school participation rates, pure participation rates, drinking water sources and per capita expenditure. *Cluster* 2 consisting of Selayar Regency, Bantaeng Regency, Jeneponto Regency, Takalar Regency, Gowa Regency, Pangkep Regency, Bone Regency, Wajo Regency, Tator Regency, East Luwu Regency has 3 variables whose average value is higher than *cluster* 1, namely in the labor force participation rate, electric lighting sources and the percentage of the population. If sorted by the size of the average value of the people's welfare indicators for each *cluster*. *Cluster* 1 is a *cluster* that has better people's welfare characteristics than *cluster* 2.

## 5. Conclusion [11pt, Garamond, Bold, Justified]

The results of the grouping of *observations using the algorima k-prototypes* showed the number of *clusters* formed, namely as many as 2 *clusters*. The determination of the optimum cluster is by using the *elbow* method. The selection of the optimum cluster is based on looking at the *sum of square error* results of the k value which drops drastically. The results showed that *cluster* 1 has 4 variables whose average value is higher than *cluster* 2, namely in school participation rates, pure participation rates, drinking water sources and per capita expenditure. *Cluster* 2 has 3 variables whose average value is higher than *cluster* 1, namely in the level of labor force participation, sources of electric lighting and the percentage of the population aged 15-64. Based on the size of the average value of the people's welfare indicators for each *cluster*, *cluster* 1 is a *cluster* that has better people's welfare characteristics than *cluster* 2.

## References

Amah, N., Wahyuningsih, S., Deny, F., & Amijaya, T. (2017). Analisis Cluster Non-Hirarki Dengan Menggunakan Metode K-Modes pada Mahasiswa Program Studi Statistika Angkatan 2015

FMIPA Universitas Mulawarman Non-Hierarchical Cluster Analysis Using K-Modes Method. *Eksponensial*, *8*, 9–16.

Annas, S., Rahmat, H. S., & Rais, Z. (2022). *K-Prototypes Algorithm For Clustering The Tectonic Earthquake In Sulawesi Island.* *5*(2), 191–198.

Azizah, R. N. (2013). Sistem Informasi Mengklasifikasi Pemilihan Jurusan Di Perguruan Tinggi Bagi Lulusan Sma Berbasis Web Menggunakan Algoritma K-Mean. *Journal of Chemical Information and Modeling*, *53*(9), 1689–1699.

Azwar, Saifuddin.Penyusunan Skala Psikologi / Azwar, Saifuddin .(2017)

Badruttamam, A., Sudarno, S., & Maruddani, D. A. I. (2020). Penerapan Analisis Klaster *K-Modes* dengan Validasi *Davies Bouldin* Index dalam menentukan Karakteristik Kanal Youtube di Indonesia (Studi Kasus: 250 Kanal YouTube Indonesia Teratas Menurut Socialblade). *Jurnal Gaussian*, *9*(3), 263–272. https://doi.org/10.14710/j.gauss.v9i3.28907

Badan Pusat Statistik. (2021) *Indikator Kesejahteraan Rakyat.* BPS: Sulawesi Selatan.

García Reyes, L. E. (2013). Kesejahteraan Masyarakat. *Journal of Chemical Information and Modeling*, *53*(9), 1689–1699.

Hair, J. F., Black, W. C., Babin, B. J. & Anderson, R. E., (2014). Multivariate Data Analysis 7th. USA: Pearson.

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.

Jia, Z., & Song, L. (2020). Weighted k-Prototypes Clustering Algorithm Based on the Hybrid Dissimilarity Coefficient. *Mathematical Problems in Engineering*, *2020*. https://doi.org/10.1155/2020/5143797

Johnson, R. A. & Wichern, D. W., (2002). Applied Multivariate Statistical Analysis 5th. New Jersey: Pearson.

Madhulatha, T. S. (2012). An overview on clustering methods. *arXiv preprint arXiv:1205.1117*.

Maiti, & Bidinger. (1981). Analisis Cluster. *Journal of Chemical Information and Modeling*, *53*(9), 1689–1699.

Mattjik, A. A. & Sumertajaya, I. M., 2011. Sidik Peubah Ganda dengan Menggunakan SAS. Bogor: IPB Press.

Nahak, M. (2017). Bab Ii Tinjauan Pustaka Dan Landasan Teori. *Journal of Chemical Information and Modeling*, *53*(9), 21–25. http://www.elsevier.com/locate/scp

Nooraeni, R., Supriadi, J., Si, S., & Sc, M. (2019). *K-Prototype Untuk Pengelompokan Data Campuran*. *February*, 1–6.

Nooraeni, R., Tinggi, S., & Statistik, I. (2015). Metode Cluster Menggunakan Kombinasi Algoritma Cluster K-Prototype Dan Algoritma Genetika Untuk Data Bertipe Campuran Cluster Method Using a Combination of Cluster K-Prototype Algorithm and Genetic Algorithm for Mixed Data. *Jurnal Aplikasi Statistika & Komputasi Statistik*, *7*(2), 17–17. https://jurnal.stis.ac.id/index.php/jurnalasks/article/view/23

Poerwadarminta, W. J. S. (1990). Kamus Besar Bahasa Indonesia Balai Pustaka, Jakarta p. 1158. *Go to reference in article*.

Prasetyo, E. (2010).Data Mining dan Aplikasi Menggunakan Matlab. Yogyakarta: Penerbit ANDI.

Rachmatin, D. (2014). Aplikasi Metode-Metode Agglomerative Dalam Analisis Klaster Pada Data Tingkat Polusi Udara. *Infinity Journal*, *3*(2), 133. https://doi.org/10.22460/infinity.v3i2.59

Rachmatin, D., & Sawitri, K. (2016). *Perbandingan Antara Metode Agglomeratif, Metode Divisif dan Metode K-Means Dalam Analisis Klaster. 1*, 9–17.

Setyawan, A. H., & Pratiwi, N. (2019). Penerapan Metode Two Step Cluster Untuk Pengelompokan Potensi Desa. *Jurnal Statistika Industri Dan …*, *4*(2), 41–51. https://journal.akprind.ac.id/index.php/Statistika/article/view/1923

Sopha, B. M. (2018). *Analisis Klasterisasi Industri Kecil Menengah di Kabupaten Banyuasin, Provinsi Sumatera Selatan dengan Algoritma K-Prototypes ANDIKA YUSUF PUTRA, Bertha Maya Sopha, S.T., M.Sc., Ph.D.*

Sulastri, S., Usman, L., & Syafitri, U. D. (2021). K-prototypes Algorithm for Clustering Schools Based on The Student Admission Data in IPB University. *Indonesian Journal of Statistics and Its Applications*, *5*(2), 228–242. https://doi.org/10.29244/ijsa.v5i2p228-242

Suryono, A. (2018). Kebijakan Publik Untuk Kesejahteraan Rakyat. *Transparansi Jurnal Ilmiah Ilmu Administrasi*, *6*(2), 98–102. https://doi.org/10.31334/trans.v6i2.33

Usman, H., & Sobari, N. (2013). Aplikasi multivariate untuk riset pemasaran. *Jakarta: PT. RajaGrafindo.*