

Practical No. 6

Roll No.: K041	Name: Anish Sudhan Nair
Class: B.Tech CSE Cybersecurity	Batch: K2
Date of Practical: 26/02/2022	Date of Submission: 05/03/2022
Grade:	

Aim: To implement K Means Clustering Algorithm and use Visualization tab in WEKA to visualize the clusters created

Prerequisite:

- Working of K Means Clustering algorithm
- Understanding of fundamental programming constructs in C/C++/Java
- Basic features of WEKA tool

Outcome: After successful completion of this experiment students will be able to

- Understand the working of K Means Clustering algorithm through implementation and the various measures used for checking its performance
- Use Cluster tab in WEKA , apply K Means clustering algorithm and analyze the model created.
- Use the Visualization tab in WEKA and visualize the clusters created

Theory:

k-means is one of the simplest unsupervised learning algorithms that solve the well-known clustering problem. The procedure follows a simple and easy way to classify a given data set through a certain number of clusters (assume k clusters) fixed apriori. The main idea is to define k centers, one for each cluster. These centers should be placed in a cunning way because of different location causes different result. So, the better choice is to place them as much as possible far away from each other. The next step is to take each point belonging to a given data set and associate it to the nearest center. When no point is pending, the first step is completed and an early group age is done. At this point we need to re-calculate k new centroids as barycenter of the clusters resulting from the previous step. After we have these k new centroids, a new binding has to be done between the same data set points and the nearest new center. A loop has been generated. As a result of this loop we may notice that the k centers change their location step by step until no more changes are done or in other words centers do

not move any more. Finally, this algorithm aims at minimizing an objective function known as squared error function given by:

$$J(V) = \sum_{i=1}^c \sum_{j=1}^{c_i} (\|x_i - v_j\|)^2$$

where,

' $\|x_i - v_j\|$ ' is the Euclidean distance between x_i and v_j .

' c_i ' is the number of data points in i^{th} cluster

' c ' is the number of cluster centers.

Algorithmic steps for k-means clustering

Let $X = \{x_1, x_2, x_3, \dots, x_n\}$ be the set of data points and $V = \{v_1, v_2, \dots, v_c\}$ be the set of centers.

- 1) Randomly select ' c ' cluster centers.
- 2) Calculate the distance between each data point and cluster centers.
- 3) Assign the data point to the cluster center whose distance from the cluster center is minimum of all the cluster centers.
- 4) Recalculate the new cluster center using:

$$v_i = (1 / c_i) \sum_{j=1}^{c_i} x_i$$

where, ' c_i ' represents the number of data points in i^{th} cluster.

- 5) Recalculate the distance between each data point and new obtained cluster centers.
- 6) If no data point was reassigned then stop, otherwise repeat from step 3).

Advantages

- 1) Fast, robust and easier to understand.
- 2) Relatively efficient: $O(tknd)$, where n is # objects, k is # clusters, d is # dimension of each object, and t is # iterations. Normally, $k, t, d \ll n$.
- 3) Gives best result when data set are distinct or well separated from each other.

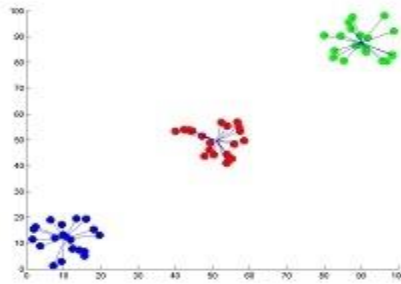


Fig I: Showing the result of k-means for ' N ' = 60 and ' c ' = 3

Disadvantages

1. The learning algorithm requires apriori specification of the number of cluster centers.
2. The use of Exclusive Assignment - If there are two highly overlapping data then k-means will not be able to resolve that there are two clusters.
3. The learning algorithm is not invariant to non-linear transformations i.e., with different representation of data we get different results (data represented in form of cartesian co-ordinates and polar co-ordinates will give different results).
4. Euclidean distance measures can unequally weight underlying factors.
5. The learning algorithm provides the local optima of the squared error function.
6. Randomly choosing of the cluster center cannot lead us to the fruitful result. Pl. refer [Fig.](#)
7. Applicable only when mean is defined i.e. fails for categorical data.
8. Unable to handle noisy data and outliers.
9. Algorithm fails for non-linear data set.

