



**Vidyavardhini's College of Engineering and Technology**  
**Department of Artificial Intelligence & Data Science**

<b>Name:</b>	
<b>Roll No:</b>	
<b>Class/Sem:</b>	TE/V
<b>Experiment No.:</b>	6
<b>Title:</b>	Using open source tools Implement Clustering Algorithms.
<b>Date of Performance:</b>	
<b>Date of Submission:</b>	
<b>Marks:</b>	
<b>Sign of Faculty:</b>	



# Vidyavardhini's College of Engineering and Technology

## Department of Artificial Intelligence & Data Science

**Aim:** To implement k-means Algorithm on large dataset using Open source tool WEKA.

**Objective:** To make students well versed with open source tool like WEKA to implement k-means algorithm.

### Theory:

- The process of grouping a set of physical or abstract objects into classes of similar objects is called clustering.
- A cluster is a collection of data objects that are similar to one another within the same cluster and are dissimilar to the objects in other clusters.
- A cluster of data objects can be treated collectively as one group and so may be considered as a form of data compression.
- Cluster analysis is an important human activity. Cluster analysis has been widely used in numerous applications including market research, pattern recognition, data analysis and image processing.
- Clustering is also called data segmentation in some applications because clustering partitions large data sets into groups according to their similarity.
- Clustering can also be used in outlier detection where outliers may be more interesting than common case.

WEKA contains "clusterers" for finding groups of similar instances in a dataset. The clustering schemes available in WEKA are k-means, Cobwebs, DBSCAN, OPTICS. Clusters can be visualized and compared to true clusters. Evaluation is based on log likelihood if clustering scheme produces a probability distribution. In 'preprocess' window click on 'open file...' button to select data file. Choosing Clustering scheme: In the 'clusterer' box click on 'choose' button. In pull-down menu select WEKA Clusteres, and select the cluster scheme 'simple K means'. Some implementations of K -means only allow numerical values for attributes ; therefore we do not need to use a filter.

Once the clustering algorithm is chosen, right click on algorithm, 'weak.gui.GenericObjectEditor' comes up to the screen. Set the value in 'numclusters' box to number of clusters required. The seed value is used in generating a random number, which is used for making the initial assignments of instances to clusters. Before we run the clustering algorithm,



# Vidyavardhini's College of Engineering and Technology

## Department of Artificial Intelligence & Data Science

we need to select 'cluster mode'. Click on 'Classes to cluster evaluation' radio-button in 'Cluster mode' box. Click the start button to run the program. When training set is complete, the 'Cluster' output area on the right panel of 'Cluster' window is filled with text describing the results of training and testing. A new entry appears in the 'Result list' box on the left of the result. Run information gives the information about : the clustering scheme used, the relation name, the number of instances, number of attributes. The clustering model shows the centroid of each cluster and statistics on the number and percentage of instances assigned to different clusters. Cluster centroid is the mean vector of each cluster so each dimension value and centroid represents mean value for that dimension in the cluster. Thus centroids can be used to characterize the cluster.

Another way of representation of results of clustering is through visualization. Right click on the entry in the 'Result list' and select ' Visualize cluster assignments' in the pull-down window. This brings up Weka clusterer visualize window. This window displays clusters in different colors for better visibility.

