Cyberspace: Jurnal Pendidikan Teknologi informasi Volume 7, Nomor 2, Oktober 2023, hal. 118 - 125 ISSN 2598-2079 (print) | ISSN 2597-9671 (online)

PRINCIPAL COMPONENT K-MEANS SOFT CONSTRAINT BASED ON WELL-BEING INDICATORS IN ACEH PROVINCE

Winny Dian Safitri

Department of Economics, Faculty of Islamic Economics and Business, Ar-Raniry State Islamic University, Banda Aceh, 23111 E-mail: winny.diansafitri@ar-raniry.ac.id

Abstract

The success of government policies can be from the state of the well- being indicators. This research was conducted to obtain district/city groupings based on the similarity of characteristics of the well-being indicators of each district/city in Aceh Province in 2022. The data used in the Aceh well-being indicator data for 2022 consists of 29 variables. The analysis method used is the principal component kmeans soft constrain method. The background information data can be used as a provision to streamline the clustering algorithm by creating soft constraints which is found as the most appropriate algorithm. The results of this study indicate there are four district/city clusters in Aceh Province. The characteristics of the first cluster are that kindergarten and elementary school facilities are adequate, while the school enrollment rate needs to be improved. The characteristics of the second cluster are superior to the Gross Enrollment Rate (GER) and the population of university graduates, but still very lacking in school facilities. The third cluster is the cluster that is the center of well-being in Aceh, so this cluster is the cluster with the best well-being level. The characteristic of the fourth cluster is that it is very good in the school participation rate indicator, but it must increase early childhood school participation.

Keywords: Well-being, Clustering, Principal Component, K-means, Constrain

Abstrak

Keberhasilan kebijakan pemerintah dapat digambarkan dari keberhasilan pembangunan yang tergambarkan dari kesejahteraan masyarakat. Penelitian ini dilakukan untuk memperoleh pengelompokan kabupaten/kota berdasarkan kesamaan karakteristik indikator kesejahteraan khususnya sektor pendidikan masing-masing kabupaten/kota di Provinsi Aceh pada tahun 2022. Data yang digunakan dalam penelitian ini yaitu data indikator kesejahteraan Aceh tahun 2022 terdiri dari 29 variabel. Metode analisis yang digunakan adalah analisis komponen utama K-means dengan batasan. Data informasi latar belakang dapat digunakan sebagai bekal untuk memproses algoritma clustering dengan membuat soft constraint yang ditemukan sebagai algoritma yang paling tepat. Hasil penelitian ini menunjukkan terdapat empat klaster kabupaten/kota di Provinsi Aceh. Karakteristik klaster pertama adalah fasilitas TK dan SD memadai, sedangkan tingkat pendaftaran sekolah perlu ditingkatkan. Karakteristik klaster kedua lebih unggul dibandingkan Gross Enrollment Rate (GER) dan populasi lulusan universitas, namun masih sangat kurang di fasilitas sekolah. Klaster ketiga adalah klaster yang menjadi pusat kesejhateraan di Aceh, sehingga klaster ini merupakan klaster dengan jenjang kesejahteraan terbaik. Karakteristik klaster keempat adalah sangat baik dalam indikator tingkat partisipasi sekolah, tetapi harus meningkatkan partisipasi sekolah anak usia dini.

Kata Kunci: Kesejahteraan, Klaster, Komponen Utama, K-mean, Batasan

1. Introduction

Prosperity is a dream for all countries of the world. In the eyes of a country is said to be prosperous, if the education is good. Education is one aspect that becomes the benchmark in development success. In the education system, there is a transfer of knowledge between humans. According to Law No. 20 of 2003, the definition of education is an effort and planning to create a controlled teaching and learning process by including spiritual abilities to create independent, moral, and knowledgeable humans. Education will be successful if it has a transparent system and facilities that are met.

The United Nations Development Programs (UNDP) since 1990 have started to issue annual reports on human development in various countries, namely the Human Development Report. Some of the approaches used in measuring poverty levels include proper education. Several developed countries, such as Indonesia, should already have a sound education system. The Indonesian Central Bureau of Statistics formulates education indicators, which are present the development of Indonesian education over time and compare provinces and areas of residence.