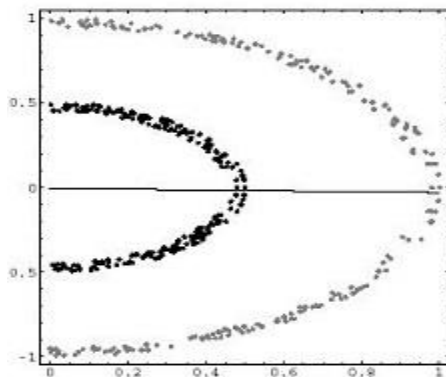


Fig II: Showing the non-linear data set where k-means algorithm fails

(TO BE COMPLETED BY STUDENTS)

Roll No.: K041	Name: Anish Sudhan Nair
Class: B.Tech CSE Cybsersecurity	Batch: K2
Date of Practical: 26/02/2022	Date of Submission: 05/03/2022
Grade:	

- a. Implement a code to create clusters using K Means Clustering algorithm for the data set indicating the marks scored by 7 students in two subjects A and B. Create 2 clusters.

	A	B
1	1.0	1.0
2	1.5	2.0
3	3.0	4.0
4	5.0	7.0
5	3.5	5.0
6	4.5	5.0
7	3.5	4.5

Code:

```
1  #include <stdio.h>
2  #include <math.h>
3  #include <stdbool.h>
4
5  float A[7]={1.0,1.5,3.0,5.0,3.5,4.5,3.5};
6  float B[7]={1.0,2.0,4.0,7.0,5.0,5.0,4.5};
7  int row_ID[7]={1,2,3,4,5,6,7};
8  float distances1[7],distances2[7];
9  float m1[2]={1.0,1.0};
10 float m2[2]={1.5,2.0};
11 float oldM1[2], oldM2[2];
12 int cluster1[7],cluster2[7];
13
14 void displayMean(float m[2], int j){
15     printf("\nMean Value %d\n", j);
16     for (int i = 0; i < 2; i++) {
17         printf("%f\t",m[i] );
18     }
19     printf("\n");
20 }
21
22 void displayCluster(int cluster[]){
23
24     for (int i = 0; i < 7; i++) {
```

```

25     if (cluster[i]==1) {
26         printf("%d\t",row_ID[i] );
27     }
28 }
29 printf("\n");
30 }
31
32 float distanceCalculator(float m[2], int i){
33     float distance=0.0;
34     distance=fabsf(m[0]-A[i])+fabsf(m[1]-B[i]);
35     return distance;
36 }
37
38 void meanDistanceCalculator(float m1[2], float m2[2]){
39
40     printf("\nDistance Matrix\n\n");
41
42     for (int i = 0; i < 10; i++) {
43         distances1[i]=distanceCalculator(m1,i);
44         distances2[i]=distanceCalculator(m2,i);
45         printf("%f\t%f\n",distances1[i],distances2[i]);
46     }
47 }
48
49 void clusterValues(){
50
51     for (int i = 0; i < 7; i++) {
52         if (distances1[i]<distances2[i]) {
53             cluster1[i]=1;
54         }
55         else {
56             cluster2[i]=1;
57         }
58     }
59
60     printf("\nCluster 1: \n");
61     displayCluster(cluster1);
62     printf("\nCluster 2: \n");
63     displayCluster(cluster2);
64
65 }
66
67 void preserveMeanValues(float m[2], float oldM[2]) {
68     for (int i = 0; i < 3; i++) {
69         oldM1[i]=m1[i];
70         oldM2[i]=m2[i];
71     }
72 }
73
74 bool meanEqual(){
75
76     int count=0;

```

```

77     for (int i = 0; i < 2; i++) {
78         if ((m1[i]==oldM1[i])&&(m2[i]==oldM2[i])) {
79             count++;
80         }
81         if (count==2) {
82             return true;
83         }
84     }
85     return false;
86 }
87
88 void newMeanValue(int cluster[7], float m[]){
89
90     float sumA=0, sumB=0;
91     int number=0;
92
93     for (int i = 0; i < 7; i++) {
94         if (cluster[i]==1) {
95             sumA+=A[i];
96             sumB+=B[i];
97             number++;
98         }
99     }
100     sumA/=number;
101     sumB/=number;
102
103     m[0]=sumA;
104     m[1]=sumB;
105 }
106
107
108 int main(){
109
110     do {
111         displayMean(m1,1);
112         displayMean(m2,2);
113         meanDistanceCalculator(m1,m2);
114         clusterValues();
115         preserveMeanValues(m1, oldM1);
116         preserveMeanValues(m2, oldM2);
117         newMeanValue(cluster1, m1);
118         newMeanValue(cluster2, m2);
119     } while(!(meanEqual()));
120
121     displayMean(m1,1);
122     displayMean(m2,2);
123
124     printf("Hence we verified \n");
125
126     return 0;
127 }
128