### Output:

**Step 1: Open Weka, the following GUI should appear on your screen**

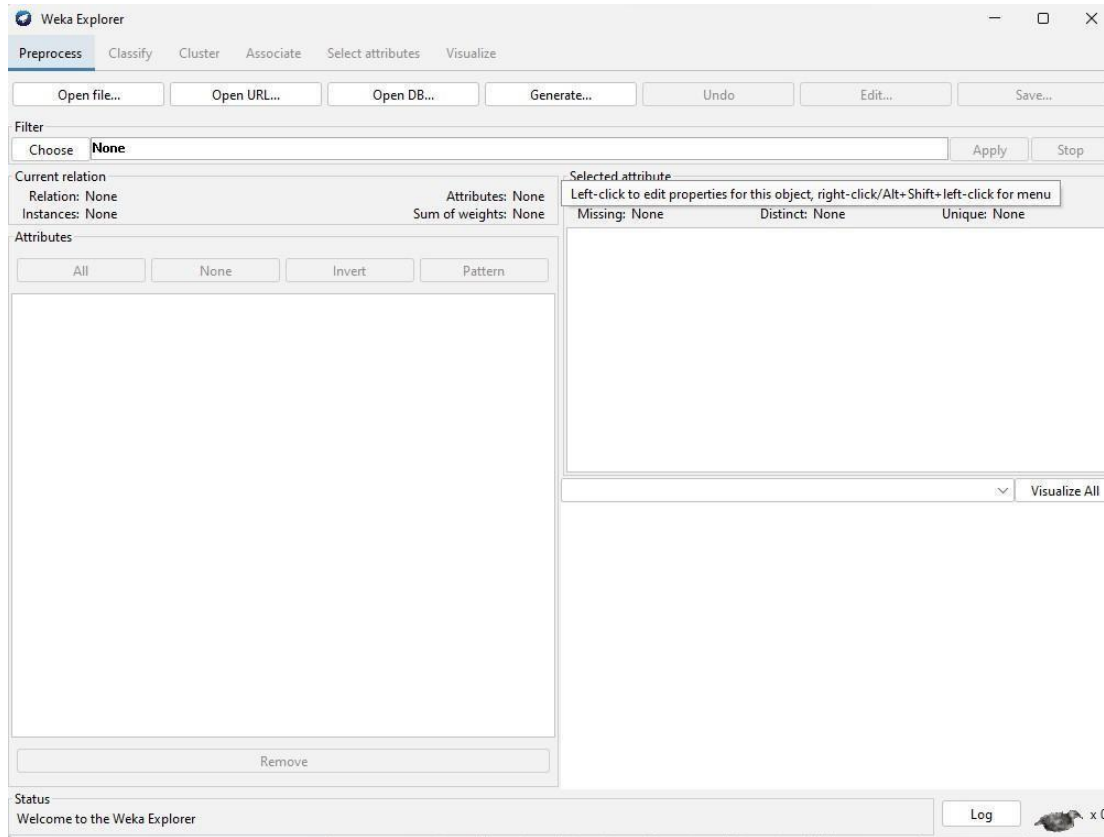




# Vidyavardhini's College of Engineering and Technology

## Department of Artificial Intelligence & Data Science

### Step 2: Click on the Explorer





# Vidyavardhini's College of Engineering and Technology

## Department of Artificial Intelligence & Data Science

**3: Cluster tab after selecting the dataset. We have selected weather.nominal here.**

Weka Explorer

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

Open file... | Open URL... | Open DB... | Generate... | Undo | Edit... | Save...

Filter: Choose: None | Apply | Stop

Current relation: weather.symbolic | Instances: 14 | Attributes: 5 | Sum of weights: 14

Attributes: All | None | Invert | Pattern

No. | Name:

No.	Label	Count	Weight
1	sunny	5	5
2	overcast	4	4
3	rainy	5	5

Selected attribute: Name: outlook | Missing: 0 (0%) | Distinct: 3 | Type: Nominal | Unique: 0 (0%)

Class: play (Nom) | Visualize All

Status: OK | Log | x 0

**Step 4: Select the algorithm to be applied. Cluster -> SimpleKmeans.**

Weka Explorer

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

Clusterer: SimpleKmeans

Clusterer output:

Status: OK | Log | x 0

**Step 5: Click on the start, the following output will appear.**



# Vidyavardhini's College of Engineering and Technology

## Department of Artificial Intelligence & Data Science

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 2 -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) play  
☒ Store clusters for visualization

Ignore attributes  
Start Stop

Result list (right-click for options)  
12:36:51 - SimpleKMeans

Clusterer output

Within cluster sum of squared errors: 26.0

Initial starting points (random):

Cluster 0: rainy,mild,normal,FALSE,yes  
Cluster 1: overcast,cool,normal,TRUE,yes

Missing values globally replaced with mean/mode

Final cluster centroids:

Attribute	Full Data	Cluster# 0	Cluster# 1
	(14.0)	(10.0)	(4.0)

=====  
outlook sunny sunny overcast  
temperature mild mild cool  
humidity high high normal  
windy FALSE FALSE TRUE  
play yes yes yes

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

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

Clustered Instances

0	10 ( 71%)
1	4 ( 29%)

### Conclusion:

In conclusion, implementing clustering algorithms using open source tools is a practical and powerful approach for discovering patterns, grouping similar data points, and making datadriven decisions. Whether you are working on data exploration, customer segmentation, or anomaly detection, these tools provide the flexibility and resources needed for effective clustering. It's important to choose the right algorithm and fine-tune parameters based on your data and use case to obtain meaningful and actionable results.