

St. Francis Institute of Technology, Mumbai-400 103

**Department Of Information Technology**

A.Y. 2024-2025

Class: TE-ITA/B, Semester: VI

Subject: **Business Intelligence Lab**

**Experiment – 7: To implement K-means clustering algorithm using open source tools, WEKA & ORANGE**

1. **Aim:** : To implement K-means clustering algorithm using open source tools, WEKA & ORANGE

1. **Objectives:** After study of this experiment, the students will be able to Implement K Means

1. **Outcomes:** After study of this experiment, the students will be able to

**CO 4:** Design and Implement various clustering data mining techniques such as Partitioning methods, Hierarchical Methods, Density - Based methods along with identification and analysis of outlier.

1. **Prerequisite:** Introduction to all the three clustering algorithms & Problem solving approach.

1. **Requirements:** Personal Computer, Windows XP operating system/Windows 7, Internet Connection, Microsoft Word, WEKA tool, ORANGE tool.

1. **Theory:**
  - a. Explain K means (graph) algorithm
  - b. Explain K medoids algorithm

1. **Laboratory Exercise:** Implementation of K means clustering algorithm using WEKA & ORANGE along with screenshots.

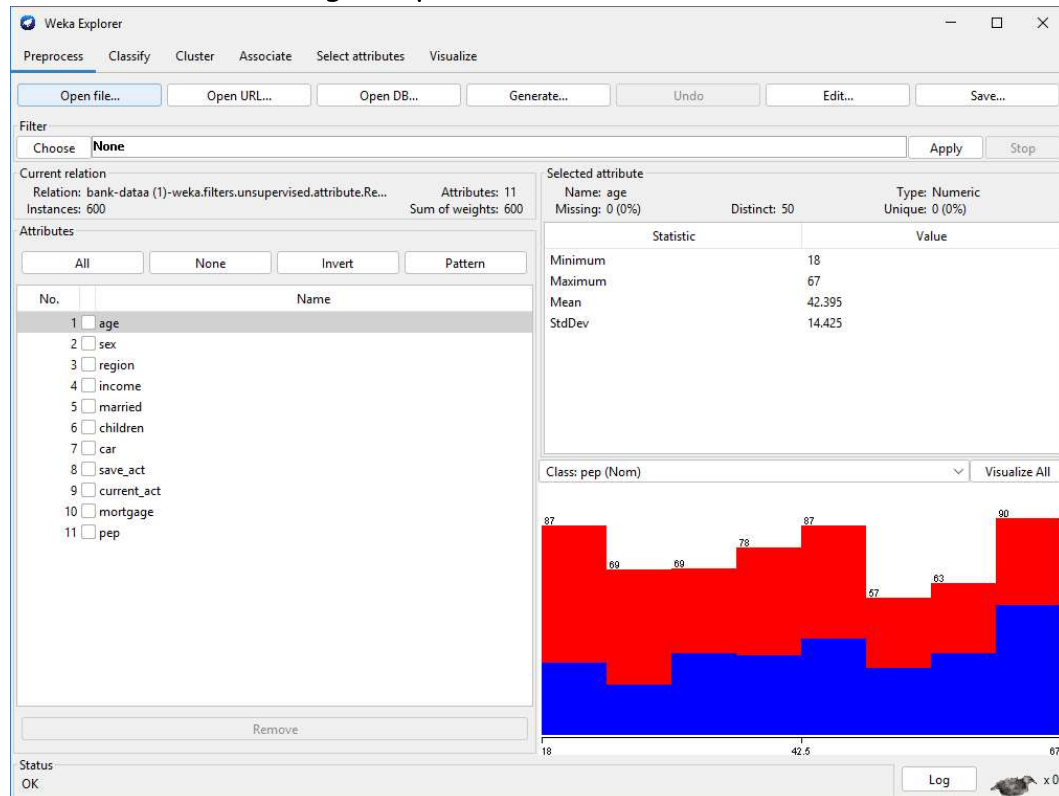
1. **Post-Experiments Exercise**

- a. **Questions:**
  - Explain advantages and disadvantages of K means
  - K means (graph) solved numerical

- a. **Conclusion:**
  - Summary of Experiment
  - Importance of Experiment
  - Application of Experiment