```

Output:

```
Lab6 — -zsh — 80x60
[(base) anish@PotatoBook lab6 % clang k_means.c -o kmeans]
[(base) anish@PotatoBook lab6 % ./kmeans]

Mean Value 1
1.000000      1.000000

Mean Value 2
1.500000      2.000000

Distance Matrix
0.000000      1.500000
1.500000      0.000000
5.000000      3.500000
10.000000     8.500000
6.500000      5.000000
7.500000      6.000000
6.000000      4.500000
1.000000      2.500000
2.000000      2.500000
4.000000      4.500000

Cluster 1:
1      2

Cluster 2:
3      4      5      6      7

Mean Value 1
1.250000      1.500000

Mean Value 2
3.900000      5.100000

Distance Matrix
0.750000      7.000000
0.750000      5.500000
4.250000      2.000000
9.250000      3.000000
5.750000      0.500000
6.750000      0.700000
5.250000      1.000000
1.750000      8.000000
2.250000      7.000000
4.250000      5.200000

Cluster 1:
1      2

Cluster 2:
3      4      5      6      7

Mean Value 1
1.250000      1.500000

Mean Value 2
3.900000      5.100000
Hence we verified
(base) anish@PotatoBook lab6 %
```

- b. Using WEKA Tool – Create a K Means clustering model in WEKA for the Iris data set (iris.arff) and analyze the clusters with the Visualization tab.

Clusterer

Choose SimpleKMeans -init 0 -max-candidates 100 -periodic-pruning 10000 -min-density 2.0 -t1 -1.25 -t2 -1.0 -N 3 -A "weka.core.EuclideanDistance" -R first-last" -I 500 -num-slots

Cluster mode

☒ Use training set
☐ Supplied test set Set...
☐ Percentage split % 66
☐ Classes to clusters evaluation (Nom) class
☒ Store clusters for visualization

Ignore attributes

Start Stop

Result list (right-click for options)

12:30:06 - SimpleKMeans

Clusterer output

=== Run information ===

Scheme: weka.clusterers.SimpleKMeans -init 0 -max-candidates 100 -periodic-pruning 10000 -min-density 2.0
Relation: iris
Instances: 150
Attributes: 5
sepalwidth
sepalwidth
petalwidth
petalwidth
class

Test mode: evaluate on training data

=== Clustering model (full training set) ===

kMeans

Number of iterations: 3
Within cluster sum of squared errors: 7.817456892309574

Initial starting points (random):

Cluster 0: 6.1,2.9,4.7,1.4,Iris-versicolor
Cluster 1: 6.2,2.9,4.3,1.3,Iris-versicolor
Cluster 2: 6.9,3.1,5.1,2.3,Iris-virginica

Missing values globally replaced with mean/mode

Final cluster centroids:

