

# Implementation of K-means Clustering on SIPP-KLING Dashboard Applications

Fatona Fadilla Rohma, Iklima Ermis Ismail, Yoyok S Waluyo

Program Studi Teknik Informatika

Jurusan Teknik Informatika dan Komputer

Politeknik Negeri Jakarta

[fatonafadilla@gmail.com](mailto:fatonafadilla@gmail.com), [iklimaermis.ismail@tik.pnj.ac.id](mailto:iklimaermis.ismail@tik.pnj.ac.id), [yoyok\\_sw@yahoo.com](mailto:yoyok_sw@yahoo.com)

**Abstract** - This research focuses on grouping health house (rumah\_sehat) data into five clusters, namely Very Unhealthy, Unhealthy, Not Healthy Yet, Healthy, Very Healthy. There were 17 criteria as input parameters for K-Means calculation. These research aims to grouping 8969 houses into the clusters. The results of these clustering can help decision maker (government) to analyze which parts of the houses whose need the attention more, or which areas that lower than healthy standards. The test result shows that there were 3308 Very Healthy, 2496 Healthy, 792 Not healthy Yet, 1706 Unhealthy, and 667 Very Unhealthy houses. The accuracy of this method was found 87.05% using confusion matrix, with precision of 95.64% and 75, 81% , and recall of 83.82% and 92, 98% . Based on ROC the level of diagnostic value accuracy of 87.05% includes good clustering.

**Keywords** : k-means, healthy\_home, clustering, data mining

## I. INTRODUCTION

SIPP-KLING is an environmental health mapping profile information system for UPT Puskesmas Limo. This application aims to facilitate service, health, coaching and even assistance to the community to implement PHBS (Clean and Healthy Life Behavior) in the surrounding environment. In the current SIPP-KLING system, there were 2 categories of final results, namely Healthy and Unhealthy which was obtained from the final total value. Because it only has 2 categories, the analysis obtained from each village is not specific. So that supervision and evaluation cannot be carried out optimally because the gap in the same region is too wide, hence the houses in the same region (categoric) have different health problems and require different handling. With a lack of analysis it also has an impact on the funds

that will be given to each region. The amount of funds channeled will not reach the areas that need it most.

Based on the condition above, there is a need to improve the grouping of the SIPP-KLING datasets to get more specific groups. This research use K-Means Clustering to grouping 8969 houses into the cluster, to extracts the information and important patterns of interest in the SIPP-KLING application. K-Means is a most widely used and well studied method in data mining [1]. Refer to [2], the clustering analysis is useful to draw meaningful information or drawing interesting patterns form data sets and used in many fields like bioinformatics, pattern recognition, image processing, data mining, marketing, economics, etc., to get the hidden knowledge.

## II. K-MEANS CLUSTERING

Clustering is a process of grouping data objects into disjointed group called clusters, so that the data in the same cluster are similar and different to other cluster [3]. The main aim of clustering is to offer a combination of imilar objects. Between classification and clustering are confuse to be different, but in classification objects is assigned in predefined classes while in clustering classes is created [1].

The K-Means algorithm uses the process repeatedly to get the cluster database. It takes the desired number of initial clusters as input and produces the number of end clusters as output. If the algorithm is needed to generate cluster K then there will be an initial K and a final K. The K-Means method will randomly select the k pattern as the starting point of the centroid. The number of iterations to reach the centroid cluster will be affected by the prospective random initial centroid cluster where the position of the new centroid does not change. The K value chosen as the initial center will

be calculated using the Euclidean Distance formula, which is to find the closest distance between the centroid point and the data / object. Data that has a short distance or closest to the centroid will form a cluster [1].

In its completion, the K-Means algorithm will produce the centroid point which is the purpose of the K-Means algorithm. After the iteration of K-Means stops, each object in the dataset becomes a member of a cluster. The cluster value is determined by looking for all objects to find clusters with the closest distance to the object. The K-means algorithm will group data items in a dataset into a cluster based on the closest distance [4]. Following are the steps of K-means algorithm[5]:

INPUT : Number of desired clusters K Data Object D={d1,d2,...dn}

Step:

- randomly elevate K data objects (as initial centers) from data set D.
- Repeat;
- Calculate the distance of each data object  $d_i$  ( $1 \leq i \leq n$ ) from all k clusters  $C_j$  ( $1 \leq j \leq k$ ) and then assign data object  $d_i$  to the nearest cluster.
- For each cluster  $j$  ( $1 \leq j \leq k$ )
- Recalculate the cluster center until no change in the center of clusters.

