

LAB Manual
PART A
(PART A: TO BE REFERRED BY STUDENTS)

Experiment No.07

A.1 Aim:

Implementation of K-means clustering using JAVA or WEKA.

A.2 Prerequisite:

Familiarity with the WEKA tool and programming languages.

A.3 Outcome:

After successful completion of this experiment students will be able to

Use classification and clustering algorithms of data mining.

A.4 Theory:

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.

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:

$$\mathbf{v}_i = (1 / c_i) \sum_{j=1}^{c_i} \mathbf{x}_i$$

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).

PART B

(PART B: TO BE COMPLETED BY STUDENTS)

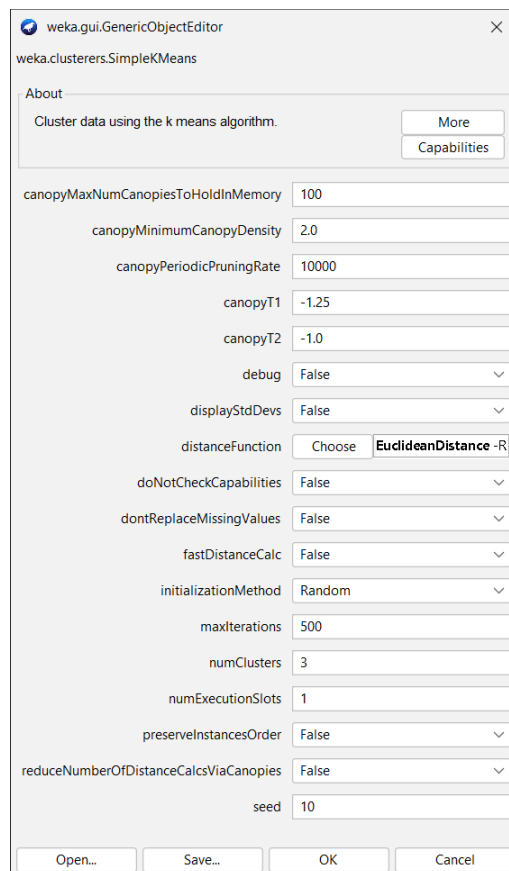
(Students must submit the soft copy as per following segments within two hours of the practical. The soft copy must be uploaded on the Blackboard or emailed to the concerned lab in charge faculties at the end of the practical in case there is no Black board access available)

Roll. No. A12	Name: SUFIYAN KHAN
Class: TE-A (AI&DS)	Batch: A1
Date of Experiment:	Date of Submission:
Grade:	

B.1 Software Code written by student:

We are using the WEKA tool and iris.arff dataset to demonstrate the K-Means clustering.

After selecting the k-means clusterer we set the number of clusters to 3.



The centroid of each cluster is shown in the result window, along with statistics on the number and percent of instances allocated to each cluster. Each cluster centroid is represented by a mean vector. This cluster can be used to describe a cluster.

Weka Explorer

Preprocess Classify **Cluster** Associate Select attributes Visualize

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 1

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)

- 15:00:00 - SimpleKMeans
- 15:01:58 - HierarchicalClusterer
- 15:29:54 - SimpleKMeans
- 15:30:26 - SimpleKMeans
- 15:34:42 - SimpleKMeans**

Clusterer output

kMeans

=====

Number of iterations: 6

Within cluster sum of squared errors: 6.998114004826762

Initial starting points (random):

Cluster 0: 6.1,2.9,4.7,1.4

Cluster 1: 6.2,2.9,4.3,1.3

Cluster 2: 6.9,3.1,5.1,2.3

Missing values globally replaced with mean/mode

Final cluster centroids:

Attribute	Full Data	Cluster# 0	Cluster# 1	Cluster# 2
	(150.0)	(61.0)	(50.0)	(39.0)
sepalwidth	5.8433	5.8885	5.006	6.8462
sepalwidth	3.054	2.7377	3.418	3.0821
petalwidth	3.7587	4.3967	1.464	5.7026
petalwidth	1.1987	1.418	0.244	2.0795

Time taken to build model (full training data) : 0.01 seconds

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

Clustered Instances

Status OK Log x0

Weka Explorer

Preprocess Classify **Cluster** Associate Select attributes Visualize

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 1

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)