## K-Means Clustering in Weka

### Loading Pre-processed Dataset in WEKA:



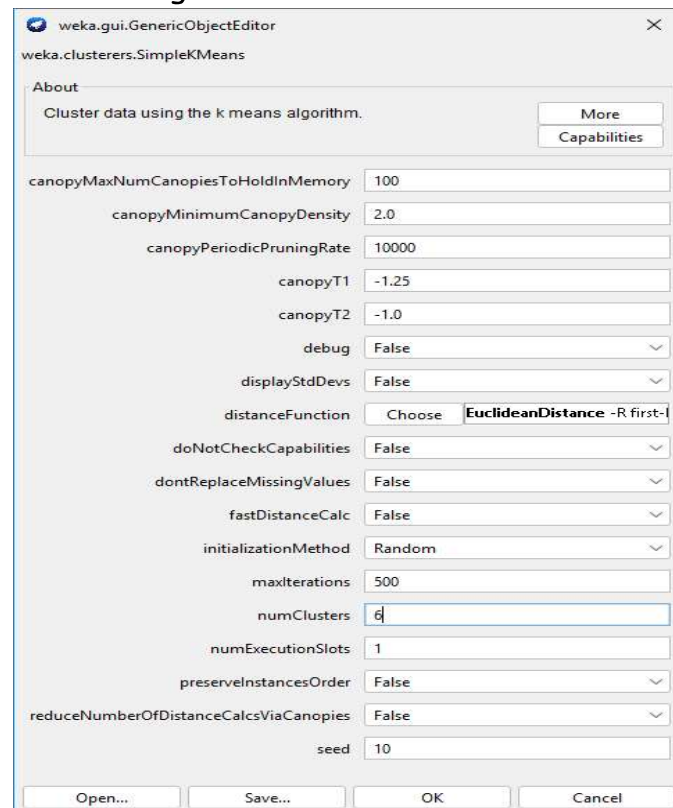
### Pre-processed Dataset:

```
@relation 'bank-dataaa (1)-
weka.filters.unsupervised.attribute.Remove-R1'

@attribute age numeric
@attribute sex {FEMALE,MALE}
@attribute region {INNER_CITY,TOWN,RURAL,SUBURBAN}
@attribute income numeric
@attribute married {NO,YES}
@attribute children numeric
@attribute car {NO,YES}
@attribute save_act {NO,YES}
@attribute current_act {NO,YES}
@attribute mortgage {NO,YES}
@attribute pep {YES,NO}

@data
48,FEMALE,INNER_CITY,17546,NO,1,NO,NO,NO,NO,YES
40,MALE,TOWN,30085.1,YES,3,YES,NO,YES,YES,NO
51,FEMALE,INNER_CITY,16575.4,YES,0,YES,YES,YES,NO,NO
23,FEMALE,TOWN,20375.4,YES,3,NO,NO,YES,NO,NO
57,FEMALE,RURAL,50576.3,YES,0,NO,YES,NO,NO,NO
57,FEMALE,TOWN,37869.6,YES,2,NO,YES,YES,NO,YES
```