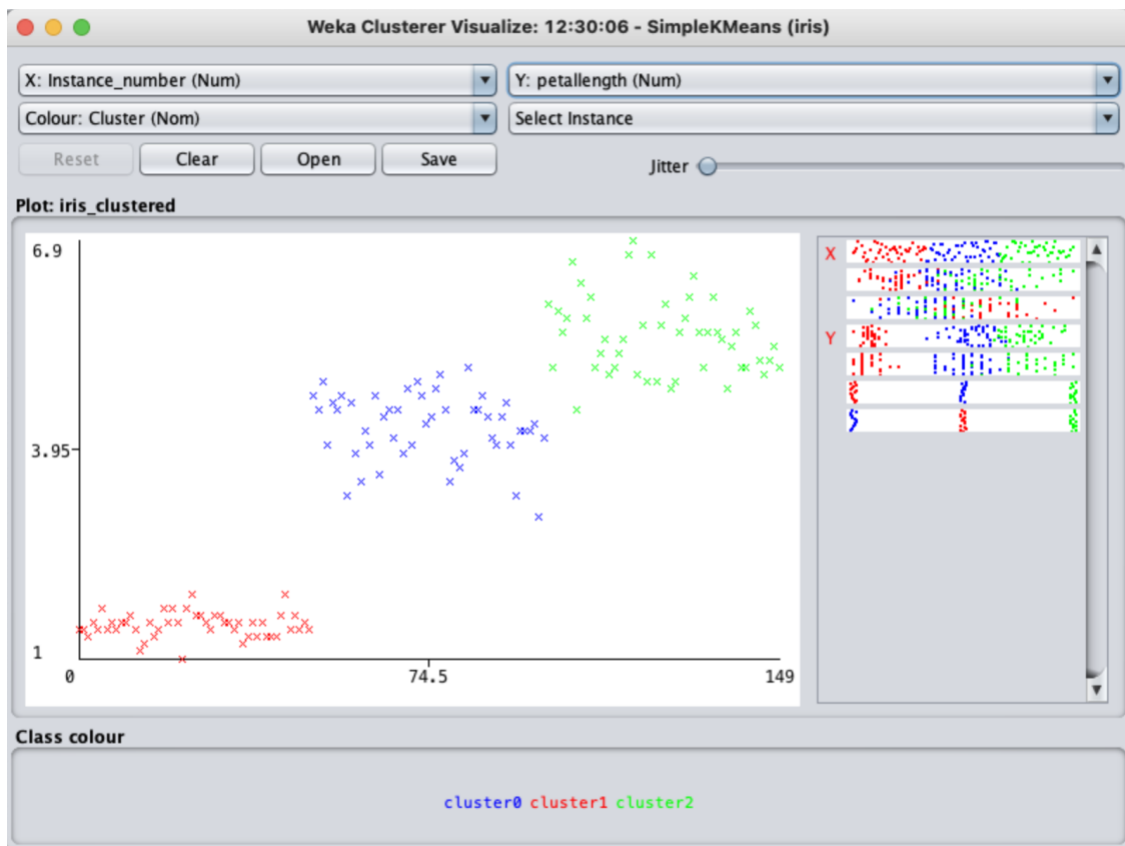
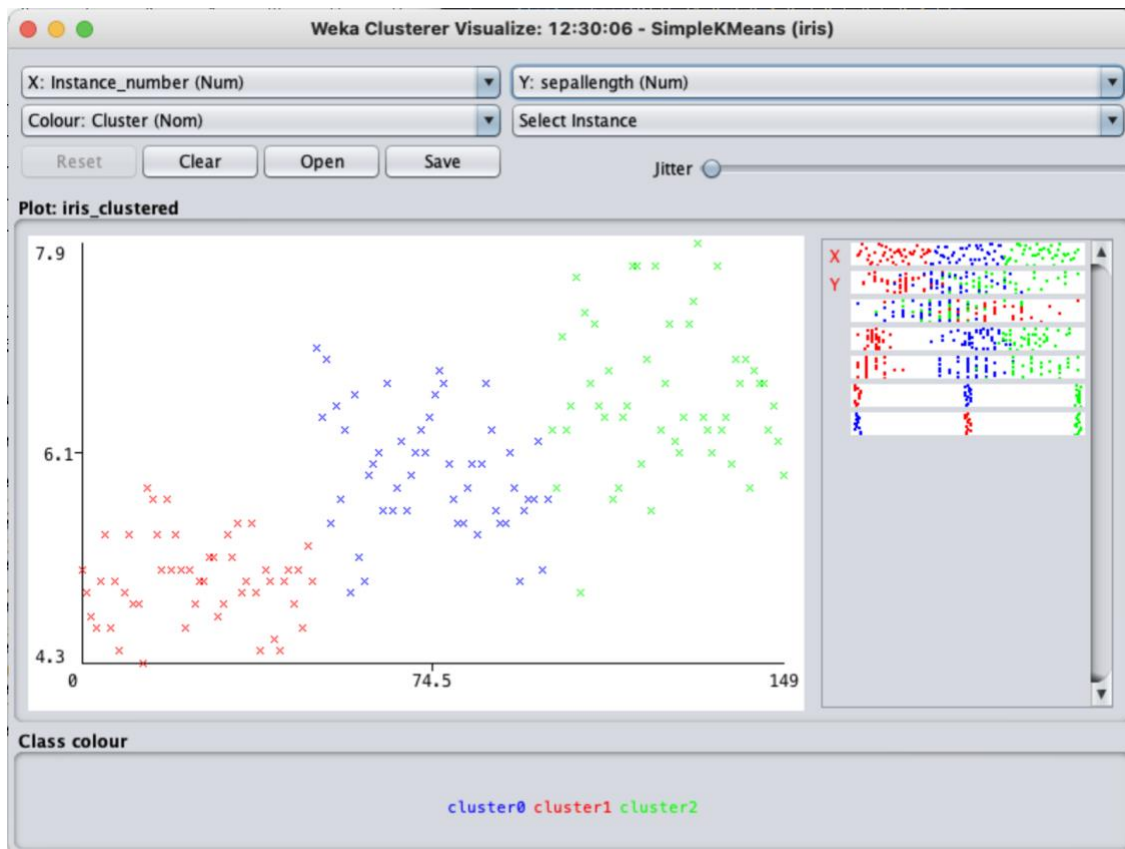
Attribute	Full Data (150.0)	Cluster# 0 (50.0)	1 (50.0)	2 (50.0)
sepalwidth	5.8433	5.936	5.006	6.588
sepalwidth	3.054	2.77	3.418	2.974
petalwidth	3.7587	4.26	1.464	5.552
petalwidth	1.1987	1.326	0.244	2.026
class	Iris-setosa Iris-versicolor	Iris-setosa	Iris-virginica	