- 15:00:00 - SimpleKMeans
- 15:01:58 - HierarchicalClusterer
- 15:29:54 - SimpleKMeans
- 15:30:26 - SimpleKMeans
- 15:34:42 - SimpleKMeans**

Clusterer output

Attribute	Full Data	Cluster# 0	Cluster# 1	Cluster# 2
sepalwidth	3.054	2.7377	3.418	3.0821
petalwidth	3.7587	4.3967	1.464	5.7026
petalwidth	1.1987	1.418	0.244	2.0795

Time taken to build model (full training data) : 0.01 seconds

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

Clustered Instances

Cluster	Count	Percentage
0	61	(41%)
1	50	(33%)
2	39	(26%)

Class attribute: class

Classes to Clusters:

```

0 1 2 <-- assigned to cluster
0 50 0 | Iris-setosa
47 0 3 | Iris-versicolor
14 0 36 | Iris-virginica

Cluster 0 <-- Iris-versicolor
Cluster 1 <-- Iris-setosa
Cluster 2 <-- Iris-virginica

```

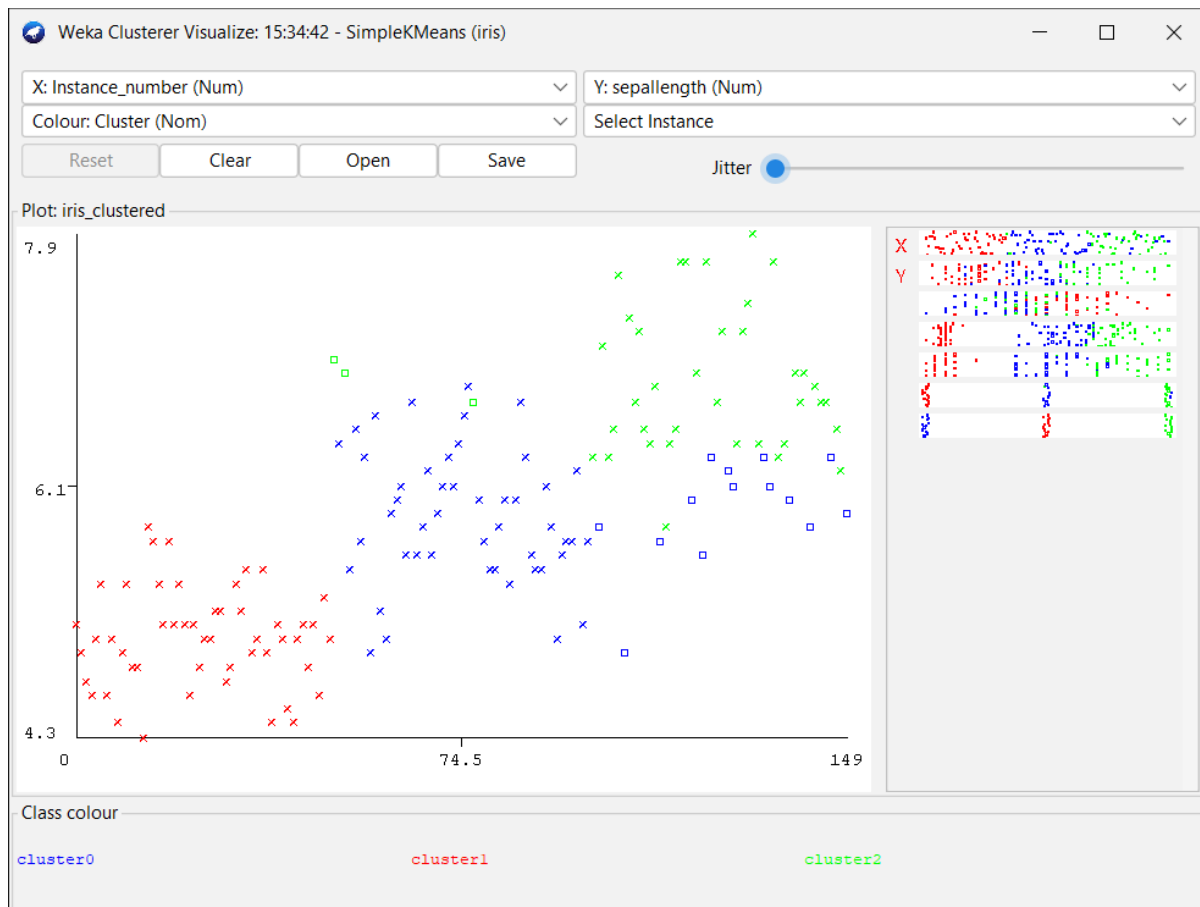
Incorrectly clustered instances : 17.0 11.3333 %

Status OK Log x0

B.2 Input and Output:

Another way to grasp the characteristics of each cluster is to visualize them.

