



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 Clusterer, 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, we need to select 'cluster mode'. 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.



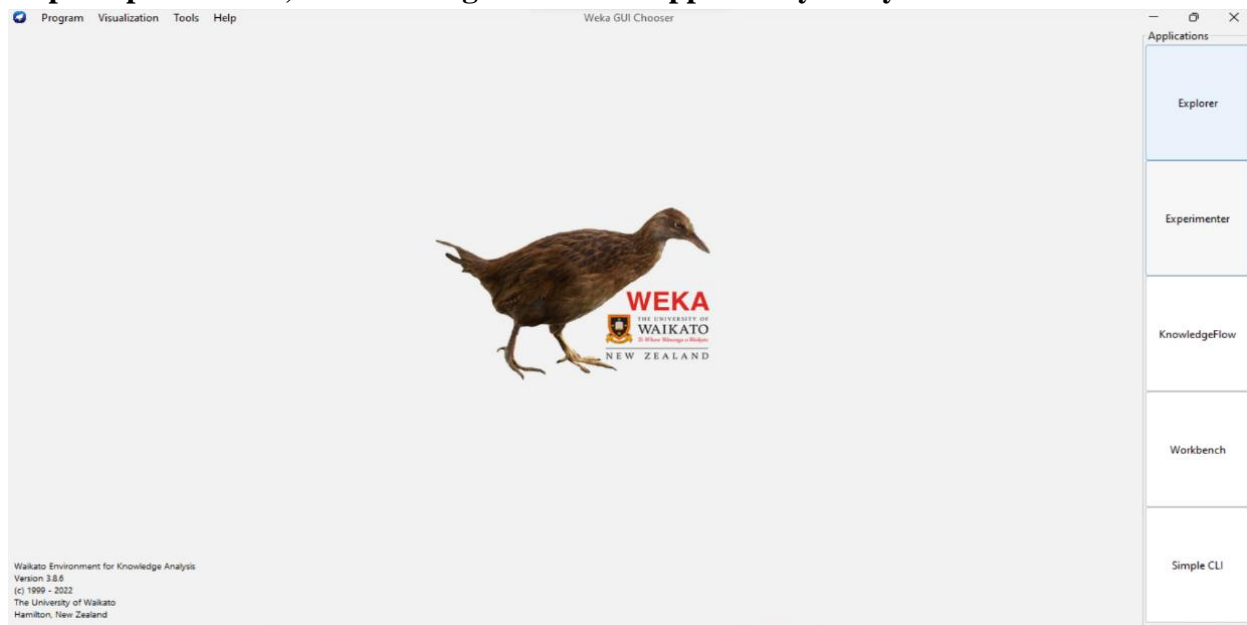
Vidyavardhini's College of Engineering and Technology