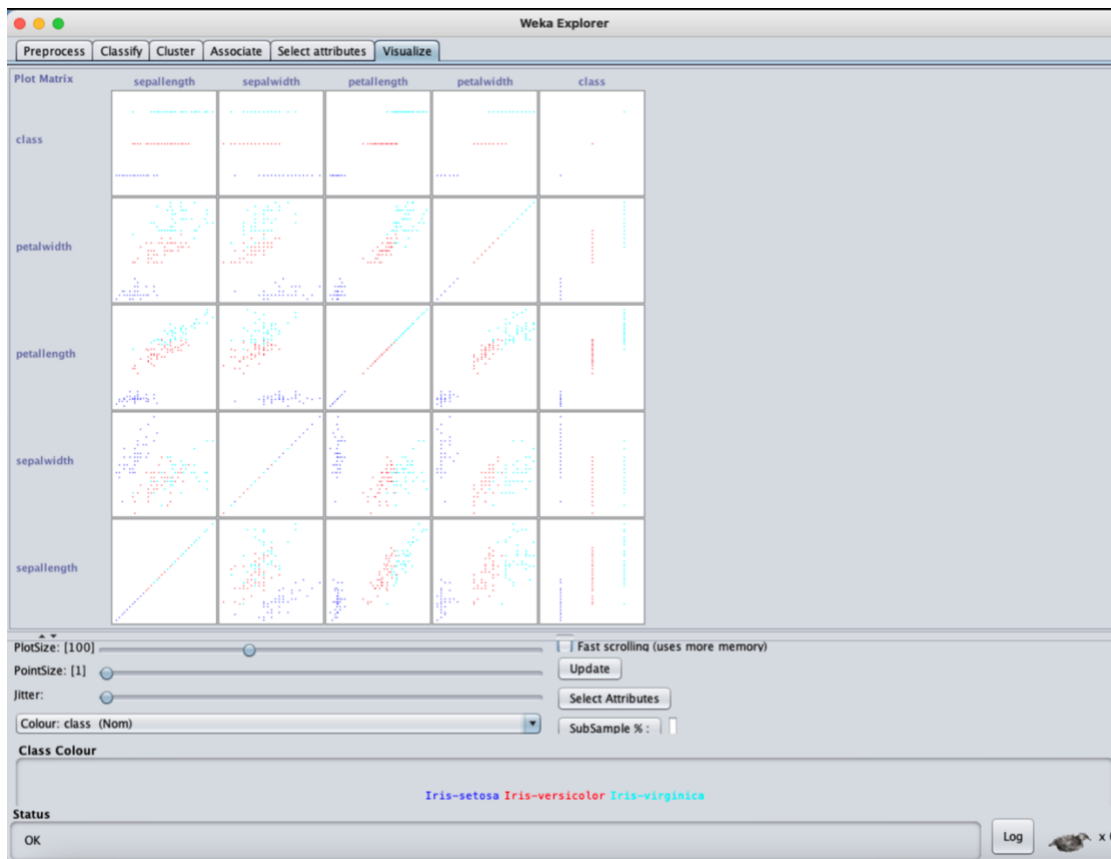
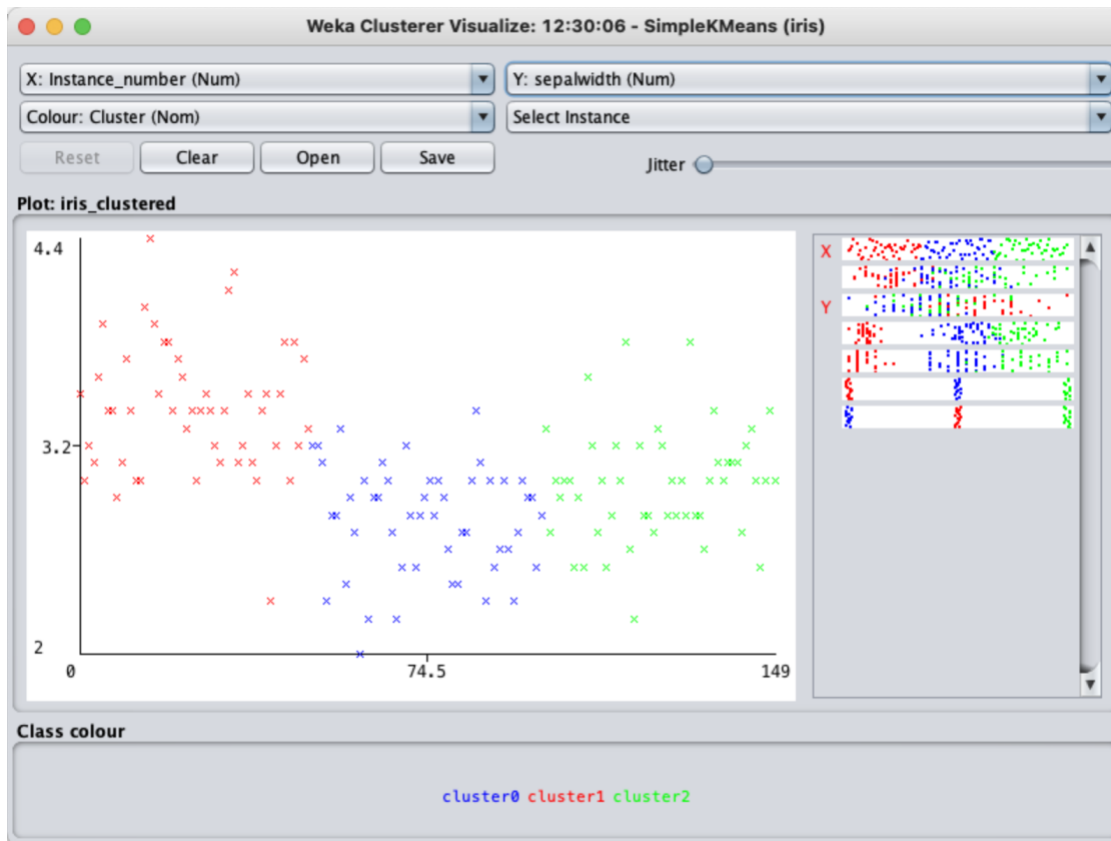
Time taken to build model (full training data) : 0.02 seconds

