

## **Experiment 02: Implementation of K-means algorithm**

### **DESCRIPTION:**

K-means algorithm aims to partition  $n$  observations into “ $k$  clusters” in which each observation belongs to the cluster with the nearest mean, serving as a prototype of the cluster. This results in partitioning of the data into Voronoi cells.

### **ILLUSTRATION:**

As a simple illustration of a k-means algorithm, consider the following data set consisting of the scores of two variables on each of the five variables.

I	X1	X2
A	1	1
B	1	0
C	0	2
D	2	4
E	3	5

This data set is to be grouped into two clusters: As a first step in finding a sensible partition, let the A & C values of the two individuals furthest apart (using the Euclidean distance measure), define the initial cluster means, giving:

Cluster	Individual	Mean Vector(Centroid)
Cluster1	A	(1,1)
Cluster2	C	(0,2)

The remaining individuals are now examined in sequence and allocated to the cluster to which they are closest, in terms of Euclidean distance to the cluster mean. The mean vector is recalculated each time a new member is added. This leads to the following series of steps:

	A	C
A	0	1.4
B	1	2.5
C	1.4	0
D	3.2	2.82
E	4.5	4.2

Initial partitions have changed, and the two clusters at this stage having the following characteristics.

	Individual	Mean vector( Centroid)
Cluster 1	A,B	(1,0.5)
Cluster 2	C,D,E	(1.7,3.7)

But we cannot yet be sure that each individual has been assigned to the right cluster. So, we compare each individual's distance to its own cluster mean and to that of the opposite cluster. And, we find:

I	A	C
A	0.5	2.7
B	0.5	3.7
C	1.8	2.4
D	3.6	0.5
E	4.9	1.9

The individual C is now relocated to Cluster 1 due to its less mean distance with the centroid points. Thus, it is relocated to cluster 1 resulting in the new partition

	Individual	Mean vector(Centroid)
Cluster 1	A,B,C	(0.7,1)
Cluster 2	D,E	(2.5,4.5)

The iterative relocation would now continue from this new partition until no more relocation occurs. However, in this example each individual is now nearer its own cluster mean than that of the other cluster and the iteration stops, choosing the latest partitioning as the final cluster solution. Also, it is possible that the k-means algorithm won't find a final solution. In this case, it would be a better idea to consider stopping the algorithm after a pre-chosen maximum number of iterations.

#### Checking the solution in weka:

In order to check the result in the tool we need to follow a procedure.

#### **Step 1:**

Create a csv file with the above table considered in the example. the csv file will look as shown below:

Clipboard		Font		
A1		$f_x$	i	
	A	B	C	D
1	i	x1	x2	
2	A	1	1	
3	B	1	0	
4	C	0	2	
5	D	2	4	
6	E	3	5	
7				

## Step 2:

Now open weka explorer and then select all the attributes in the table.

**Filter**

Choose **None** Apply

**Current relation**

Relation: menas Attributes: 3  
Instances: 5 Sum of weights: 5

**Attributes**

All None Invert Pattern

No.	Name
1	<input checked="" type="checkbox"/> i
2	<input checked="" type="checkbox"/> x1
3	<input checked="" type="checkbox"/> x2

Remove


**Selected attribute**

Name: i  
Missing: 0 (0%)  
Distinct: 5  
Type: Nominal  
Unique: 5 (100%)

No.	Label	Count	Weight
1	A	1	1.0
2	B	1	1.0
3	C	1	1.0
4	D	1	1.0

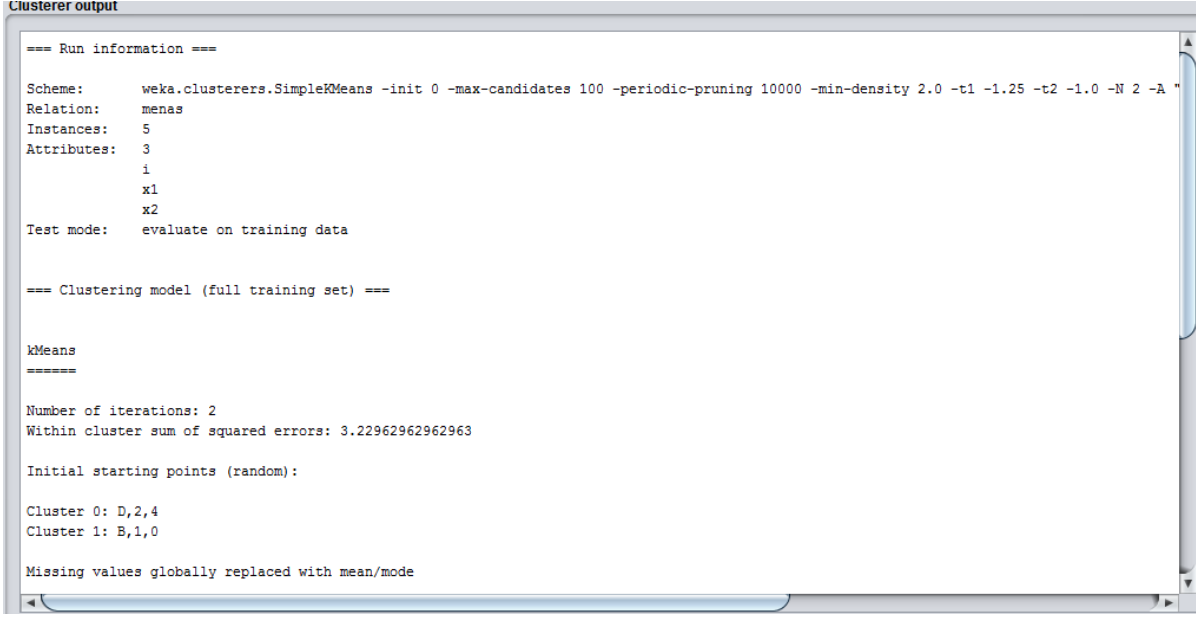
Class: x2 (Num) Visualize All

**Status**

OK Log  x 0

### **Step 3:**

Select the cluster tab in the tool and choose normal k-means technique to see the result as shown below.



```
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 "
Relation:     menas
Instances:    5
Attributes:   3
              i
              x1
              x2
Test mode:    evaluate on training data

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

kMeans
=====

Number of iterations: 2
Within cluster sum of squared errors: 3.22962962962963

Initial starting points (random):

Cluster 0: D,2,4
Cluster 1: B,1,0

Missing values globally replaced with mean/mode
```

Final cluster centroids:

Attribute	Full Data (5.0)	Cluster#	
		0 (2.0)	1 (3.0)
i	A	D	A
x1	1.4	2.5	0.6667
x2	2.4	4.5	1

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

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

Clustered Instances

0	2 ( 40%)
1	3 ( 60%)

## Exercise

1. Implement of K-means clustering using crime dataset.