**Academic Year: 2025-26**                                  **Semester: V**
**Class/Branch: T.E. DS**                                   **Subject: DWMLab**

## EXPERIMENT NO. 8

1. **Aim:** To apply clustering algorithms on a given dataset using WEKA to discover natural groupings within the data.

2. **Objectives:** From this experiment, the student will be able to

   ● To implement Simple K-Means Clustering using Weka.
   ● To implement the Hierarchical Clustering using Weka.

3. **Theory:**

## Clustering using WEKA:

Clustering in WEKA is a technique used to group similar data instances based on their attributes, without using predefined class labels. It falls under unsupervised learning, where the goal is to discover hidden patterns or natural groupings within the dataset. WEKA provides an easy-to-use graphical interface where users can load datasets, select clustering algorithms, and visualize the results to interpret how data points are organized.

Simple K-Means is one of the most commonly used clustering algorithms in WEKA. It partitions the dataset into a specified number of clusters (K) by minimizing the distance between data points and their cluster centers. The algorithm iteratively adjusts centroids until the clusters become stable. It is fast and suitable for large datasets with clearly separated groups.

Hierarchical Clustering builds a hierarchy of clusters using either an agglomerative (bottom-up) or divisive (top-down) approach. It produces a dendrogram in WEKA, which shows how clusters merge or split at different similarity levels. This method does not require specifying the number of clusters beforehand and helps explore relationships among data points.

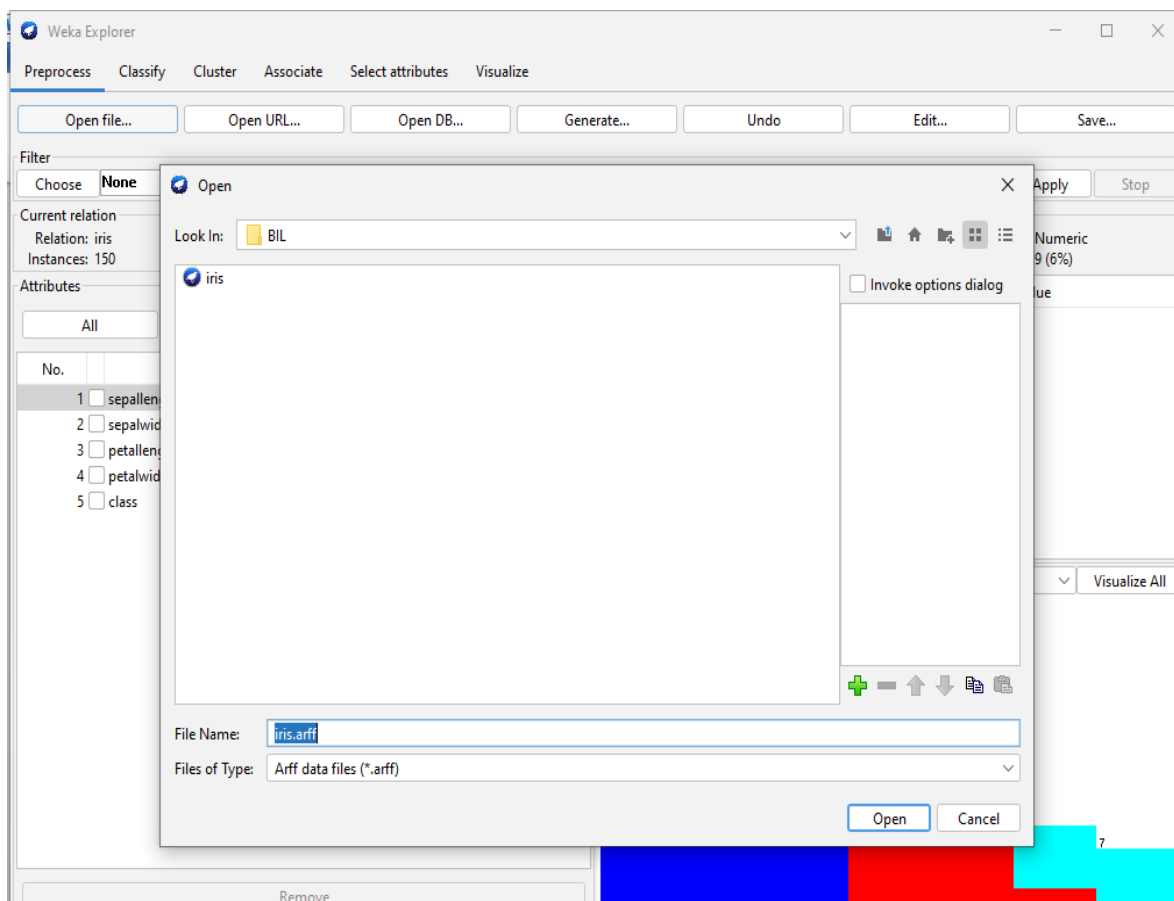## Steps for Performing Clustering in WEKA (Common for K-Means and Hierarchical)

1. Download the Iris dataset and load it into WEKA.
2. Click on the Cluster tab in the Explorer interface.
3. Click on Choose, then select either Simple K-Means (for K-Means) or HierarchicalClusterer (for Hierarchical).
4. Right-click on the selected algorithm in the choose box to open the Properties window.

.

5. Set the Number of clusters = 3 (or as required).
6. Select the -use training set option.
7. Click Start to run the clustering process.
8. To visualize the results, right-click on the algorithm name in the result list → choose View Cluster Assignments.

**The process for k-means is as shown below. The same process should be followed for Hierarchical Clustering.**

**1.        Download Iris dataset**

2.    **Click on cluster then on choose**



3.  **Right clik on simple k means in choose text box to open these properties.**
4.  **Make changes Number of clusters=3**

**5.** **Select use training set and then start**



**6.** **To visualise right click on simplekmeans in result list →view cluster assignments**



**7.** **Conclusion:**