| Experiment No: 08 | TE AI&DS |
|---|---|
| Date of Performance: | Roll No: 9696 |
| **Aim**: : To analyze and evaluate the performance of different Clustering techniques using WEKA tool ||
| **Related CO5:** To analyze and evaluate perform of data mining techniques applied on large dataset using open-source tool for data mining ||
| **Objective:** <br> To learn WEKA(Data Mining tool) and analyze and evaluate clustering algorithm performance. ||

**Rubrics for assessment of Experiment:**

| Sr. No | Parameters | Exceed Expectations(EE) | Meet Expectations (ME) | Below Expectations (BE) |
|---|---|---|---|---|
| 1 | Timeline (2) | Early or on time (2) | One session late (1) | More than one session late (0) |
| 2 | Preparedness (2) | Knows the basic theory related to the experiment very well. (2) | Managed to explain the theory related to the experiment. (1) | Not aware of the theory to the point. (1) |
| 3 | Effort (3) | Done expt on their own. (3) | Done expt with help from other. (2) | Just managed. (1) |
| 4 | Documentation(2) | Lab experiment is documented in proper format and maintained neatly. (2) | Documented in proper format but some formatting guidelines are missed. (1) | Experiments not written in proper format (0.5) |
| 5 | Result (1) | Specific conclusion.(1) | Partially specific conclusion. (0.5) | Not specific at all. (0) |

**Assessment Marks:**

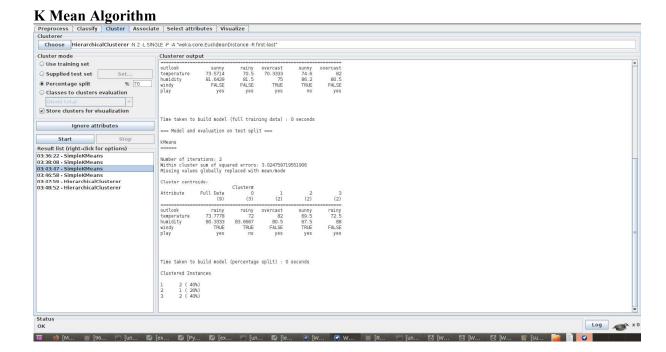| Timeline(2) | Preparedness(2) | Effort(3) | Documentation(2) | Result(1) | Total(10) |
|---|---|---|---|---|---|
|  |  |  |  |  |  |

**Theory:**

WEKA is a data mining system developed by the University of Waikato in New Zealand that implements data mining algorithms. It is a collection of machine learning algorithms for data mining tasks. The algorithms are applied directly to a dataset.

WEKA implements algorithms for data preprocessing, classification, regression, clustering, association rules; it also includes a visualization tools. When 'WEKA GUI Chooser' window appears on the screen, we can select one of the four options.
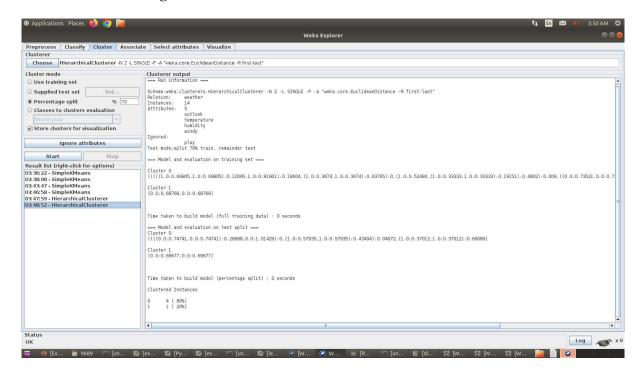
The options are 1. Simple CLI : provides a simple command line interface and allows direct execution of WEKA commands. 2. Explorer : is an environment for exploring data. 3. Experimenter : is an environment for performing experiments and conducting statistical tests between learning schemes. 4. Knowledge Flow : is a Java-Beans-based interface for setting up and running machine learning experiments. Classifiers in WEKA are the models for predicting nominal or numeric quantities. The learning schemes available in WEKA include decision trees, Bayes net, neural network, support vector machine and so on. In this experiment analysis and evaluation of three classification algorithm: Naive Bayesian algorithm, C4.5 algorithm, zero R is done. Before running the classification algorithm it is required to set test options. The selected test options were : 1. Use training set : Evaluates the classifier on how well it predicts the class of the instances it was trained on. 2. Cross-validation : Evaluates the classifier by cross-validation, using the number of folds (10) 3. Percentage split : Evaluates the classifier on how well it predicts a certain percentage of the data, which is held out for testing. In the classifier evaluation options following options are checked : 1. Output model : The output is the classification model on the full training set. 2. Output per-class stats : The precision/recall and true/false statistics for each class output 3. Output confusion matrix : The confusion matrix of one classifier's prediction is included in the output When training set is complete, the 'classifier' output area on the right panel of the classifier window is filled with text describing the results of training and testing.
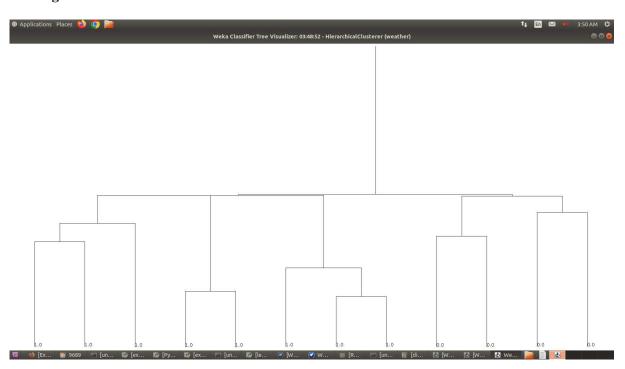
**Practical Exercise:**