OUTPUT : A set of K clusters

### III. DESIGN AND REALIZATION

#### A. Design

Dashboard SIPP-KLING (Profile Information System Mapping Environmental Health) is a web-based system developed by applying K-means clustering methods in data processing to become information that will be easy to learn knowledge. The actors involved in this application are the admin and cadre, the admin manages the SIPP-KLING web dashboard while the cadres are actors who collect the data through the citizen. This system serves admin to control the activities of data processing, and facilitate the admin in making reports and making decisions.

#### B. Implementation of K-Means Algorithm

This stage will explain the steps to operate the K-Means algorithm manually:

1. Determination of the number of clusters is five ( $k = 5$ ) to be made, the determination of clusters is based on conceptual observations by experts grouped into five, namely, Healthy, Unhealthy, Unhealthy, Very Healthy, Very Unhealthy. The amount of data used is 8969 SIPP-KLING data

and 17 attributes from the criteria for determining the Dinkes Health House. The criteria are taken based on interview process with dr. Tiur. The criteria are initialized by r (r1 up to r17) namely: ceiling, wall, floor, the bed room window, the family room window, ventilation, kitchen exhaust fan, lighting, toilet, safe drinking water, sewer system, rubbish management, opening the bed room window, opening the family room window, keep the house and the yard clean, throw the feces of the baby and toddler into the toilet, and throw the garbage in its place.

2. Determination of the centroid of each cluster can be seen in TABLE 1.

TABLE 1. The Centroid Value of Clusters

	Criteria																
	r	r	r	r	r	r	r	r	r	r	r	r	r	r	r	r	r
	1	2	3	4	5	6	7	8	9	1	1	1	1	1	1	1	1
C	6	6	6	3	3	6	6	6	1	1	1	7	8	8	8	8	8
1	2	2	2	1	1	2	2	2	0	0	0	5	8	8	8	8	8
C	6	6	6	3	3	6	6	6	7	7	7	7	8	8	8	8	8
2	2	2	2	1	1	2	2	2	5	5	5	5	8	8	8	8	8
C	6	6	6	3	3	6	6	6	5	5	5	5	8	8	8	8	8
3	2	2	2	1	1	2	2	2	0	0	0	0	8	8	8	8	8
C	3	3	3	3	3	3	3	3	2	2	2	2	4	4	4	4	4
4	1	1	1	1	1	1	1	1	5	5	5	5	4	4	4	4	4
C																	
5	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

3. Calculate the distance between data and the centroid.

Measuring the distance between the data and the centroid used Euclidian Distance formula (eq 1)

$$d(x, y) = |x - y| = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

Then the distance matrix will be obtained, namely C1, C2, C3, C4 and C5 as follows:

Data distance of cluster 1 is:

- $d(X1, C1) = \sqrt{(62-62)^2 + (62-62)^2 + (62-62)^2 + (0-31)^2 + (0-31)^2 + (0-62)^2 + (31-62)^2 + (31-62)^2 + (100-100)^2 + (100-100)^2 + (50-100)^2 + (50-75)^2 + (0-88)^2 + (88-88)^2 + (88-88)^2 + (88-88)^2 + (44-88)^2} = 143.153763485$

- $d(X_2, C_2) = \sqrt{(62-62)^2 + (62-62)^2 + (62-62)^2 + (31-31)^2 + (31-31)^2 + (31-62)^2 + (62-62)^2 + (62-62)^2 + (100-100)^2 + (75-100)^2 + (100-100)^2 + (75-100)^2 + (88-88)^2 + (88-88)^2 + (88-88)^2 + (88-88)^2} = 39.8246155035$
- $d(X_3, C_3) = \sqrt{(62-62)^2 + (62-62)^2 + (62-62)^2 + (0-31)^2 + (31-31)^2 + (31-62)^2 + (31-62)^2 + (31-62)^2 + (100-100)^2 + (100-100)^2 + (100-100)^2 + (25-100)^2 + (0-88)^2 + (44-88)^2 + (88-88)^2 + (88-88)^2 + (88-88)^2} = 126.585939188$
- $d(X_4, C_4) = \sqrt{(62-62)^2 + (62-62)^2 + (62-62)^2 + (31-31)^2 + (31-31)^2 + (31-62)^2 + (31-62)^2 + (62-62)^2 + (100-100)^2 + (100-100)^2 + (100-100)^2 + (50-100)^2 + (88-88)^2 + (88-88)^2 + (88-88)^2 + (88-88)^2 + (88-88)^2} = 50.4678115238$
- $d(X_5, C_5) = \sqrt{(62-62)^2 + (62-62)^2 + (31-62)^2 + (31-31)^2 + (31-31)^2 + (31-62)^2 + (31-62)^2 + (31-62)^2 + (100-100)^2 + (100-100)^2 + (50-100)^2 + (50-100)^2 + (44-88)^2 + (44-88)^2 + (88-88)^2 + (88-88)^2 + (44-88)^2} = 113.035392687$

The calculation is done to another data distance for clusters 2, 3, 4 and 5.