Education in Indonesia, Aceh Province, is notably still inferior, from readiness to availability of facilities. Based on the National Examination results issued by the Ministry of Education and Culture of the Republic of Indonesia, Aceh is ranked 27th out of 34 Provinces in Indonesia. It very concerns about the education world in Aceh.

The government's efforts to improve the education system are not yet optimal, so it is necessary to have an education strategy in studies with existing data. The study results were included in policymaking to achieve an intelligent community by the "Aceh Carong" program promoted by the Aceh government to know the obstacles in each district/city in the education system with different levels of difficulty.

Research related to factors that influence education, including planning, regulations, human resources, technical, coordination, and procurement of goods and services, affect the realization of the education budget in Aceh Province. The factors that influence the success of district/city education in Aceh must be different, given the region's geographic location and the various cultures, so it is necessary to have an analysis that can identify each region. One of the statistical methods used is cluster analysis with the primary condition that there must be no relationship between variables (multicollinearity).

When analyzing the data, there were several problems, including a high correlation between variables. This high correlation can cause a high bias value from the analysis results. Principal Component Analysis is a solution to overcome high correlation problems between variables by reducing the dimensions of large variables. The hope is to provide more accurate results from clusters. The use of PCA in k-means to reduce the high dimensional data showed that the PCA can be able to be a solution to produce the better grouping result by using k-means.

The other problem that could occur when analyzing the data is missing value. In order to continue the analysis, the missing values must be addressed first. One of the cluster analysis methods that can handle the problem of missing value without imputation is Kmeans Soft Constraint method. Hence, the Principal Component Analysis K-means Soft Constraint method is used in this research.

PRINCIPAL COMPONENT K-MEANS SOFT CONSTRAINT BASED ON WELL-BEING INDICATORS IN ACEH PROVINCE

2. Literature Review

Principal Component Analysis

The principal component analysis is a reasonably good method for obtaining estimator coefficients in regression equations with multicollinearity problems. The independent variable in principal component regression is a linear combination of the original Z variable called the principal component. This method's estimation coefficient is obtained from the shrinkage of the main component dimensions, with the subset of the main components selected having to maintain the most remarkable diversity.

The method of eliminating the principal component starts from the procedure for selecting the root feature of an equation:

$$|AX - \lambda I| = 0 \qquad (1)$$

If the root of the feature λ_j is ordered from the most massive value to the smallest value, then the effect of the main component W_j corresponds to λ_j . These components explain the proportion of diversity to the dependent variable Y, which is getting smaller and smaller. The main components of W_j are orthogonal to each other and are formed through a relationship:

$$W_j = v_{1j} Z_1 + v_{2j} Z_2 + v_{3j} Z_3 + \dots + v_{pj} Z_p \quad (2)$$

where *p* is the number of variables used. The vector v_j is obtained from each feature root λ_j which satisfies a homogeneous system of equations:

$$/AX - \lambda_j I / v_j = 0 \qquad (3)$$

where $v_j = (v_{1j}, v_{2j}, v_{3j}, ..., v_{pj})$.

There are three methods commonly used to determine the number of main components, namely:

1. If the number of main components produced is q where $q \le p$, then what has been transformed (main component score data) has as many variables as q. Suppose the proportion for the root of trait i^{th} is:

$$\frac{\lambda_i}{\sum_{i=1}^p \lambda_i} \quad (4)$$

The determination of the number of principal components (q) is based on the cumulative proportion of its characteristic roots. If the cumulative proportion of q, the first feature root is 80% or more, then the number of principal components is q [4].

- 2. The main components' selection is based on the variety of the main components, which are none other than the root features. According to, after doing a study, the better cut off is 0.7.
- 3. The scree plot is a plot between the root features λ_k and k. Using this plot, the number of principal components selected is k. If at that point k, the plot is steep to the left but not steep to the right.

K-means Soft Constraints

The background information data can be used as a provision to streamline the clustering algorithm by creating soft constraints. Soft constraints are a function made as initial information from members of a group. The use of constraint becomes important for several clustering algorithms. Several researchers have shown that constraints can improve the results of a variety of clustering algorithms. Thus, the soft constraint addition in k-means algorithm becomes beneficial.

