

Hybrid Method Air Quality Classification Analysis Model

By Syafri Arlis

Hybrid Method Air Quality Classification Analysis Model

Musli Yanto ¹, Syafril Arlis ², Devia Kartika ³, Deri Marse Putra ⁴

Abstract – This paper aims to present a discussion of air quality based on the classification analysis developed by hybrid method. Problems that occur, where the longer the air quality is getting worse and can cause serious problems for human life. The method in performing the analysis used consisted of K-Means clustering, Multiple Linear Regression (MRL), Artificial Neural Network (ANN), and Decision Tree Algorithm C.45. The results of MRL measurement show that correlation relationship is same as the variable with an output of 70.7%. Then the results of determination with ANN showed an MSE value of 0.0018197 and output accuracy of 99.99%.

Keywords – Analysis Model, Classification, Hybrid Method, Knowledge Based, Quality Air.

1. Introduction

Air pollution often experiences a high enough increase to have an impact on environmental pollution and human health[1]. Broadly speaking, the perceived health impacts are respiratory problems[2]. Related to this, the negative effects that will be felt in the future such as damage to the lungs, heart, and other organs[3]. For the results to be generated, and analysis process is needed in determining the status of air conditions that will occur next in the short and long term, so that adverse impacts can be anticipated early on.

In the classification of air quality, previous analytical models are generated, such as the model developed to perform a multivariable air pollution prediction process[4]. Another discussion also states that the model using the MultiLayer Perceptron (MLP) gives good accuracy results in determining air quality[5]. Other studies also explain that air quality by adopting mathematical calculations of variables that affect air pollution[6]. The analytical model of the Multiple Linear Regression (MRL) and Artificial Neural Network (ANN) methods can be used to predict air quality[7]. MRL and ANN are performance models that are quite good at making predictions and providing accurate results[8].

This paper discusses the air quality classification analysis model. The update contained in this study presents the results of grouping air quality status based on output for the evaluation of knowledge-

based system analysis. The evaluation results are taken into consideration in the control and monitoring process. The hybrid method is used to enhance the analysis model that already exists. The stages of analysis in the proposed model start from (I) the preprocessing analysis stage with K-Means dataset clustering to produce a classification pattern, (II) Multiple Linear Regression (MRL) is used to measure the correlation of variables in the resulting pattern, (III) Artificial Neural Network (ANN) learning process. (IV) Classification using the C.45 decision tree method to provide an overview of the classification results on air quality status in the form of a decision tree.

The clustering process can group objects originating from facts or events to find stored information [9],[10]. Discussion in the case of clusters is also able to group factors that affect air quality[11]. Multiple Linear Regression (MRL) is a method to test the correlation between variables and outputs. This correlation test process will be can to provide an overview of the results measurement variables[12]. In previous studies, MRL can to the relationship between two or more variables and was able to measure the relationship between the predictor variables used[13]. Previous studies have explained that MRL can analyze factors that affect air quality with a fairly minimal error rate[14].

ANN is a supervised learning method used in solving problems with optimal solutions[15]. ANN can do learning to provide a form of knowledge-based evaluation[16]. In other cases, ANN is used to identify based on network variables and models[17]. In the same case, ANN has shown significant and superior results in determining air quality[18]. C.45 is an algorithm in the Decision Tree data mining method used in the making form of patterns that contain information and knowledge[19]. This algorithm is a process used to convert facts into information[20]. The C.45 algorithm is used to find new knowledge[21].

From the explanation above, this paper presents a good concept and analysis model from the previous model. This study also presents another update, namely a new analytical model in the classification of air quality status. The working representation of this model is seen based on the output of the data mining clustering process forming a classification

pattern. After that, the MRL will be used to test the accuracy of the pattern in the given correlation relationship. The hybrid method in the classification analysis model provides precise and accurate results of air quality status. The purpose of this paper is to present the concept of a more structured classification analysis model on the status of air quality and be taken into consideration in decision making. So that the developed model will answer all the doubts of the previous model.