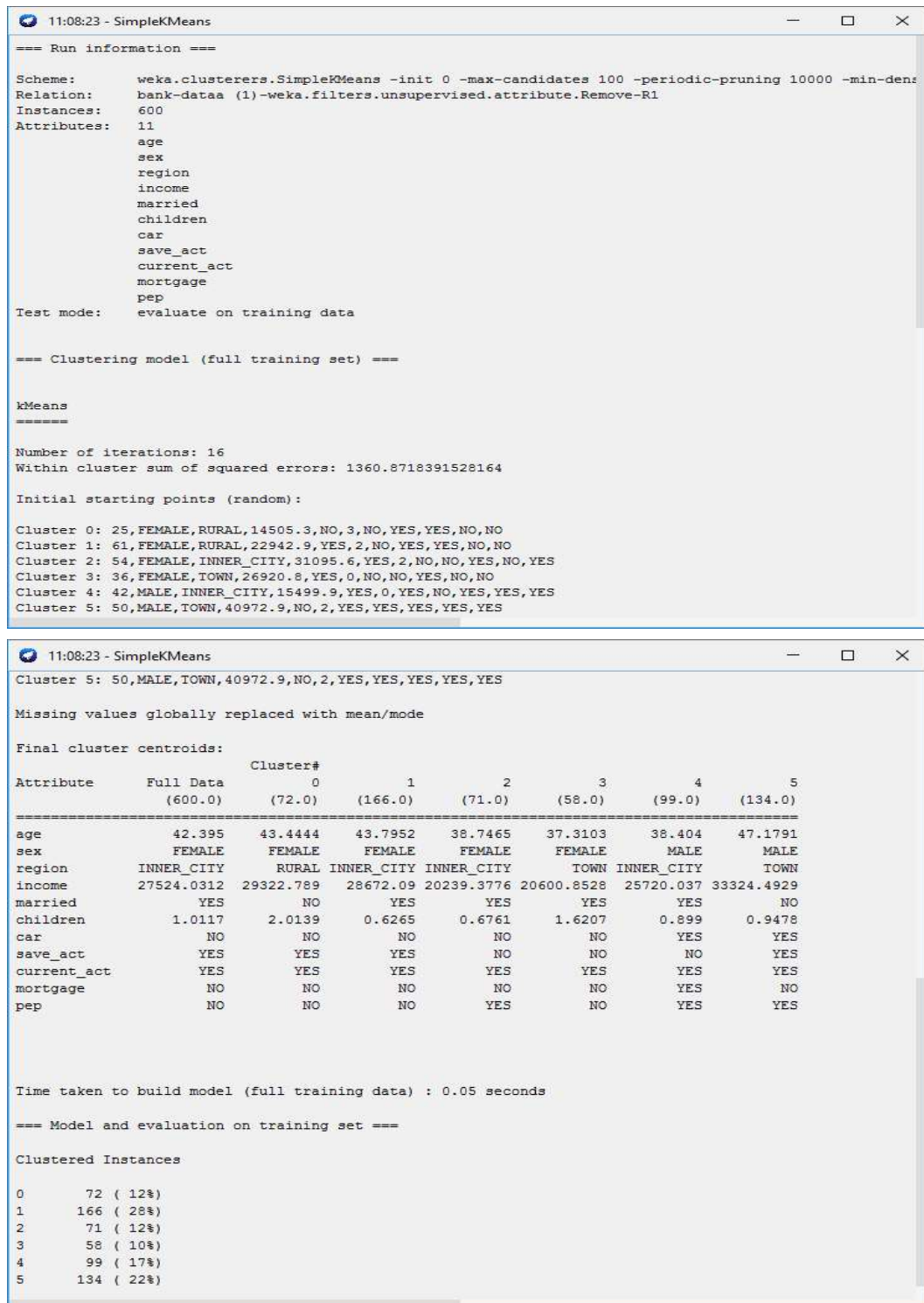
### Setting the Number of Clusters to 6