- Then grouping the data according to its cluster, which is the data that has the shortest distance. TABLE 2 shows the result of 1<sup>st</sup> iteration calculation.

TABLE 2. The result of 1<sup>st</sup> Iteration

Iteration 1				
C1	C2	C3	C4	C5
143.153 763485	140.953 893171	11.484819 124332	158.46135 1754931	252.23401 8324254
39.8246 155035	47.0212 71782	84.917607 125967	175.31400 4004244	303.43533 0836737
126.585 939188	133.787 144375	147.13599 1518051	168.70981 0028937	266.77518 6252395
50.4678 115238	66.4981 202742	97.066987 178958	176.13631 0850432	300.73243 9221312
113.035 392687	110.236 110236	121.04544 6010992	135.25531 4128503	244.43813 1231606

- The process returns again step 2. For the next repetition (1st repetition to completion), the new centroid is calculated by calculating the average value of the data in each cluster, like shown in TABLE 3.

TABLE 3. Average Centroid Value in 1st Repeat

	The Criteria																
	r	r	r	r	r	r	r	r	r	r	r	r	r	r	r	r	r
	1	2	3	4	5	6	7	8	9	1	1	1	1	1	1	1	1
C1	5	6	6	2	3	5	4	5	9	9	9	6	6	8	7	8	7
1	9	1	1	7	0	4	9	7	9	1	7	1	5	7	5	9	1
	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.
	8	4	7	7	5	0	5	7	0	8	5	0	7	4	9	8	9
	7	8	0	5	4	5	7	7	3	6	7	3	8	3	7	0	2
C2	3	5	5	5	9	9	5	6	6	6	8	8	7	3	5	5	5
2	0	4	1	9	7	6	0	0	6	6	6	6	4	6	0	4	1
	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.
	1	5	2	3	1	5	7	1	4	8	7	0	9	1	5	2	3
	7	5	9	8	3	4	1	6	9	9	8	9	6	7	5	9	8
C3	3	5	4	5	8	8	4	1	5	6	8	8	6	3	5	4	5
3	0	4	7	5	7	4	7	6	5	0	6	1	0	0	4	7	5
	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.
	7	5	9	8	2	7	6	2	7	3	8	2	4	7	5	9	8
	2	0	5	7	1	4	8	3	5	1	3	6	8	2	0	5	7
C4	6	3	2	4	7	7	5	2	2	3	7	4	2	2	3	2	4
4	.	2	7	5	4	4	0	9	6	3	7	5	8	6	2	7	5
	4	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.
	9	0	0	2	3	5	3	0	1	2	2	3	8	4	0	0	2
	5	0	7	3	3	6	7	9	4	5	0	8	5	9	0	7	3
C5												1	2				
	0	0	0	0	0	0	0	5	0	0	4	6	4	0	0	0	0

If the new centroid is different from the previous centroid, then the process continues the next step. However, if the new centroid is calculated the same as the previous centroid, then the clustering process is complete.

In this calculation, the processes stop at 12<sup>th</sup> iteration, like shown in TABLE 4 and TABLE 5. After getting the cluster label for each data, the average value is searched by adding up all members of each cluster and dividing the number of members.

TABLE 4. Data Grouping in the 12th Repetition

Iteration 11				
C1	C2	C3	C4	C5
131.9010 7031862 50	120.9553 2084984 30	114.7087 2330186 70	85.13399 6868447 0	123.3096 3521650 50
39.07710 5049157 2	70.46932 1117856 9	112.7033 5871371 70	109.2693 1585632 40	138.9591 6512706 30
108.2886 0321233 00	114.6481 4719265 70	133.8441 1196885 80	47.83863 9555203 3	126.9050 5074306 00
44.57767 4052077 6	70.20964 2990202 1	102.7415 7594922 60	98.55140 1356583 8	134.6284 6808694 00
97.92594 2945531 5	82.62801 1956220 2	77.29028 7574655 2	65.02538 2327006 4	97.25940 5221802 3