2. Research Method

The classification analysis model that is carried out in the process of determining air quality with this hybrid method provides a structured form of analysis in determining air quality. The following analysis model can be seen in Figure.1 below:

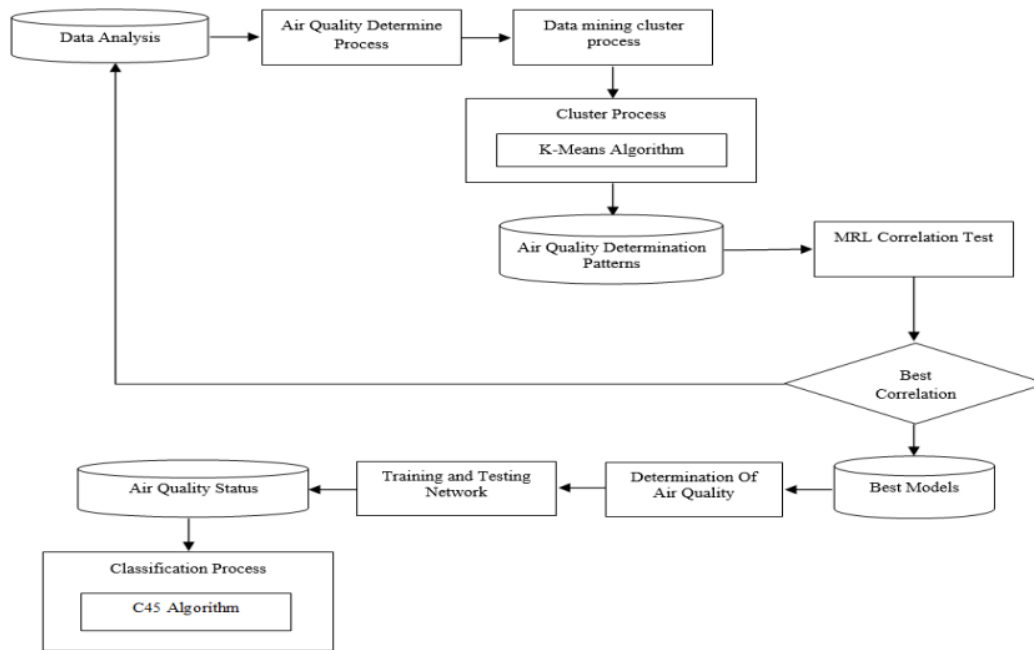


Figure.1 Air Quality Analysis Model

Figure.1 above describes the analysis model that will be carried out in determining air quality. The following steps are carried out:

1. The data analysis stage is the stage in determining the variables that affect air quality.
2. Stages of clusters with data mining methods. At this stage is the preprocessing analysis stage. The clustering process will use the K-means algorithm.
3. Stages of variable correlation test using MRL. This stage aims to test the relationship level of the previously generated rule pattern variables. The output obtained provides the measurement results of the variables used in determining the status of air quality.
4. Stages of ANN learning. This stage will carry out the training and testing process on the previously formed rule pattern. The initial process starts from the formation of the network architecture to arrive at the final process by classifying air quality status. The

output results given will be measured by the percentage level of accuracy, sensitivity, and the resulting error value.

5. The final stage in the classification analysis model uses the decision tree algorithm C.45 method. The results given from this process provide an overview of the air quality status rules in the form of a decision tree.

K-Means Cluster

A cluster process is an approach that is widely used in data mining[22]. The purpose of the K-Means method is to use grouping by mining data to generate information[23]. Discussion, K-means is used to analyze air quality status data. Cluster analysis begins by determining the number of clusters (C1, C2, C3, and C4). After the number of clusters is determined, the process of calculating the distance from the cluster can be continued by using Formula.1 below[24]:

$$D_e = \sqrt{(M_{ix} - C_{ix})^2 + (M_{iy} - C_{iy})^2} \quad (1)$$

Where:

- D_e = Euclidean Distance
- M_{ix}, M_{iy} = Data Object Coordinates
- C_{ix}, C_{iy} = Centroid Coordinate Center

Multiple Regression Linear (MRL)

The multiple Linear Regression (MRL) method is used to measure the level of relationship between variables[25]. MRL in principle performs a knowledge-based evaluation based on the correlation between variables[26]. The equations in the MRL can be seen in the Formula.2 & 3 below[27]:

$$E(Y|X=x_1, X_2=x_2, \dots, X_m=x_m) = \beta_0 + \beta_1 x_1 + \dots + \beta_m x_m \quad (2)$$

$$Y = \beta_0 + \beta_1 x_1 + \dots + \beta_m x_m + \epsilon, E\epsilon = 0, D\epsilon = DY = \alpha^2 \quad (3)$$

Artificial Neural Network (ANN)

The air quality classification process uses ANN, the approach used is backpropagation. The approach to learning in a network is in the form of a series of feedforward processes or is called reverse inference[28]. Learning is applied to mathematical calculation models and acceptable logic to produce decisions[29]. The network architecture pattern that

is built consists of an input layer, a hidden layer, and an output layer[30]. This process aims to produce an optimal network by learning variations of the number of hidden layers[31].