Here is the visualization.



B.3 Observations and learning:

The k-means clustering algorithm is a popular unsupervised machine learning technique used for partitioning a dataset into K distinct, non-overlapping subgroups or clusters.

B.4 Conclusion:

K-means clustering partitions data into K clusters based on proximity to centroids, with sensitivity to initializations, a need for determining optimal K, and an assumption of spherical clusters, making it important to preprocess and interpret results judiciously.

B.5 Question of Curiosity

Q1: What is Clustering? Types of clustering? Explain advantages and disadvantages of clustering.

Clustering is an unsupervised machine learning technique used to group similar data points together in a dataset. The goal is to discover inherent patterns or structures within the data without any predefined labels.

Types of Clustering:

1. K-means Clustering
2. Hierarchical Clustering.
3. DBSCAN.
4. GMM (distribution based).

Advantages of Clustering:

- **Pattern Discovery:**

Clustering helps identify natural groupings or patterns in data that may not be immediately apparent.

- **Data Compression:**

By representing clusters with their centroids, data can be effectively compressed.

- **Anomaly Detection:**

Outliers often don't fit well into any cluster, making them easier to identify.

- **Simplifies Complex Data:**

Clustering reduces the complexity of large datasets, making them easier to analyze and interpret.

- **Facilitates Decision-Making:**

Clusters can be used as a basis for making decisions or for further analysis.

Disadvantages of Clustering:

- **Sensitivity to Initial Conditions:**

Some clustering algorithms are sensitive to the initial placement of centroids, potentially leading to different results for different initializations.

- **Determining Optimal K:**

Selecting the right number of clusters (K) can be subjective and challenging, and an inappropriate choice can lead to suboptimal results.

- **Assumption of Cluster Shape:**

Some algorithms, like K-means, assume that clusters are spherical. If clusters have complex shapes, this assumption may not hold.

- **Difficulty with Noisy Data:**

Outliers or noisy data can significantly impact clustering results, especially in methods like K-means.

- Scalability:

Some clustering algorithms may not scale well to large datasets.

- Lack of Ground Truth:

In unsupervised learning, there's no ground truth to evaluate the quality of the clustering. Evaluation metrics are often heuristic.

Q2: Give the advantages and disadvantages of K- means clustering.

Advantages of K-means Clustering:

- Efficiency:

K-means is computationally efficient and is suitable for large datasets.

- Ease of Implementation:

It's relatively simple to implement and understand, making it a good starting point for clustering analysis.

- Scalability:

It can handle a large number of data points and features.

- Interpretability:

Results are easily interpretable, as each cluster can be described by its centroid.

- Well-Suited for Spherical Clusters:

It works well when clusters are roughly spherical in shape.

Disadvantages of K-means Clustering:

- Sensitive to Initializations:

Results can be highly dependent on the initial placement of centroids, which can lead to suboptimal solutions.

- Dependent on Number of Clusters (K):

Choosing an inappropriate value of K can lead to poor clustering results.

- Assumption of Spherical Clusters:

K-means assumes that clusters are spherical and have roughly equal variance in all dimensions, which may not always be the case in real-world data.

- Vulnerability to Outliers:

Outliers can significantly affect the clustering results, as the algorithm aims to minimize the sum of squares.

- Difficulty with Non-Linearly Separable Data:

It struggles with clusters that are not well-separated or have complex shapes.

- Lack of Robustness to Noise:

Noise or irrelevant features can impact the clustering results.

Q3: How is the number of cluster chosen?

The elbow method is used to determine the number of cluster which is actually the 'k' in k-means clustering.

This method involves running the k-means clustering algorithm on the dataset for a range of k values and calculating the sum of squared distances (inertia or WCSS - Within-Cluster Sum of Squares) for each k.

The "elbow point" in the plot of WCSS against k is the point where the rate of change of WCSS starts to slow down significantly. This can indicate a good value for k.