K-means soft constraint is a developed k-means algorithm that has robustness for grouping a set of data without any imputation process required. Missing value issue on a dataset may limit the use of clustering methods. Hence, a method that resistance for clustering the dataset which contains missing values is required. It's obtained that kmeans soft constraint on dataset Glass, Wine, Iris, and Breast Cancer outperformed kmeans for all datasets. Beside the robustness for dealing with missing value, k-means soft constraint also showed its performance for dealing with dataset that contains multicollinearity. It's found that the use of k-means soft constraint results in high accuracy on dataset contains multicollinearity.

It is explained that k-means clustering is an algorithm that is most appropriate for grouping with soft constraints, so that soft constraints are used as information in the grouping. K-means with soft constraints is done by dividing the data set into two parts, namely the set with complete data variables (F_o) and the set with incomplete data variables (F_m). Suppose sc is the symbol for soft constraints, F_m is the set of the incomplete data variable, x_{im} is the item of the *i*th object of the incomplete data variable m, x_{jm} is the item of the jth object of the missing data variable m, f is an incomplete variable member. The soft constraints of x_{im} and x_{jm} are:

$$sc = -\sqrt{\sum_{f \in F_m} (x_{im}f - x_{jm}f)^2}$$
 (5)

where *sc* is always negative. It indicates that one object has different groups. The k-means soft constraints algorithm adopts the steps of the k-means algorithm in dividing k objects into c suitable groups. The stages of the k-means soft constraints algorithm are:

- 1. Determine the center of the c^{th} band.
- 2. Determine the member of the c^{th} band by calculating the minimum distance of an object to the k^{th} band to the center of the c^{th} band

$$C = \frac{\arg\min}{c_c} \left((1 - w) \frac{v}{v_{\max}} + w \frac{cv}{cv_{\max}} \right) \quad (6)$$

by calculating the distance from the k^{th} object to the c^{th} center of the complete data variable *d* is as follows:

$$v = \sum_{d=1}^{D} (x_{kd} - c_{cd})^2$$
(7)

where:

 c_{cd} = center of the c-th band based on the-d variable.

w = weighting factor determined subjectively by value w $\in [0,1]$, in this study using the value of w = 0.5.

 v_{max} = the maximum distance from all objects to the center of the cluster on the complete data variable.

 $cv = sum of the squares of soft constraints containing the value of sc <math>cv_{max} = sum of squares of all soft constraints$

3. Repeat steps 1 to 2 through $\max_{l} \left(|c_{cd}^{(r)} - c_{cd}^{(r-1)}| \le 10^{-4} \right)$ (8)

3. Material and Method

The data used in this study is secondary data, namely well-being indicators based education statistics data for 2022 sourced from the Central Statistics Agency of Aceh Province. The software that used to conduct the analysis is R version 3.4.1. The research variables consisted of indicators of the participation rate of children aged 3-6 years in the Early Childhood Education program. It consisted of 3 variables, the School Participation Rate, which consisted of 3 variables, the Gross Enrollment Rate, which consisted of 3

PRINCIPAL COMPONENT K-MEANS SOFT CONSTRAINT BASED ON WELL-BEING INDICATORS IN ACEH PROVINCE

variables, the net enrollment rate, which consists of 3 variables, the percentage of the population ten years and over is detailed according to the highest diploma held which consists of 6 variables and the number of schools consisting of 11 variables.

The step carried out in this research is started by presenting the overview of Aceh literacy rate. The next step is the implementation of PCA. This is aimed to reduce the dimensions of 29 variables into the smaller dimensions. The result of the dimension reduction is then grouped by using cluster analysis namely k-means soft constraints. The clusters obtained were identified based on the characteristics of each cluster.

4. Result and Discussion

Aceh's well-being-based education indicators have changed significantly from year to year. Along with these changes, an analysis of the education indicator variables in Aceh will be carried out and classifying the districts in Aceh Province based on their similar characteristics.

From the results of the correlation analysis, if it's looked at the correlation in the correlation matrix R measuring p x p = 29×29 (p is the number of observed variables), there are several high correlations between independent variables that indicate multicollinearity, which may be due to different units of measurement. The correlation matrix is shown in Figure 1.