Apply and evaluate the result for different classification techniques on various datasets using WEKA.

**K Mean Algorithm**

# Hierarchical Clustring:



# Dendogram:

**Post lab Questions:**

1. Apply agglomerative hierarchical clustering technique (single linkage) on following dataset.
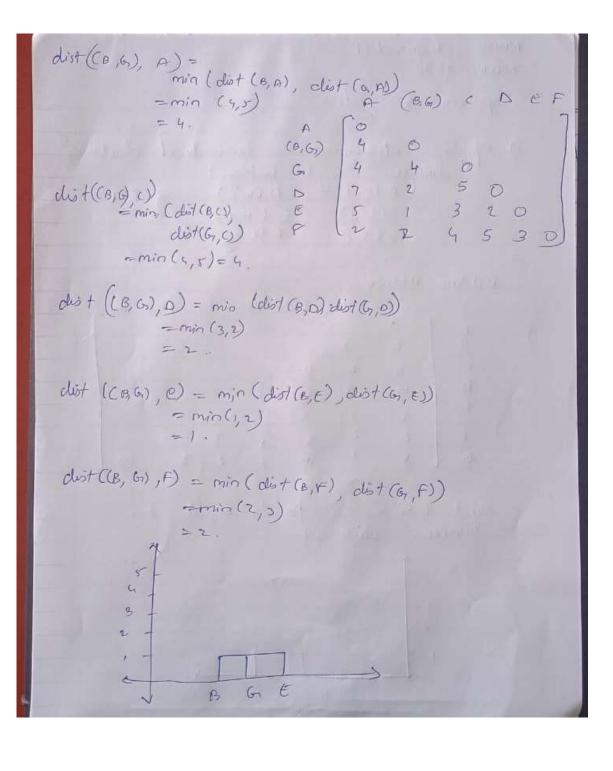
   A(1,2)
   B(3,4)
   C(5,2)
   D(4,6)
   E(4,4)
   F(2,3)
   G(3,5)

NAME= Shivprasad.cp
ROLLNO- 9696

Experiment 8   Postlab-

1) Apply agglomerative hierarchical clustering technique (single linkage)

   A(1,2)        D(4,6)        g(3,5)
   B(3,4)        E(4,4)
   C(5,2)        F(2,3)

→

Distance Matrix

|    | A | B | C | D | E | F | G |
|----|---|---|---|---|---|---|---|
| A  | 0 |   |   |   |   |   |   |
| B  | 4 | 0 |   |   |   |   |   |
| C  | 4 | 4 | 0 |   |   |   |   |
| D  | 7 | 3 | 5 | 0 |   |   |   |
| E  | 5 | 1 | 3 | 2 | 0 |   |   |
| F  | 2 | 2 | 4 | 5 | 3 | 0 |   |
| G  | 5 | 1 | 5 | 2 | 2 | 3 | 0 |

Point E & B & point (G, B) has minimum distance.



B.        G.

$$\text{dist}((B,G), A) =$$
$$\min (\text{dist}(B,A), \text{dist}(G,A))$$
$$= \min (4,5)$$
$$= 4.$$

|        | A | (B,G) | G | D | E | F |
|--------|---|-------|---|---|---|---|
| A      | 0 |       |   |   |   |   |
| (B,G)  | 4 | 0     |   |   |   |   |
| G      | 4 | 4     | 0 |   |   |   |
| D      | 7 | 2     | 5 | 0 |   |   |
| E      | 5 | 1     | 3 | 2 | 0 |   |
| F      | 2 | 2     | 4 | 5 | 3 | 0 |

$$\text{dist}((B,G), C)$$
$$= \min (\text{dist}(B,C),$$
$$\text{dist}(G,C))$$
$$= \min (4,5) = 4.$$

$$\text{dist}((B,G), D) = \min (\text{dist}(B,D), \text{dist}(G,D))$$
$$= \min (3,2)$$
$$= 2.$$

$$\text{dist}((B,G), E) = \min (\text{dist}(B,E), \text{dist}(G,E))$$
$$= \min (1,2)$$
$$= 1.$$