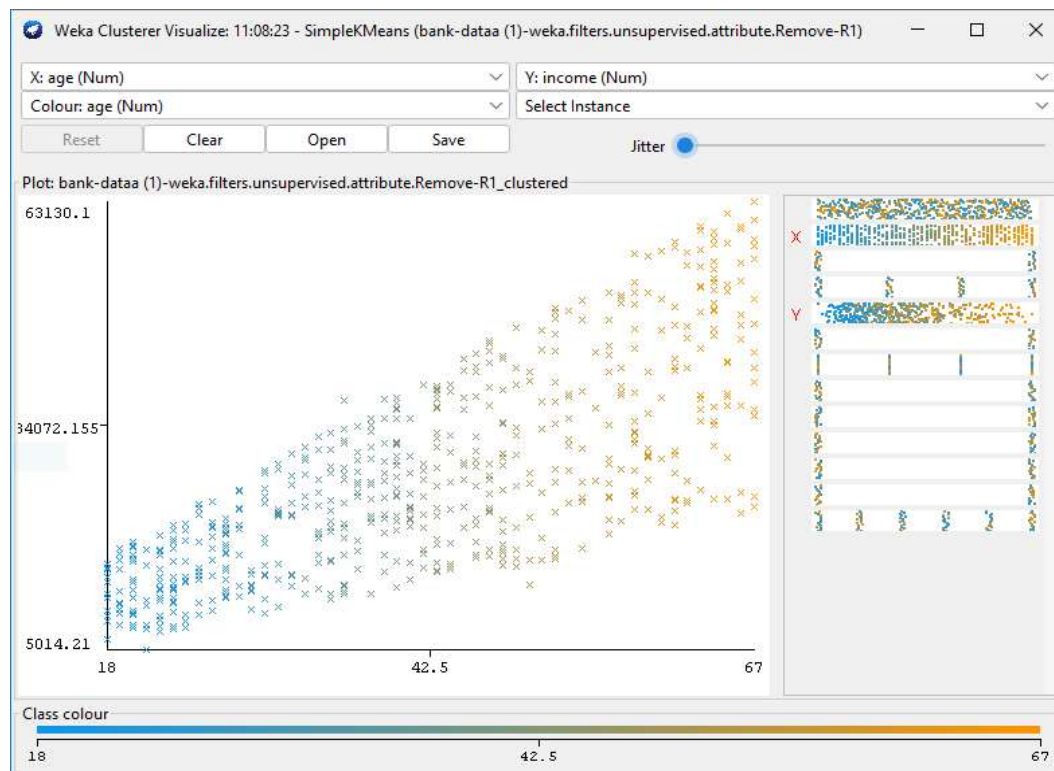
### K-Means Clustering Result in WEKA:

The configuration for the K-Means clustering scheme is as follows:

```
"weka.clusterers.SimpleKMeans -init 0 -max-candidates 100 -periodic-pruning 10000  
-min-density 2.0 -t1 -1.25 -t2 -1.0 -N 6 -A "weka.core.EuclideanDistance -R first-last"  
-I 500 -num-slots 1 -S 10"
```



The scatter plot represents the relationship between age (on the X-axis) and income (on the Y-axis), with colors differentiating various age groups. In the right panel, additional visualizations display the clustered data. The color bar at the bottom shows the different cluster assignments. The Jitter slider is used to adjust the plot by adding slight randomness, which helps to resolve overlapping data points.



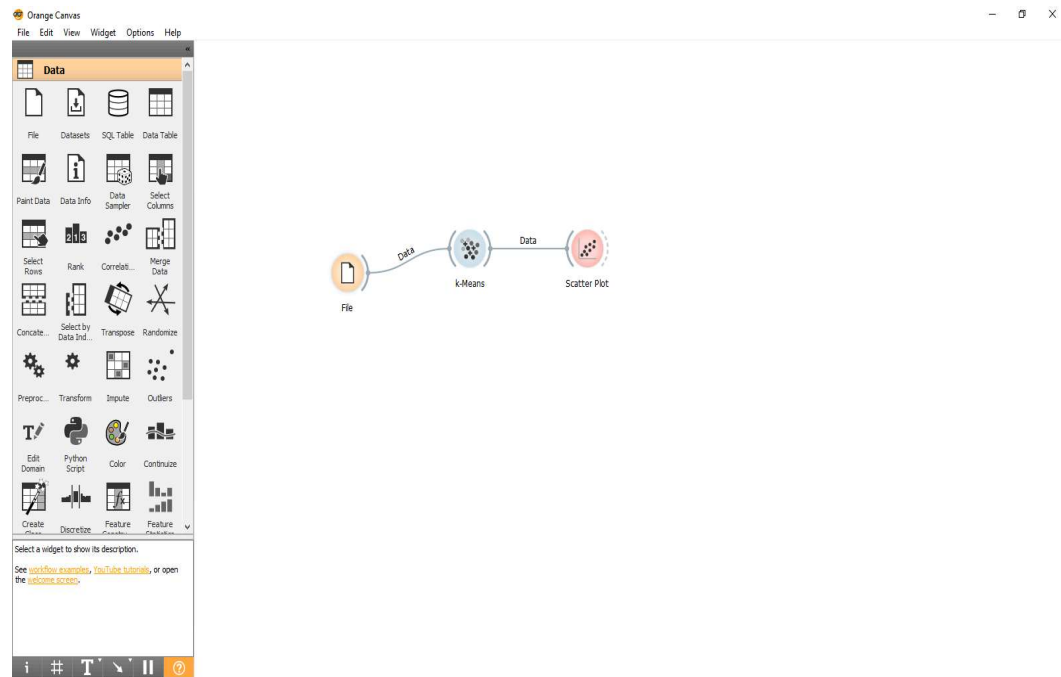
### Clustered Dataset:

```
@relation 'bank-dataa (1)-
weka.filters.unsupervised.attribute.Remove-R1_clustered'

@attribute Instance_number numeric
@attribute age numeric
@attribute sex {FEMALE,MALE}
@attribute region {INNER_CITY,TOWN,RURAL,SUBURBAN}
@attribute income numeric
@attribute married {NO,YES}
@attribute children numeric
@attribute car {NO,YES}
@attribute save_act {NO,YES}
@attribute current_act {NO,YES}
@attribute mortgage {NO,YES}
@attribute pep {YES,NO}
@attribute Cluster
{cluster0,cluster1,cluster2,cluster3,cluster4,cluster5}

@data
0,48,FEMALE,INNER_CITY,17546,NO,1,NO,NO,NO,NO,YES,cluster2
1,40,MALE,TOWN,30085.1,YES,3,YES,NO,YES,YES,NO,cluster4
2,51,FEMALE,INNER_CITY,16575.4,YES,0,YES,YES,YES,NO,NO,cluster1
3,23,FEMALE,TOWN,20375.4,YES,3,NO,NO,YES,NO,NO,cluster3
4,57,FEMALE,RURAL,50576.3,YES,0,NO,YES,NO,NO,NO,cluster1
5,57,FEMALE,TOWN,37869.6,YES,2,NO,YES,YES,NO,YES,cluster3
6,22,MALE,RURAL,8877.07,NO,0,NO,NO,YES,NO,YES,cluster2
7,58,MALE,TOWN,24946.6,YES,0,YES,YES,YES,NO,NO,cluster5
8,37,FEMALE,SUBURBAN,25304.3,YES,2,YES,NO,NO,NO,NO,cluster3
9,54,MALE,TOWN,24212.1,YES,2,YES,YES,YES,NO,NO,cluster5
10,66,FEMALE,TOWN,59803.9,YES,0,NO,YES,YES,NO,NO,cluster1
11,52,FEMALE,INNER_CITY,26658.8,NO,0,YES,YES,YES,YES,NO,cluster1
12,44,FEMALE,TOWN,15735.8,YES,1,NO,YES,YES,YES,YES,cluster2
13,66,FEMALE,TOWN,55204.7,YES,1,YES,YES,YES,YES,YES,cluster5
14,36,MALE,RURAL,19474.6,YES,0,NO,YES,YES,YES,NO,cluster1
15,38,FEMALE,INNER_CITY,22342.1,YES,0,YES,YES,YES,YES,NO,cluster
1
```

## K-Means Clustering in Orange



### Loading Dataset in ORANGE:

The screenshot shows the 'File' widget in the Orange Data Mining software. The 'File' tab is selected, and the dataset 'bank-data.csv' is loaded. The 'Info' section indicates that the dataset contains 600 instances, 11 features, and 1 meta-attribute, with no target variable. The 'Columns' section displays a table with the following data:

	Name	Type	Role	Values
1	age	N numeric	feature	
2	sex	C categorical	feature	FEMALE, MALE
3	region	C categorical	feature	INNER_CITY, RURAL, SUBURBAN, TOWN
4	income	N numeric	feature	
5	married	C categorical	feature	NO, YES
6	children	N numeric	feature	
7	car	C categorical	feature	NO, YES
8	save_act	C categorical	feature	NO, YES
9	current_act	C categorical	feature	NO, YES
10				

At the bottom of the window, there is a 'Browse documentation datasets' button and an 'Apply' button.

