### Exercise 5 :

## Comparison of Clustering Algorithms Using WEKA Tool:

### Aim:

The objective of this lab experiment is to compare the results of K-Means, COBWEB, Canopy, and Hierarchical clustering algorithms using the Iris dataset. The comparison will be based on the following factors: Number of clusters.,Clusters instances,Time taken to build the model.

### Algorithm:

1. Launch WEKA and Select Explorer.

2. Load the Iris dataset.

**For each clustering algorithm (K-Means, COBWEB, Canopy, Hierarchical):**

a. Select the algorithm under the "Cluster" tab.

b. Configure the algorithm parameters:

i. **K-Means:**

- numClusters: 3

- distanceFunction: EuclideanDistance

- initializationMethod: kMeansPlusPlus

- maxIterations: 100

- fastDistanceCalc: true

- seed: Any integer

ii. **COBWEB:**

- Acuity and Cutoff: Adjust if needed or leave default

iii. **Canopy:**

- T1: 1.0

- T2: 0.5

- numClusters: 3

- seed: Any integer

iv. **Hierarchical:**

- linkType: Ward

- distanceFunction: Euclidean

- numClusters: 3

c. Execute the algorithm by clicking "Start".

d. Record execution time, number of clusters, and instance distribution.

**3. Analyze Data:**

- Use external tools like Excel or Python for visualization:

a. Bar graphs for number of clusters.

b. Pie charts/stacked bar graphs for instance distribution.

c. Line graphs/bar graphs for execution times.

- Compare algorithm performance, time efficiency, and clustering effectiveness.

**4. Conclude:**

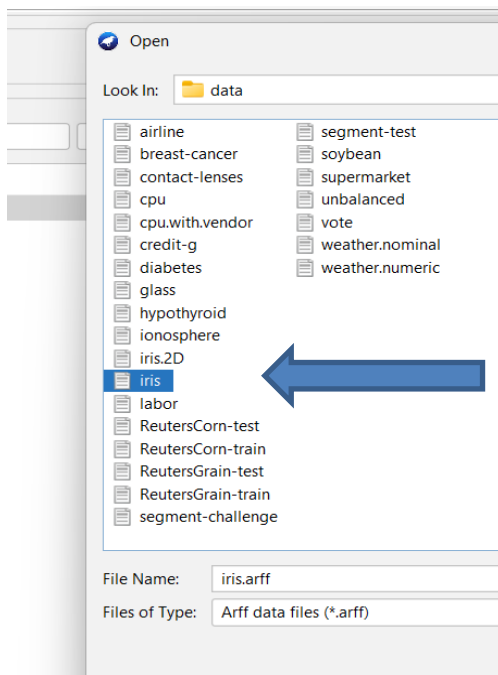- Summarize findings, compare strengths and weaknesses of each algorithm.

- Provide recommendations based on analysis., time efficiency, and clustering effectiveness.

# Implementation:

# Data Preparation:

Launch WEKA: Open the WEKA GUI Chooser by clicking on the WEKA icon.

Load Dataset: Go to the "Explorer" -> click "Open file..." and navigate to the Iris dataset. The Iris dataset is typically included in WEKA's default datasets.