$$\text{dist}((B,G), F) = \min (\text{dist}(B,F), \text{dist}(G,F))$$
$$= \min (2,3)$$
$$= 2.$$

$$\text{dist}(((B,G),E),A)$$
$$= \min(\text{dist}(B,A),$$
$$\text{dist}(G,A),$$
$$\text{dist}(E,A))$$
$$= \min(4,5,5)$$
$$= 4.$$

|  | A | (B,G), E | C | D | F |
|---|---|---|---|---|---|
| A | 0 |  |  |  |  |
| ((B,G),E) | 4 | 0 |  |  |  |
| C | 4 | 3 | 0 |  |  |
| D | 7 | 2 | 5 | 0 |  |
| F | 2 | 2 | 4 | 5 | 0 |

$$\text{dist}(((B,G),E),C) = \min(\text{dist}(B,C), \text{dist}(G,C), \text{dist}(E,C))$$
$$= \min(4,5,3)$$
$$= 3.$$

$$\text{dist}(((B,G),E),D) = \min(\text{dist}(B,D), \text{dist}(G,D),$$
$$\text{dist}(E,D))$$
$$= \min(3,5,2)$$
$$= 2$$

$$\text{dist}(((B,G),E),F) = \min(\text{dist}(B,F), \text{dist}(G,F),$$
$$\text{dist}(E,F))$$
$$= \min(2,3,3)$$
$$= 2.$$

$dist ((A,F), C)$

$= min (dist (A,C), dist (F,C))$

$= min (4, 2)$

$= 2.$

$$\begin{array}{c|ccc} & (A,F) & C & ((B,G,E),D) \\ \hline (A,F) & 0 & & \\ C & 2 & 0 & \\ ((B,G,E),D) & 2 & 3 & 0 \end{array}$$

$dist (A,F), (((B,G),E),D))$

$= dist\,min (dist (A,D), dist (F,D), dist (A,B),$
$\qquad dist (A,G), dist (A,E), dist (F,B)$
$\qquad dist (F,G), dist (F,E))$

$= min (7, 5, 4, 5, 5, 2, 3, 3)$

$= 2.$

$dist (C, (((B,G),E),D)) = min (dist (C,D), dist (C,B), dist (C,G)$
$\qquad\qquad dist (C,E)).$

$= min (5, 4, 5, 3)$

$= 3.$

$$(A,F,C) \quad (B,G,E,D)$$

$$\begin{array}{c} (A,F,G) \\ ((B,G,E),D) \end{array} \begin{bmatrix} 0 & \\ 2 & 0 \end{bmatrix}$$

$dist\big((A,F,c),((B,G),E),D\big)$

$= min\big(dist(A,B), dist(A,G), dist(A,E), dist(A,D)$
$\quad dist(F,B), dst\ dist(F,G), dist(F,E),$
$\quad dist(F,D), dist(C,B), dist(G,G)$
$\quad dist(C,E), dist(C,D)\big)$
$= min(4,5, 5, 7, 2, 3, 5, 4, 5, 3, 5)$
$\approx 2$



2. Apply k-mean clustering technique on following dataset.
   1,2,6,7,8,10,15,17,20
   K=3

2) 1, 2, 6, 7, 8, 10, 15, 17, 20

K = 3

→ Let Centroids be $K_1 = 2$, $K_2 = 8$, $K_3 = 15$

Iteration 1

$K_1$  $m_1 = 2$

$\boxed{2}$

$K_2$  $m_2 = 8$

$\boxed{8}$

$K_3$  $m_3 = 15$

$\boxed{15}$

$K_1 = \{1, 2\}$  $K_2 = \{6, 7, 8, 10\}$  $K_3 = \{15, 17, 20\}$

update the Centroid

$C_1 = \dfrac{1+2}{2}$  $C_2 = \dfrac{6+7+8+10}{4}$  $C_3 = \dfrac{15+17+20}{3}$

$= 2$  $= 7.75$  $= 17.33$

$= 8$  $= 17$

∴ new centroids ($C_3$) & assumed centroid are not equal, therefore repeat the step.

Iteration 2

$K_1$  $m_1 = 2$

$\boxed{2}$

$k_2$  $m_2 = 8$

$\boxed{8}$

$K_3$  $m_3 = 17$

$\boxed{17}$

$K_1 = \{1, 2\}$ $K_2 = \{6, 7, 8, 10\}$ $K_3 = \{15, 10\}$

update the centroid

$c_1 = \dfrac{1+2}{2}$ $c_2 = \dfrac{6+7+8+10}{4}$ $c_3 = \dfrac{15+17+20}{3}$

$= 2$ $= 8$ $= 17$

∴ The centroids did not change in this iteration.

∴ The clusters are
$K_1 = \{1, 2\}$
$K_2 = \{6, 7, 8, 10\}$
$K_3 = \{15, 17, 20\}$

**Conclusion:**
**We are able to understand different clustering algorithm and am able to implement the algorithm using WEKA tool.**