Metode Decision Tree

The Decision Tree method is a concept that performs the classification process with the results of the decision tree as the output[32]. The output of this method forms a series of trees and is widely used in describing information[33]. A decision tree has nodes and nodes. Each node represents a decision based on the attributes or variables of the dataset. In each of them, there will be a label and a value from the range of attribute values[34].

3. Results And Discussions

Analysis Cluster K-means

In the analysis air quality status classification, the variables used include: (X1) Ozone (O3), (X2) Carbon Monoxide (CO), (X3) Air Particles (PM10), (X4) Nitrogen Dioxide (NO2), (X5) Sulfur Dioxide (SO2), and (X6) Air Pollution Standard Index (ISPU). This variable is found from data on air content conditions in the Padang City area, West Sumatra Province, Indonesia, which occurs for 1 year. The results of the sample K-means cluster analysis process can be seen in Table.1 below:

Table. 1 Air Quality Status Classification Pattern

X1	X2	X3	X4	X5	X6	Y	X1	X2	X3	X4	X5	X6	Y
58.89	22.284	166.5	14.1	0	54	Medium	46.2	0.573	271.2	9	0.082	46.2	Dangerous
78.96	23.947	148.6	17.2	0	64.48	Medium	51.2	0.576	221.7	10.5	0.065	50.6	Dangerous
82.46	22.617	112.7	17.7	0	66.23	Medium	59.4	0.516	551.5	10.4	0.079	54.7	Healthy
55.63	23.501	96.4	14	0	52.83	Medium	94.7	0.531	671.8	17.2	0.07	72.4	Healthy
38.4	22.428	115.5	13.5	0	38.4	Medium	100.9	0.558	414.5	18.6	0.028	75.5	Healthy
35.69	24.171	81.1	12.8	0	35.69	Medium	44.7	0.54	245.2	7.7	0.054	44.7	Dangerous
72.38	25.639	137.1	18.2	0	61.23	Medium	9.6	0.544	368.2	9.1	0.04	9.6	Dangerous
77.46	20.705	279.4	17.4	0	63.73	Not Healthy	34.7	0.551	376.6	11.3	0.062	34.7	Dangerous
70.25	20.473	254	14.9	0	60.1	Not Healthy	55.7	0.52	254.8	18.8	0	52.9	Dangerous
57.12	22.962	264.1	11.1	0.037	53.56	Not Healthy	22.8	0.516	281.9	9.3	0.075	22.8	Dangerous
103.46	18.437	228.2	13.4	0.085	76.73	Not Healthy	34	0.545	289.6	8.6	0.052	34	Dangerous
159.21	24.355	242.7	24.7	0.03	105	Not Healthy	59.4	0.766	317.7	14.1	0.054	54.7	Dangerous
275.21	22.266	247.4	34.8	0	225.21	Dangerous	63.3	0.82	314.7	19	0	56.7	Dangerous
73.75	0.5693	221.3	19.8	0.013	61.88	Dangerous	45.2	0.844	329.2	9.1	0.041	45.2	Dangerous
41.04	0.6157	147.6	9.3	0.04	41.04	Medium	39.9	0.761	227.9	10.4	0.053	39.9	Dangerous
41.92	0.5986	133.3	9.5	0.084	41.92	Medium	11.7	0.856	236.2	14.1	0	11.7	Dangerous
54.04	0.5719	116.6	13.1	0.082	52.05	Medium	21.2	0.933	260.5	10.2	0.041	21.2	Dangerous
60.25	0.5241	106.7	11.3	0.069	55.1	Medium	20.5	0.91	226.4	10	0.081	20.5	Dangerous
46.83	0.5521	122.5	6	0.099	46.83	Medium	22.3	0.837	180.1	10.5	0.054	22.3	Not Healthy

33.58 0.6013 246 9.7 0.102 33.58 Dangerous 16.5 0.838 302.67 10.2 0.037 16.5 Dangerous

Table.1 above, shows that the clustering process forms a classification pattern of air quality. The visualization of the results of the clusterization of air quality status can be seen in Figure.2 below:

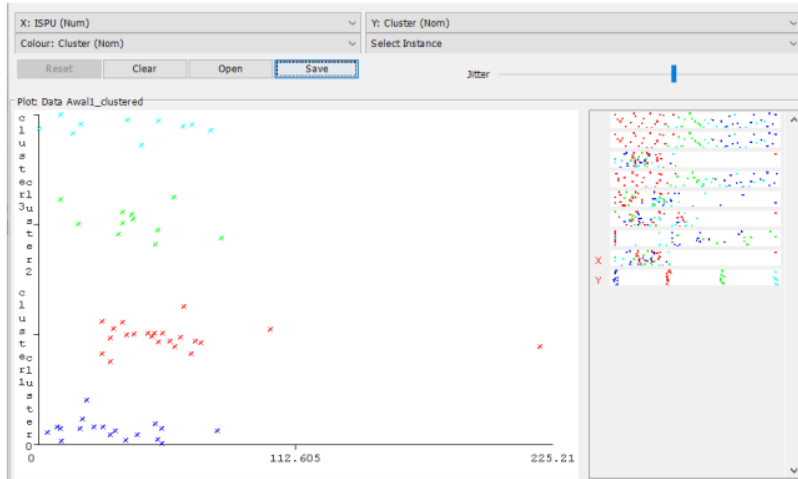


Figure.2 Visualization of Air Quality Status Clustering

Figure.2 explains that the cluster results consist of C1 = Healthy by 16%, C2 = Not Healthy by 18%, C3 = Medium by 30% and C4 = Dangerous by 36%. From the results of this cluster analysis, it can be concluded that based on the average of the dataset used, the air quality is in Dangerous status. To be able to see the relationship between variables and the output of air quality status, the analysis process is continued in the analysis process using MRL.

In this MRL analysis, the process will start from testing the level of the relationship between variables in determining air quality. The process of regression can be used in measuring the relationship between predictor variables[35]. In its implementation, MRL can model the independent variable (X) and the dependent variable (Y)[36]. This process can be used as a parameter to see the relationship based on the factors that affect the results[37]. The MRL process for the Coefficient of Determination Test that has been carried out can be seen in Table.2 below:

Analysis MRL

Table. 2 Results of the Coefficient of Determination

Model Summary				
Model	R	R Square	Adjusted R	Std. Error
1	.707 ^a	.500	.440	.40044

a. Predictors: (Constant), ISPU, NO, CO, O3, SO2, PM10

Table.2 above shows that the results of the measurement of the relationship between the variable and the dependent result produce a result of 70.7%. These results are sufficient to prove that the variables used can affect air quality. To re-test the analysis Table. 3 Variable Correlation Test Results

produced on the results of the determination test contained in Table.2, the correlation test was carried out to see the relationship between the variables and the air quality status. The results of the correlation test can be seen in Table. 3 below:

		Correlations					
Control Variables		PM10	SO2	CO	O3	NO	ISPU
Status	PM10	1,000	0,654	-0,179	0,562	-0,251	0,815
Quality	Correlation		0,000	0,172	0,000	0,053	0,000
	Significance (2-tailed)						
	df	0	58	58	58	58	58

Air	SO2	1 Correlation	0,654	1,000	-0,272	0,302	-0,532	0,540
		Significance (2-tailed)	0,000		0,036	0,019	0,000	0,000
		df	58	0	58	58	58	58
	CO	1 Correlation	-0,179	-0,272	1,000	-0,035	0,066	-0,159
		Significance (2-tailed)	0,172	0,036		0,793	0,619	0,003
		df	58	58	0	58	58	58
	O3	1 Correlation	0,562	0,302	-0,035	1,000	-0,435	0,402
		Significance (2-tailed)	0,000	0,019	0,793		0,001	0,005
		df	58	58	58	0	58	58
	NO	1 Correlation	-0,251	-0,532	0,066	-0,435	1,000	-0,196
		Significance (2-tailed)	0,053	0,000	0,619	0,001		0,002
		df	58	58	58	58	0	58
	ISPU	Correlation	0,815	0,540	-0,159	0,402	-0,196	1,000
		Significance (2-tailed)	0,000	0,000	0,003	0,005	0,002	
		df	58	58	58	58	58	0

Table. 3 explains that the relationship between each variable has shown good results with air quality status. This can be seen from the significant value which states less than 0.005. PM10 variable obtained a significant value of $0.000 < 0.005$, SO2 variable of $0.000 < 0.005$, CO variable of $0.003 < 0.005$, NO variable of $0.002 < 0.005$ and ISPU variable of $0.004 < 0.005$. This is in line with research conducted[38], which states that NO, CO, and SO2 are related to air pollution. Later in the same study also explained that PM10 and SO2 have a good relationship in influencing air quality[39]. In another paper also

stated that the element of NO compounds can also affect air quality[40]. In the research that has been done, it is also explained that CO has a significant level of coefficient on the level of air quality[41]. In this case, the MRL analysis can describe the relationship between variables and gives good enough results to test the pattern generated in the previous analysis process so that it can be used in the process of determining air quality. In the analysis that is also carried out in conducting the F-test, it can be seen in Table.4 below:

Table. 4 F Test Results

ANOVA ^a						
	Model	Sum of Squares	df	Mean Square	F	Sig.
1	Regression	8.017	6	1.336	8.333	.000 ^b
	Residual	8.018	50	.160		
13	Total	16.035	56			

a. Dependent Variable: Quality Air

b. Predictors: (Constant), ISPU, NO, CO, O3, SO2, PM10

Table. 4 illustrates that the F-test results give a significant result of 0.000 which is lower than 5%. This indicates that the variables that have been used together can be used in determining air quality. After the MRL analysis process is carried out, this identification pattern can be continued in the learning process in classifying air quality.

ANN Machine Learning

The resulting ANN learning process gives good results in analyzing the classification of air quality status. The results are given with an MSE value of 0.000333 and an accuracy rate of 99.99%. In addition, the performance value is 0.0009997 and the MAPE value is 0.000544048. The process of analyzing the status air quality that has been carried out can be seen in Figure.3 below:

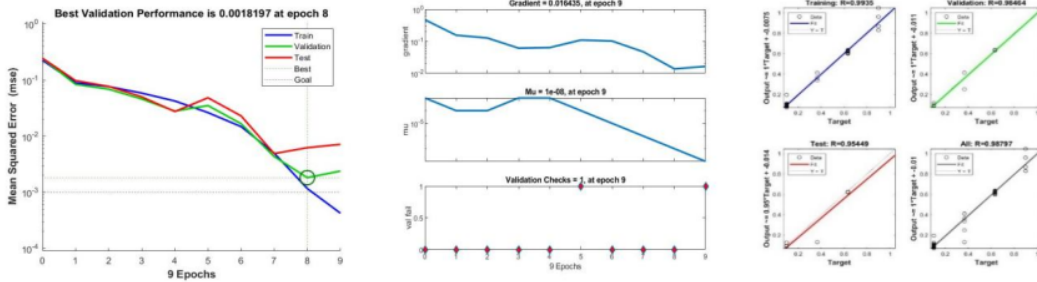


Figure.3 ANN Learning Outcomes

Figure.3 can be seen that the ANN learning output graph presents quite good results. The testing process for the classification pattern gives quite good results based on the MSE graph which has a value of 0.0018197. While the resulting validity value is 98.464%. The classification process is continuing in the decision tree analysis process to see the pattern that will be described in a decision tree in the case of air quality.

The pattern obtained will be re-analyzed by comparing the results of the cluster carried out using the C.45 approach. This process is able to describe the relationship of the grouping to the data used. The results obtained are able to provide the same results with the aim of classifying the data[42]. In this case, the C.45 method also gives results in the form of a decision tree from the formed pattern[43]. In the process carried out using the C.45 approach, the decision tree found in Figure.4:

Classification of Decision Tree

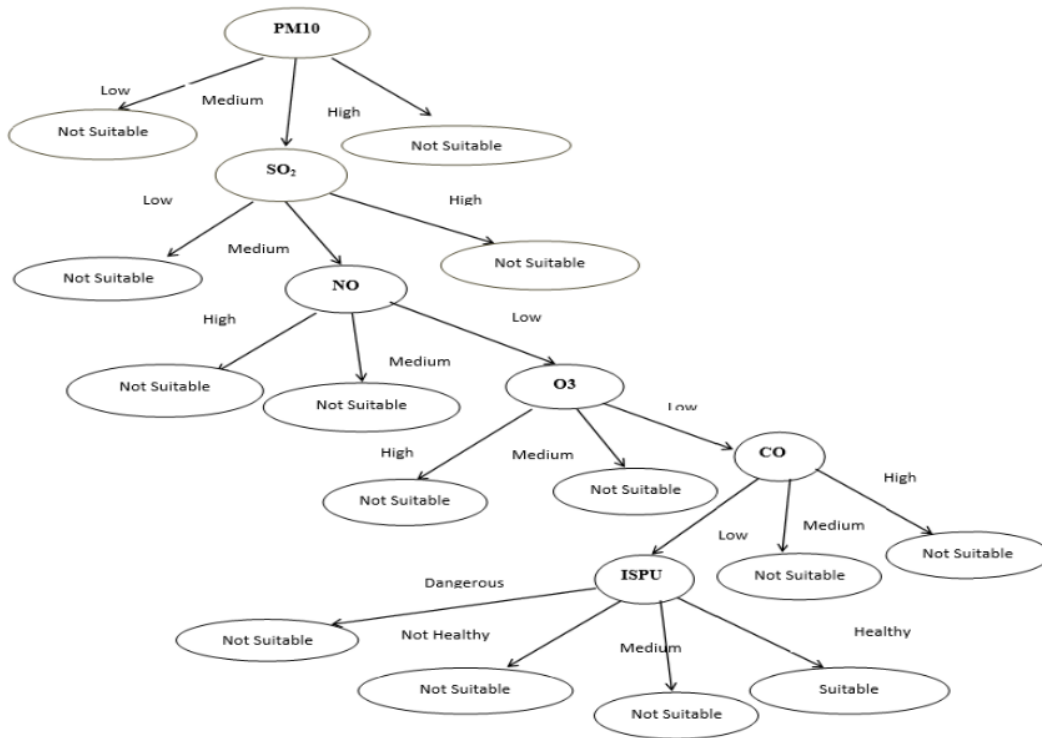


Figure.4 Decision Tree Results

14

Figure.4 explains that the decision tree generated by the decision tree algorithm C.45 can describe the pattern of air quality status. In the

picture, it can be seen that the notes and knots used are derived from the patterns obtained previously. This form of decision tree description can be used for

controlling and monitoring the development of air quality status.

4. Conclusion

The analysis that has been carried out using the hybrid method approach provides a structured and systematic analysis model in the process of determining air quality. The proposed method can optimize the classification analysis process previously seen from the presented process. The results are given also have a fairly good level of accuracy so that they can provide precise and accurate output. The application of the data mining clustering process can group data to form a classification pattern. The MRL analysis developed plays an important role in measuring the accuracy of the variables as well as the correlation with the output which is quite significant. ANN learning on the classification analysis model can be applied to train and test networks using backpropagation algorithms. The results given are quite precise in providing output in determining air quality. In this case, based on the output obtained, it is explained that the research objective is to produce a much better form of the classification analysis process. The results given are used as consideration for decision making.

Acknowledgements

Thanks to Dr. Hj, Zerni Melmusi, SE, MM, AK, CA as Yayasan Perguruan Tinggi Komputer (YPTK) Padang. To the Chancellor of the University of Putra Indonesia YPTK Padang who has supported this research.

Hybrid Method Air Quality Classification Analysis Model

ORIGINALITY REPORT

5%

SIMILARITY INDEX

PRIMARY SOURCES

1	hrcak.srce.hr Internet	28 words — 1%
2	media.neliti.com Internet	20 words — 1%
3	repository.au.edu Internet	19 words — 1%
4	lppm.upiypk.ac.id Internet	14 words — < 1%
5	www.slideshare.net Internet	14 words — < 1%
6	idoc.pub Internet	9 words — < 1%
7	revistas.ufpr.br Internet	9 words — < 1%
8	www.econstor.eu Internet	9 words — < 1%
9	www.tandfonline.com Internet	9 words — < 1%
10	epdf.pub Internet	

8 words — < 1%

11 www.opf.slu.cz
Internet

8 words — < 1%

12 Syelfia Dewimarni, Rizalina Rizalina, Zefriyenni Zefriyenni. "Validitas Media Pembelajaran Statistika Berbasis Android dengan Teknik Peta Konsep untuk Meningkatkan Pemahaman Konsep Statistika", Jurnal Cendekia : Jurnal Pendidikan Matematika, 2022
Crossref

7 words — < 1%

13 Joseph Bamidele Awotunde, Roseline Oluwaseun Ogundokun, Femi E. Ayo, Gbemisola J. Ajamu et al. "Chapter 10 Social Media Acceptance and Use Among University Students for Learning Purpose Using UTAUT Model", Springer Science and Business Media LLC, 2020
Crossref

6 words — < 1%

14 T.K. Abdel-Galil, R.M. Sharkawy, M.M.A. Salama, R. Bartnikas. "Partial discharge pulse pattern recognition using an inductive inference algorithm", IEEE Transactions on Dielectrics and Electrical Insulation, 2005
Crossref

6 words — < 1%

EXCLUDE QUOTES OFF

EXCLUDE MATCHES OFF

EXCLUDE BIBLIOGRAPHY OFF