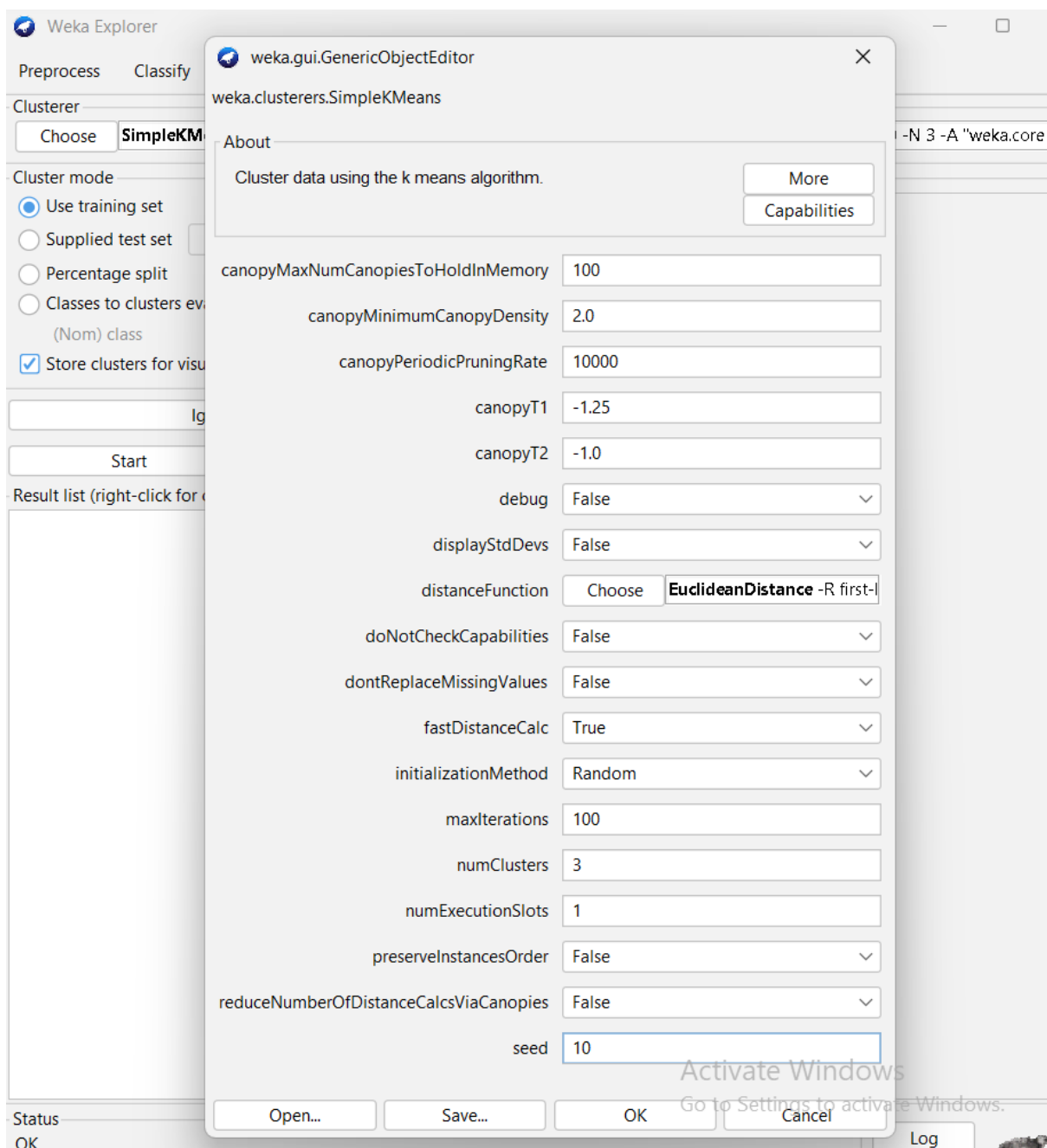## Selection of Clustering Algorithms :

a) K-Means

In the Explorer, go to the "Cluster" tab.

Click the "Choose" button, select "SimpleKMeans".

Click on the SimpleKMeans algorithm name to open its options.

Set numClusters to 3, distanceFunction to EuclideanDistance, initializationMethod to kMeansPlusPlus, maxIterations to 100, fastDistanceCalc to true, and seed to any integer.

Other parameters are set as default unless specified otherwise.

```
Clusterer output

=== Run information ===

Scheme:      weka.clusterers.SimpleKMeans -init 0 -max-candidates 100 -periodic-pruning 10000 -min-density 2.0 -t1 -1.25 -t2 -1.0 -N 3 -A "weka.core.EuclideanD:
Relation:    iris
Instances:   150
Attributes:  5
             sepallength
             sepalwidth
             petallength
             petalwidth
             class
Test mode:   evaluate on training data


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


kMeans
======

Number of iterations: 3

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:
                              Cluster#
Attribute        Full Data         0             1             2
                  (150.0)       (50.0)        (50.0)        (50.0)
=======================================================================
sepallength        5.8433        5.936         5.006         6.588
sepalwidth         3.054         2.77          3.418         2.974
```

|                              | Cluster#       |               |               |
|---------------|---------------|---------------|---------------|
| Attribute     | Full Data     | 0             | 1             | 2             |
|               | (150.0)       | (50.0)        | (50.0)        | (50.0)        |
| sepallength   | 5.8433        | 5.936         | 5.006         | 6.588         |
| sepalwidth    | 3.054         | 2.77          | 3.418         | 2.974         |
| petallength   | 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) : 0.01 seconds