	1 1	a 13		(4)	5 x6	5 x7							13 X	14 x	15 K											126 X2			29
	1	0.11	-0.42	0	-0.01	0	0.23	-0.19	0.23	0.07	0.14	0.28	0.01	0.02	0.03	0.01	-0.09	-0.05	0.01	-0.17	-0.07	-0.05	-0.09	-0.09	-0.03	0.06	0.05	-0.05	-0.1
2	0.11	1	-0.95	0.01	-0.2	0.27	-0.29	-0.17	0.47	0.17	-0.03	0.45	-0.14	-0.05	-0.5	0.23	0.26	0.19	-0.27	-0.38	-0.41	-0.27	-0.41	-0.41	-0.27	-0.39	-0.42	-0.08	-0.1
3	-0.42	-0.95	1	-0.01	0.19	-0.25	0.19	0.22	-0.5	-0.18	-0.02	-0.51	0.12	0.04	0.45	-0.21	-0.21	-0.16	0.24	0.4	0.4	0.27	0.4	0.4	0.26	0.33	0.37	0.09	0.2
1	0	0.01	-0.01	1	0.42	0.29	-0.3	0.06	0.18	-0.18	0.18	0.37	-0.49	-0.5	-0.3	0.53	0.38	0.53	0.05	0.21	-0.26	-0.16	-0.2	-0.02	-0.12	0.01	0.04	0.43	0.3
5	-0.01	-0.2	0.19	0.42	1	0.34	0.04	-0.04	0.28	0.23	0.33	0.32	-0.18	-0.39	-0.25	0.24	0.44	0.44	-0.15	0.04	-0.31	-0.24	-0.17	-0.13	-0.1	-0.1	-0.14	0.36	0.
5	0	0.27	-0.25	0.29	0.34	1	0.18	0.14	0.5	0.44	0.32	0.74	-0.13	-0.66	-0.47	0.39	0.66	0.63	-0.59	-0.26	-0.72	-0.66	-0.71	-0.62	-0.57	-0.62	-0.57	0.22	0.1
7	0.23	-0.29	0.19	-0.3	0.04	0.18	1	0.28	-0.1	0.1	0.34	0.06	0.48	0.07	0.37	-0.38	-0.14	-0.27	-0.1	-0.01	0.05	0.06	0.03	0.03	-0.05	-0.08	0.15	-0.26	-0.1
3	-0.19	-0.17	0.22	0.06	-0.04	0.14	0.28	1	-0.46	-0.24	0.79	-0.26	0.21	0.29	0.09	-0.21	-0.12	-0.33	-0.02	0.19	0	-0.22	-0.02	0	-0.25	-0.18	-0.08	-0.36	-0.1
9	0.23	0.47	-0.5	0.18	0.28	0.5	-0.1	-0.46	1	0.28	-0.17	0.68	-0.3	-0.57	-0.33	0.36	0.55	0.65	-0.19	-0.28	-0.24	-0.05	-0.25	-0.2	0.05	-0.15	-0.13	0.49	0.3
10	0.07	0.17	-0.18	-0.18	0.23	0.44	0.1	-0.24	0.28	1	0.09	0.35	0.19	-0.38	-0.28	0.05	0.36	0.38	-0.68	-0.53	-0.64	-0.56	-0.65	-0.65	-0.51	-0.37	-0.63	0.17	-0.1
11	0.14	-0.03	-0.02	0.18	0.33	0.32	0.34	0.79	-0.17	0.09	1	0.17	0.26	0.1	-0.04	-0.22	0.07	-0.07	-0.2	-0.03	-0.3	-0.39	-0.25	-0.24	-0.39	-0.28	-0.27	-0.18	-0.0
12	0.28	0.46	-0.51	0.37	0.32	0.74	0.06	-0.26	0.68	0.35	0.17	1	-0.06	-0.6	-0.38	0.23	0.51	0.7	-0.41	-0.29	-0.61	-0.38	-0.61	-0.52	-0.37	-0.48	-0.41	0.33	0.2
13	0.01	-0.14	0.12	-0.49	-0.18	-0.13	0.48	0.21	-0.3	0.19	0.26	-0.05	1	0.39	0.3	-0.87	-0.5	-0.34	-0.2	-0.07	0.13	0.02	0	-0.02	-0.11	0.04	-0.09	-0.29	-0.1
14	0.02	-0.05	0.04	-0.5	-0.39	-0.66	0.07	0.29	-0.57	-0.38	0.1	-0.6	0.39	1	0.29	-0.6	-0.79	-0.86	0.29	0.02	0.46	0.25	0.46	0.32	0.22	0.21	0.17	-0.7	-0.