TABLE 5. Centroid Average Value in 12th Repetition

	r	r	r	r	r	r	r	r	r	r	r	r	r	r	r	r	r	r	r	r
	1	2	3	4	5	6	7	8	9	1	1	1	1	1	1	1	1	1	1	1
C																				
1	6	6	6	0	3	5	5	6	9	9	9	6	7	7	8	8	8	8	8	8
	0	1	1	.	0	8	5	0	9	0	7	3	8	4	6	0	6	.	.	.
	.	.	.	6	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.
	7	7	8	1	8	0	3	0	2	4	8	7	1	1	6	3	5	9	5	9
	7	6	8	3	3	4	1	3	5	9	3	7	2	4	4	9	9	9	9	9
C	6	6	6	3	3	5	5	6	9	9	4	5	7	6	8	8	8	8	8	8
2	0	1	1	0	0	9	7	1	8	6	7	6	2	9	7	6	7	.	.	.
	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.
	8	6	7	4	8	1	9	2	2	7	3	9	9	6	2	3	4	4	4	4
	3	2	8	6	6	7	1	5	2	9	0	0	0	0	7	5	0	0	0	0
C	4	5	5	2	3	4	4	5	9	9	5	3	7	7	8	8	2	.	.	.
3	8	7	9	9	0	9	0	8	1	4	6	4	4	2	5	4	.	.	.	.
	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.
	8	9	3	5	5	1	3	1	9	5	2	0	5	1	3	9	6	9	6	6
	0	4	2	0	6	9	3	7	9	3	2	5	2	2	1	8	3	3	3	3
C	5	5	5	1	2	3	2	4	9	9	7	4	1	3	8	8	7	.	.	.
4	2	8	9	4	7	5	5	4	2	1	4	7	4	6	3	3	7	.	.	.
	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.
	9	1	0	3	0	0	6	8	7	4	4	6	5	5	2	3	4	4	4	4
	1	3	9	7	7	6	7	6	1	6	2	8	2	6	0	3	3	3	3	3
C	5	5	5	2	2	4	3	5	6	6	3	3	4	7	1	2	.	.	.	.
5	4	9	8	5	9	3	9	4	7	9	5	1	9	4	9	7	8	8	8	8
	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.
	2	8	8	2	7	8	0	8	.	7	.	5	6	8	0	1	3	3	3	3
	0	7	5	8	6	4	6	2	5	1	5	3	6	8	7	6	4	4	4	4

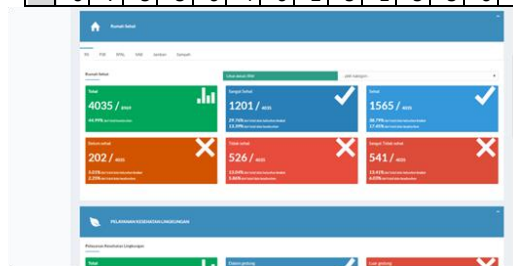


Figure 1 Main Page



Figure 2 Analysis Page

D. Testing Data Accuracy

The accuracy of the algorithm is tested using Confusion Matrix [6]. The testing aims to determine the performance of the K-Means Clustering algorithm in classifying the data into a predetermined cluster.

TABLE 6. Value of Confusion Matrix

Confusion Classification	Value
True Positive (TP)	4865
True Negative (TN)	2943
False Positive (FP)	222
False Negative (FN)	939

a. Precision

Because the centroid doesn't change (same as the previous centroid), the clustering process is complete, then the data will be grouped based on the closest distance to the cluster, from 8969 K-Means data managed to group 3308 into Very Healthy categories, 2496 Healthy categories, 792 Unhealthy categories, 1706 Unhealthy categories, and 667 Very Unhealthy categories.

C. System Implementation

Figure 1 is the implementation of the main page who displays the SIPP-KLING dashboard. On this page the system displays the total amount of data based on the status category such as whether the place is healthy or not, feasible or not, and many or few deviations. The grouped data is the result of the K-Means process. Figure 2 is an implementation of the location analysis page for each village, analysis of the location is carried out in each attribute to see which attributes most influence the village. The highest scale on the graph is the scale that most influences and becomes a warning for a village. On the page there are several buttons, the first button to send the status of the database to the database, the second button to send the distance value of the attribute to the database.