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

Clustered Instances

0       50 ( 33%)
1       50 ( 33%)
2       50 ( 33%)
```
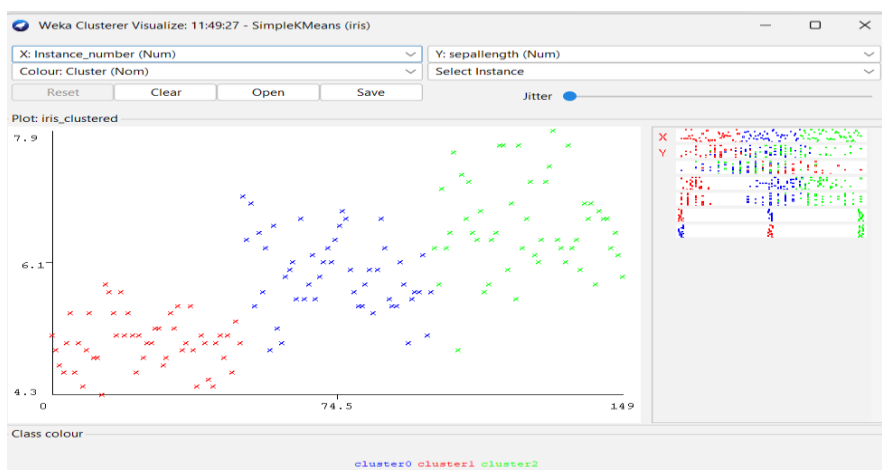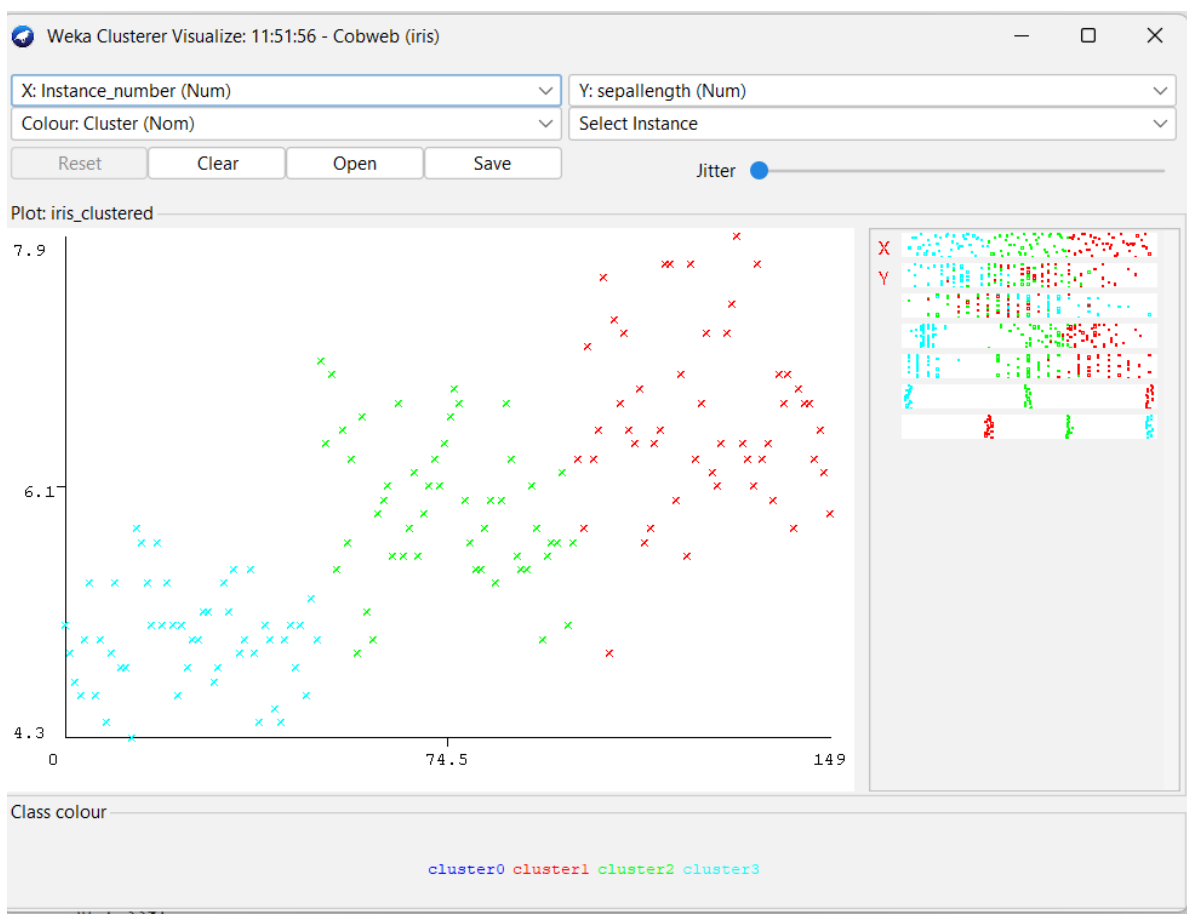
## b) COBWEB

Click "Choose", select "Cobweb".

Leave parameters as default or adjust Acuity and Cutoff as needed.

## c) Canopy

Click "Choose", select "Canopy".

Set T1 to 1.0, T2 to 0.5, seed to any integer, numClusters to 3, and other parameters as specified.