15	0.03	-0.5	0.45	-0.3	-0.25	-0.47	0.37	0.09	-0.33	-0.28	-0.04	-0.38	0.3	0.29	1	-0.56	-0.45	-0.54	0.31	0.35	0.58	0.67	0.45	0.45	0.24	0.36	0.52	-0.3	-0.1
16	0.01	0.23	-0.21	0.53	0.24	0.39	-0.38	-0.21	0.36	0.05	-0.22	0.23	-0.87	-0.6	-0.56	1	0.58	0.49	-0.06	-0.02	-0.37	-0.31	-0.23	-0.17	-0.05	-0.17	-0.08	0.33	0.1
17	-0.09	0.26	-0.21	0.38	0.44	0.66	-0.14	-0.12	0.55	0.35	0.07	0.51	-0.5	-0.79	-0.45	0.58	1	0.76	-0.24	-0.13	-0.51	-0.34	-0.48	-0.37	-0.29	-0.35	-0.34	0.6	0.4
18	-0.05	0.19	-0.16	0.53	0.44	0.63	-0.27	-0.33	0.65	0.38	-0.07	0.7	-0.34	-0.86	-0.54	0.49	0.75	1	-0.29	-0.18	-0.53	-0.38	-0.51	-0.41	-0.21	-0.26	-0.34	0.82	0.6
19	0.01	-0.27	0.24	0.05	-0.15	-0.59	-0.1	-0.02	-0.19	-0.68	-0.2	-0.41	-0.2	0.29	0.31	-0.06	-0.24	-0.29	1	0.56	0.79	0.72	0.86	0.85	0.84	0.7	0.85	0.03	0.3
20	-0.17	-0.38	0.4	0.21	0.04	-0.26	-0.01	0.19	-0.28	-0.53	-0.03	-0.29	-0.07	0.02	0.35	-0.02	-0.13	-0.18	0.56	1	0.57	0.51	0.49	0.71	0.45	0.48	0.65	0.04	0.3
21	-0.07	-0.41	0.4	-0.26	-0.31	-0.72	0.05	0	-0.24	-0.64	-0.3	-0.61	0.13	0.46	0.58	-0.37	-0.51	-0.53	0.79	0.57	1	0.88	0.92	0.94	0.85	0.78	0.87	-0.12	0.1
22	-0.06	-0.27	0.27	-0.16	-0.24	-0.66	0.06	-0.22	-0.05	-0.56	-0.39	-0.38	0.02	0.26	0.67	-0.31	-0.34	-0.38	0.72	0.51	0.88	1	0.78	0.82	0.76	0.62	0.84	-0.05	0.1
23	-0.09	-0.41	0.4	-0.2	-0.17	-0.71	0.03	-0.02	-0.25	-0.65	-0.25	-0.61	0	0.46	0.45	-0.23	-0.48	-0.51	0.86	0.49	0.92	0.78	1	0.92	0.91	0.82	0.88	-0.1	0.2
24	-0.09	-0.41	0.4	-0.02	-0.13	-0.62	0.03	0	-0.2	-0.65	-0.24	-0.52	-0.02	0.32	0.45	-0.17	-0.37	-0.41	0.85	0.71	0.94	0.82	0.92	1	0.85	0.83	0.93	-0.05	0.3
25	-0.03	-0.27	0.26	-0.12	-0.1	-0.57	-0.06	-0.25	0.05	-0.51	-0.39	-0.37	-0.11	0.22	0.24	-0.06	-0.29	-0.21	0.84	0.45	0.85	0.76	0.91	0.86	1	0.77	0.86	0.16	0.3
26	0.06	-0.39	0.33	0.01	-0.1	-0.62	-0.08	-0.18	-0.15	-0.37	-0.28	-0.48	0.04	0.21	0.36	-0.17	-0.35	-0.26	0.7	0.48	0.78	0.62	0.82	0.83	0.77	1	0.78	0.16	0.4
27	0.06	-0.42	0.37	0.04	-0.14	-0.57	0.15	-0.08	-0.13	-0.63	-0.27	-0.41	-0.09	0.17	0.52	-0.08	-0.34	-0.34	0.86	0.65	0.87	0.84	0.88	0.93	0.86	0.78	1	0.03	0.3
28	-0.05	-0.08	0.09	0.43	0.36	0.22	-0.26	-0.36	0.49	0.17	-0.18	0.33	-0.29	-0.7	-0.3	0.33	0.6	0.82	0.03	0.04	-0.12	-0.05	-0.1	-0.05	0.16	0.16	0.03	1	0.8
29	-0.13	-0.19	0.22	0.39	0.3	0.12	-0.15	-0.18	0.37	-0.12	-0.09	0.21	-0.19	-0.5	-0.11	0.14	0.43	0.61	0.33	0.32	0.19	0.16	0.23	0.31	0.39	0.45	0.31	0.86	