Precision is the amount of data that is true positive (the amount of positive data that is correctly recognized as positive) is divided by the amount of data that is recognized positively. From the test Assuming that the precision of the data is 95.63% for the Healthy House and 75.71% for Unhealthy House.

$$TP / (FP + TP) \times 100\% \leftrightarrow 4865 / (4865+222) = 0.956359 \text{ (Healthy House)}$$

$$TN / (FN + TN) \times 100\% \leftrightarrow 2943 / (2943 + 939) = 0.7571 \text{ (Unhealthy House)}$$

b. Recall

Recall is the amount of data which is true positive divided by the amount of data which is actually positive (true positive + true negative). For recall

value is 83.82% for Healthy House and 92.98% for Unhealthy House.

$$\text{TP} / (\text{FN} + \text{TP}) \times 100\% \leftrightarrow 4865 / (939 + 4865) = 0.8382 \text{ (Healthy House)}$$

$$\text{TN} / (\text{FP} + \text{TN}) \times 100\% \leftrightarrow 2943 / (222 + 2943) = 0.9298 \text{ (Unhealthy House)}$$

c. Accuracy

By knowing the amount of data grouped correctly it can be seen the accuracy of the prediction that is:

$$(\text{TP} + \text{TN}) / (\text{TP} + \text{TN} + \text{FP} + \text{FN}) \times 100\% \leftrightarrow (4865+2943) / (4865 + 939 + 222 + 2943) = 0.8705541309$$

Refer to [7], the accuracy value is 0.870, and include as good classification.

d. F-Measure

F-Measure is a value obtained from precision and recall measurements between clustered classes and actual classes. The ratio of f-measure is reversed, so the higher the fmeasure value, the smaller the difference between the same recall precision.

$$\text{F-Measure} = 2 / (1/\text{recall} + 1/\text{precision}) \text{ or } \text{F}(i,j) = (2 * \text{recall}(i,j) * \text{precision}(i,j)) / ((\text{precision}(i,j) + \text{recall}(i,j)))$$

For Healthy House :

$$2 / (1 / 0.9563 + 1 / 0.8382) = 0.893363789$$

$$(2 * 0.9563 * 0.8382) / (0.9563 + 0.8382) = 0.893363789$$

For Unhealthy House :

$$2 / (1 / 0.7581 + 1 / 0.9298) = 0.835216992$$

$$(2 * 0.7581 * 0.9298) / (0.7581 + 0.9298) = 0.83521$$

e. Data Analysis

From the above test results shows the accuracy of this method is 87, 05%, with precision of 95.64% for healthy class and 75, 81% for unhealthy class, and recall 83.82% for healthy and 92, 98 % to unhealthy. Results accuracy, precision and recall in the grouping of data Healthy House can be seen in the TABLE 7.

TABLE 7 Accuracy Result

Accuracy	True Healthy	True Unhealthy	Class Precision
87,05%	4865	939	95,64%
Pred. Healthy	4865	939	95,64%
Pred. Unhealthy	222	2943	75,81%
Class Recall	83,82%	92,98%	

## IV. CONCLUSION

The implementation of K-Means Clustering can be classified the 8969 data into 5 categories of Healthy House namely: Very Healthy, Healthy, Not Healthy, Unhealthy and Very Unhealthy The data obtained from the clustering results of K-means can help to analyze which parts of a house should be better addressed, or which areas have lower levels of health. The K-Means Clustering accuracy test obtained an accuracy of 87% that included in the category of good clustering.

## REFERENCES

- [1] Arpita A and Hitesh G 2013 Global K-Means (GKM) Clustering Algorithm: A Survey *International Journal of Computer Applications (IJCA)* vol 79 no 2 p 20-24
- [2] Chunfei Z and Zhiyi F 2013 An Improved K-Means Clustering Algorithm *Journal of Information & Computational Science (JICS)* 10: 1 p 193-199
- [3] Unnati R Raval and Chaita J 2016 Implementing & Improvisation of K-means Clustering Algorithm *International Journal of Computer Science and Mobile Computing (IJCSMC)* vol 5 issue 5 p 191-203
- [4] Bhoomi B, Prof. Nirali M, and Prof. Vimal P 2013 A survey on Efficient Enhanced K-Means Clustering Algorithm *International Journal for Scientific Research (IJSR) & Development (IJSRD)* vol 1 issue 9 p 1756-1758
- [5] Shi N, Liu X and Guan Y 2010 Research on K-means Clustering Algorithm An Improved K-means Clustering Algorithm *Third International Symposium on Intelligent Information Technology and Security Informatics (IITSI)* p 63-67
- [6] Jasmina Dj N, Alempije V, Sinisa S. I, Zeljko P and Milica T 2017 Evaluation of Classification Models in Machine Learning *Theory and Applications of Mathematics & Computer Science* 7 p 39-46
- [7] Gorunescu F 2011 Data Mining Concepts, Model and Techniques Berlin: Springer