```
Clusterer output

=== Run information ===

Scheme:        weka.clusterers.Canopy -N 3 -max-candidates 100 -periodic-pruning 10000 -min-density 2.0 -t2 0.5 -t1 1.0 -S 1
Relation:      iris
Instances:     150
Attributes:    5
               sepallength
               sepalwidth
               petallength
               petalwidth
               class
Test mode:     evaluate on training data


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


Canopy clustering
=================

Number of canopies (cluster centers) found: 3
T2 radius: 0.500
T1 radius: 1.000

Cluster 0: 6.57234,2.948936,5.531915,2.029787,Iris-virginica,{47} <0>
Cluster 1: 4.96383,3.368085,1.470213,0.242553,Iris-setosa,{47} <1>
Cluster 2: 5.864444,2.722222,4.204444,1.3,Iris-versicolor,{45} <2>


Time taken to build model (full training data) : 0 seconds

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

Clustered Instances

0       50 ( 33%)
```
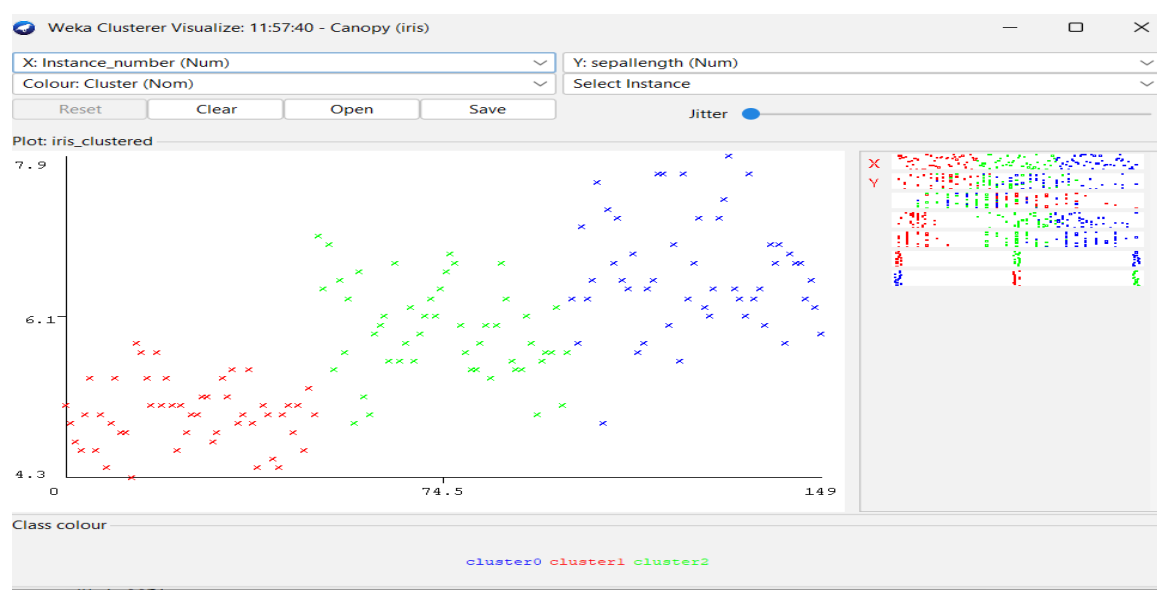
```
Clustered Instances

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

## d) Hierarchical Clustering

Click "Choose", select "HierarchicalClusterer".

Choose linkType as Ward, distanceFunction as Euclidean, numClusters to 3, and other specified settings.





## Result:

Executing the clustering algorithms (K-Means, COBWEB, Canopy, Hierarchical) on the Iris dataset in WEKA, the analysis indicates that K-Means and Hierarchical clustering effectively grouped the instances into meaningful clusters, closely reflecting the natural species division, with K-Means showing efficiency in execution time.