Department of Artificial Intelligence & Data Science

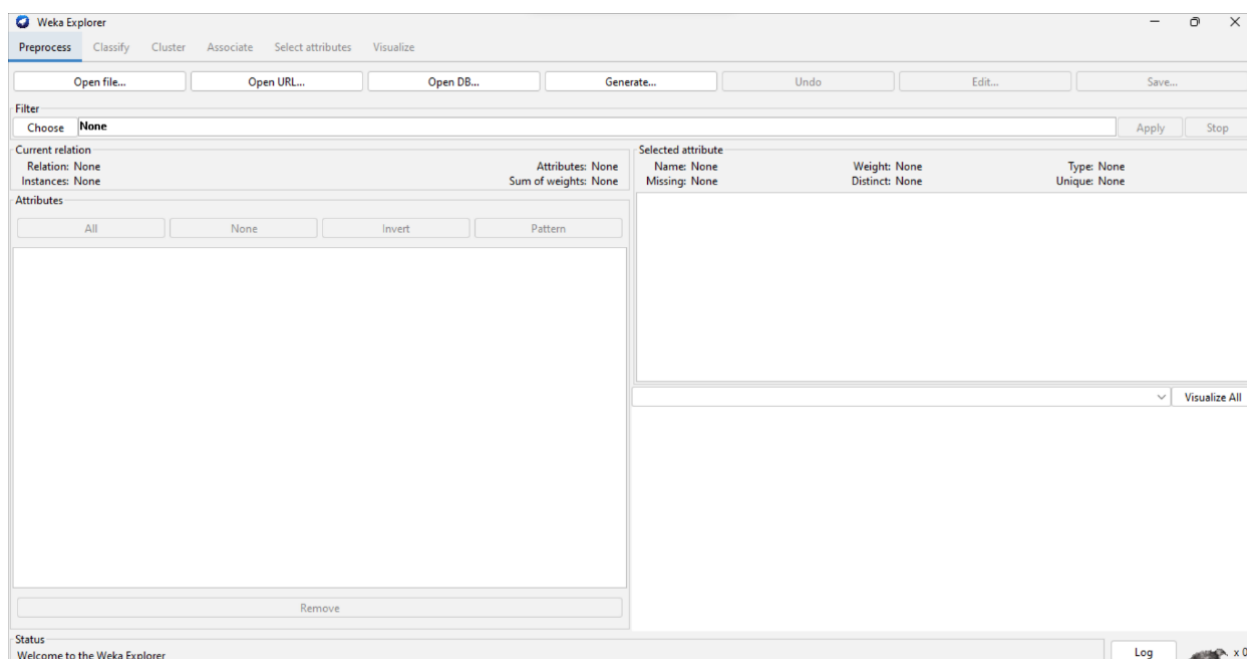
Another way of representation of results of clustering is through visualization. Right click on the entry in the 'Result list' and select ' Visualize cluster assignments' in the pull-down window. This brings up Weka clusterer visualize window. This window displays clusters in different colors for better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

Output:

Step 1:Open WEKA, the following GUI should appear on your system.



Step 2: Click on the explorer.





Vidyavardhini's College of Engineering and Technology

Department of Artificial Intelligence & Data Science

Step 3: Select 'weather.numeric.arff' dataset which already exists in the program files, from the 'Open File' section. The following screen should appear.

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
Instances: 14
Attributes: 5
Sum of weights: 14

Attributes: All None Invert Pattern

No.	Name
1	outlook
2	temperature
3	humidity
4	windy
5	play

Remove

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

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

Class: play (Nom) Visualize All

The chosen attribute will also be used as the class attribute when a filter is applied.

Status: OK Log x 0

Step 4: Select 'simpleKmeans', under cluster->Choose->simpleKmeans and click on 'Start'.

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 -S 10

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)

Clusterer output

Status: OK Log x 0



Vidyavardhini's College of Engineering and Technology

Department of Artificial Intelligence & Data Science

Step 5: Click on 'Start' , the following output will appear.

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 -S 10

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)

02:03:05 - EM

02:09:43 - SimpleKMeans

Status OK

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 2 -A "weka.core.EuclideanDistance -R first-last" -I 500 -num-slots 1 -S 10

Relation: weather

Instances: 14

Attributes: 5

outlook

temperature

humidity

windy

play

Test mode: evaluate on training data

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

KMeans

=====

Number of iterations: 3

Within cluster sum of squared errors: 16.237456311387238

Initial starting points (random):

Cluster 0: rainy,75,80,FALSE,yes

Cluster 1: overcast,64,65,TRUE,yes

Missing values globally replaced with mean/mode

Final cluster centroids:

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 -S 10

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)

02:03:05 - EM

02:09:43 - SimpleKMeans

Status OK

Clusterer output

Initial starting points (random):

Cluster 0: rainy,75,80,FALSE,yes

Cluster 1: overcast,64,65,TRUE,yes

Missing values globally replaced with mean/mode

Final cluster centroids:

Attribute	Full Data	Cluster# 0	Cluster# 1
outlook	sunny	sunny	overcast
temperature	73.5714	75.8889	69.4
humidity	81.6429	84.1111	77.2
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

Cluster	Instances	Percentage
0	9	(64%)
1	5	(36%)

Conclusion:

Thus, we have learned to implement k-means Algorithm on large dataset using Open source tool WEKA. WEKA contains "clusterers" such as k-means, Cobwebs, DBSCAN, OPTICS for finding groups of similar instances in a dataset.