Figure 1. Correlation matrix between variables

Multicollinearity can be overcome by principal component analysis by first standardizing the X variables into Z variables and selecting the component dimensions, which must have a cumulative diversity of more than 70 percent.

Table 1. Main component selection								
Eigen value	Variance percent	Cumulative						
	, minine bei en	variance percent						
10.55463796	36.3953033	36.3953						
5.496869248	18.95472154	55.35002						
2.847149847	9.817758094	65.16778						
2.131925502	7.351467249	72.51925						
1.787848096	6.164993436	78.68424						
1.164694421	4.01618766	82.70043						
1.036020135	3.572483225	86.27291						
0.827721401	2.854211729	89.12713						
	Eigen value 10.55463796 5.496869248 2.847149847 2.131925502 1.787848096 1.164694421 1.036020135	Eigen valueVariance percent10.5546379636.39530335.49686924818.954721542.8471498479.8177580942.1319255027.3514672491.7878480966.1649934361.1646944214.016187661.0360201353.572483225						

Table 1. Main component selection

Grouping objects using the grouping method. Principal Component Analysis of K-Means Soft Constraint is to see the distance between objects, but the initial process is to deal with multicollinearity problems first. If the distance value for each object is small, it will be grouped into one cluster. The following are the district/city clusters' results in Aceh based on the 2022 education indicators.

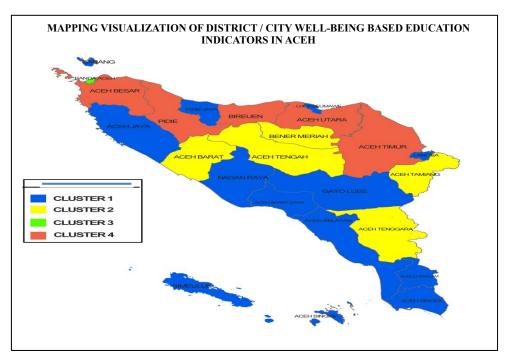


Figure 2. Results of the Aceh education cluster in 2022

Based on Figure 2, it was obtained four district/city clusters in Aceh Province based on well-being-based education indicator data for 2022, namely:

- 1. Members of the first cluster are Aceh Jaya, Pidie Jaya, Nagan Raya, Aceh Barat Daya, Aceh Selatan, Gayo Lues, Subulussalam, Aceh Singkil, Sabang, Lhokseumawe, Langsa, and Simeulue. The first cluster's characteristics are that kindergarten and elementary school facilities are adequate, while the school enrollment rate needs improvement.
- 2. Members of the second cluster are Aceh Barat, Aceh Tengah, Bener Meriah, Aceh Tamiang, and Aceh Tenggara. The second cluster characteristics are superior in the Gross Enrollment Rate and university graduates' population, but still lacking in school facilities.
- 3. Members of the third cluster, namely Banda Aceh. Banda Aceh is the center of education in Aceh Province. The third cluster characteristic is the cluster that is the center of education in Aceh, so that this cluster is the cluster with the best education level.
- 4. Members of the fourth cluster are Aceh Besar, Pidie, Bireuen, Aceh Utara, and Aceh Timur. The fourth cluster characteristic is that it is very good in the School Participation Rate indicator, but it must increase early childhood school participation.