=== Model and evaluation on training set ===

Clustered Instances

0 50 (33%)
1 50 (33%)
2 50 (33%)





Questions to be answered:

- a. What is the benefit of using a K Medoid approach than a K Means method of cluster creation?
- ➔ The primary advantage K medoid clustering (such as Partitioning Around Medoids) has over K Means clustering is the superior ability to handle outliers. Medoids are essentially a more robust estimate of the representative point than the mean since it's an actual datapoint closer to the actual points over a virtual average point, and in k means, the mean will be centered with respect to the outliers (towards them) while in K medoids one of the more centered clustered members would be chosen. Hence, the pairwise distance measure is able to interpret outliers better.
- b. Explain any one method of handling outliers in Clustering.
- ➔ The primary aspect of outlier handling is outlier identification, post which it can be removed, changed to another value, marked as irregular and removed from cluster etc.
- ➔ One of the methods of outlier handling (identification) is the Silhouette Analysis. This method essentially every data point's similarity to its own cluster as opposed to the other clusters and assigns a value between 1 to -1.
- ➔ Negative value indicates wrong cluster or outlier, 0 indicates that the data point is in the boundary decision of 2 neighbouring clusters (could belong to either) and a positive value indicates correct clustering (distance of sample from other clusters)
- ➔ Mathematically, the silhouette coefficient can be found by first finding the average distance of the chosen point with all the other points in the same cluster (a(i)), then average of the point with all the points in the closest neighbouring cluster (b(i)), and then inserting them into the formula below:

$$s(i) = \frac{b(i) - a(i)}{\max(b(i), a(i))}$$