From the cluster above, it can be seen that there should be more programs to improve the quality and quality of teaching staff as well as the distribution of teachers with national education standards evenly in every district/city in Aceh. These efforts were made

PRINCIPAL COMPONENT K-MEANS SOFT CONSTRAINT BASED ON WELL-BEING INDICATORS IN ACEH PROVINCE

primarily to reduce the gap in education level and quality between districts/cities in Aceh. The level of disparities in district/city education development in Aceh must be minimized so that every Acehnese has the same opportunity to get proper education up to the highest level. Expectations from the high level of public education will automatically increase the standard of living so that just and equitable welfare can be realized in people's lives. However, the initial process is to deal with multicollinearity problems first. If the distance value for each object is small, it will be grouped into one cluster. The following are the district/city clusters' results in Aceh based on the 2022 education indicators.

5. Conclusion

Based on the results of the clustering of districts/cities in Aceh Province using the Principal Component K-Means Soft Constraint method, it shows that the members of each cluster are strongly influenced by geographic location, the proximity between districts/cities in Aceh so that the similar characteristics of the Aceh education indicators in 2022 are formed. The results of this study indicate there are four district/city clusters in Aceh Province. The first cluster's characteristics are that the kindergarten and elementary school facilities are adequate, while the school participation rate needs to be improved. The second cluster's characteristics are superior in the Gross Enrollment Rate and the population of university graduates, but still lacking in school facilities. Characteristics of the third cluster is the cluster with the best education level. The fourth cluster's characteristic is that it is very good in the School Participation Rate indicator, but it must increase early childhood school participation.

References

- 1. Aini E N, Isnaini I, Sukamti S & Amalia L N. 2018. Pengaruh Tingkat Pendidikan Terhadap Tingkat Kesejahteraan Masyarakat di Kelurahan Kesatrian Kota Malang. *Tecnomedia Journal* 3(1) 1-15.
- 2. BPS. 2022. Aceh Dalam Angka 2022. Publication Of BPS Provinsi Aceh.
- 3. Ding, C dan He, X. 2004. K -means Clustering via Principal Component Analysis. *ICML 2004: Proceedings of the Twenty-First International Conference on Machine Learning*.
- 4. Johnson R, A & Wichern D, W. 2014. *Applied Multivariate Statistical Analisys*. New Jersey, Prentice Hall International.
- 5. Jolliffe I, T. 2016. Principal Component Analysis. New York: Springer.
- 6. Majid M S, A. 2014. Analisis Tingkat Pendidikan dan Kemiskinan di Aceh. Jurnal Pencerahan, vol 8 no 1. Jurnal Pencerahan, **8**(1) 15-27.
- 7. MacQueen J, B. 1967. Some Methods for Classification and Analysis of Multivariate Observations. *Proceedings of the Fifth Symposium on Math, Statistics and Probability: 281–297.* University of California Press.
- 8. Mesquita, D, Gomes, J dan Rodrigues, L. (2016). K-means for Datasets with MissingAttributes Building Soft Constraints with Observed and Imputed Values. *European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning*, Bruges.
- 9. Ruwaida, Darwanis, S.Abdullah. 2015.Faktor-faktor yang Mempengaruhi Realisasi Belanja Pendidikan di Provinsi Aceh. *Jurnal Akutansi, vol 4 no 4*.

- 10. Safitri W D, Nurhasanah and Rusyana A. 2012. Pengelompokan Kabupaten/ kota di Provinsi Aceh Berdasarkan Tingkat Perubahan Kesejahteraan Rakyat. Skripsi Universitas Syiah Kuala, Banda Aceh.
- 11. Safitri W D. Sartono B dan Erfiani. 2016. Kajian Penggerombolan Data Tidak Lengkap dengan Algoritma Khusus Tanpa Imputasi. Tesis. Institut Pertanian Bogor, Bogor.
- 12. Sambandam, R. 2003. Collinearity Is a Natural Problem in Clustering. 15(1), 16–21.
- 13. Wagstaff K, Cardie C, Rogers S & Schroed S. 2001. Constrained K-means Clustering With Background Knowledge. *Proceedings. of the 18th Intl. Conf. on Machine Learning : 577–584.* USA.
- 14. Wagstaff, K. 2004. Clustering with Missing Values: No Imputation Required. *Proceedings of the Meeting of the International Federation of Classification Societies* : 649–658. California.
- 15. Wagstaff, K. L., Basu, S. dan Davidson, I. 2006. When Is Constrained Clustering Beneficial, And Why? Ionosphere, 58(60.1), 62-63.