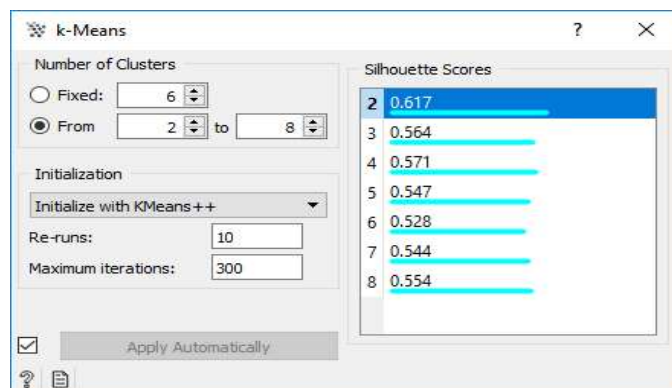


Name: Vishal Rajesh Mahajan

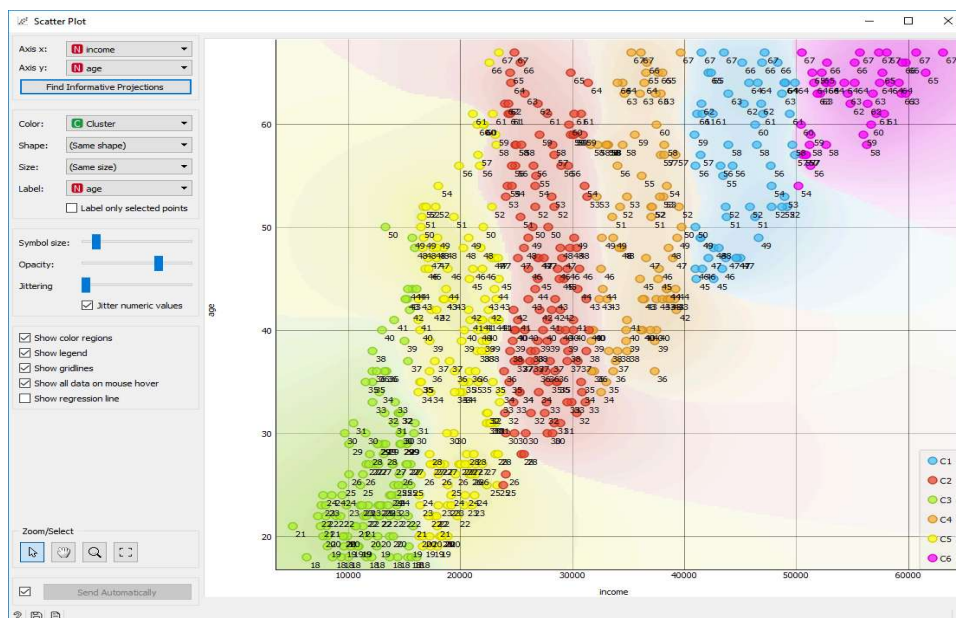
Class: TE IT A

BI Lab EXP 07

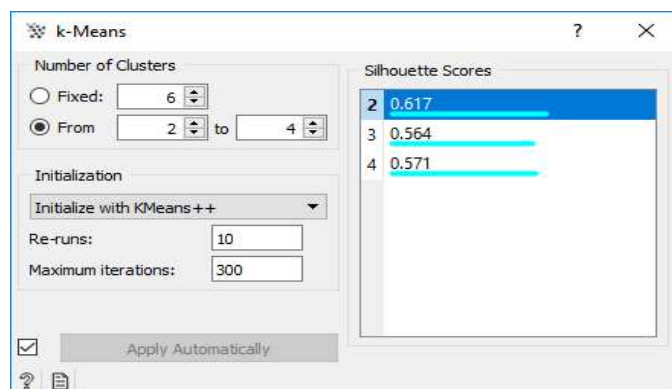
Roll No: 56



Graph of Scatter Plot in WEKA:



Determining the Number of Clusters Based on Silhouette Score:



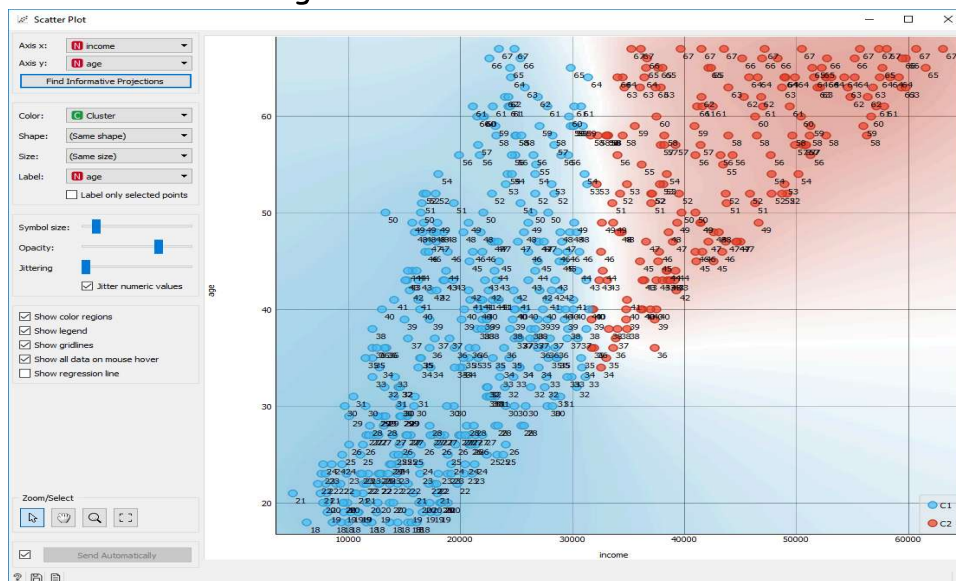
Name: Vishal Rajesh Mahajan

BI Lab EXP 07

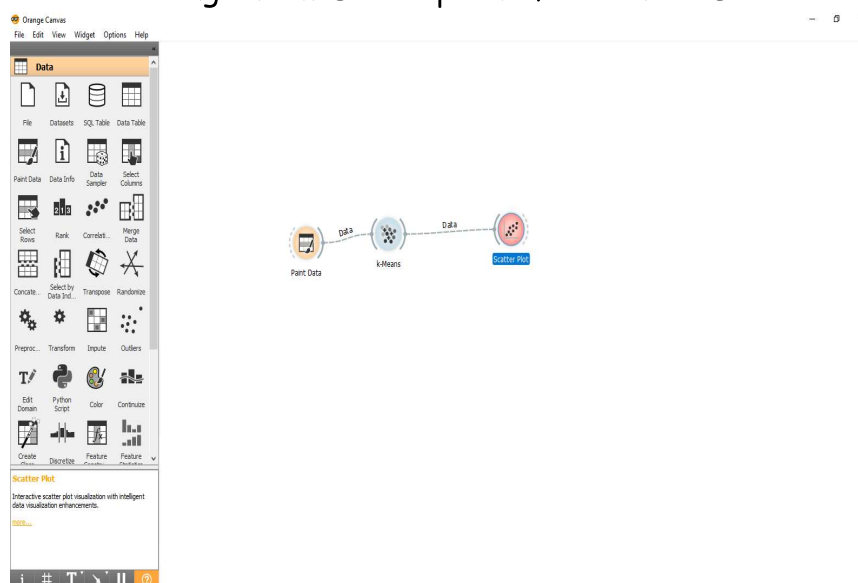
Class: TE IT A

Roll No: 56

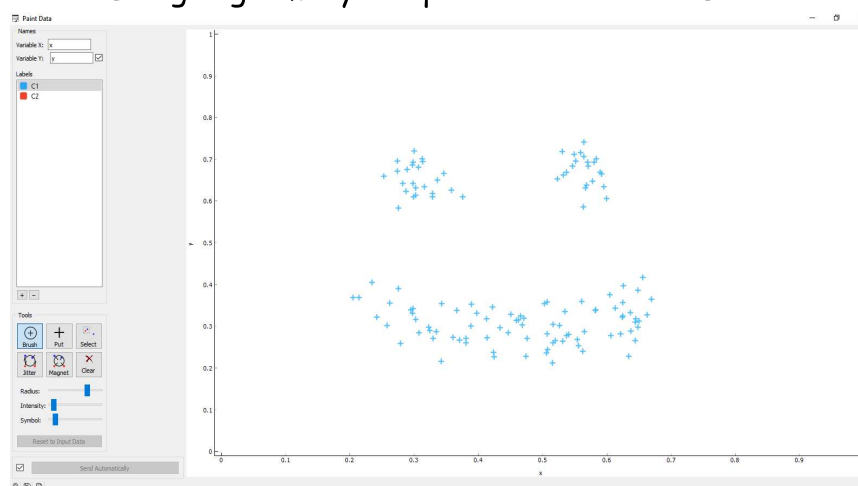
## Observing Two Clusters in the Scatter Plot:



## Creating a New Data Pipeline for Painted Data:



## Designing Smiley Shapes in the Painted Data:





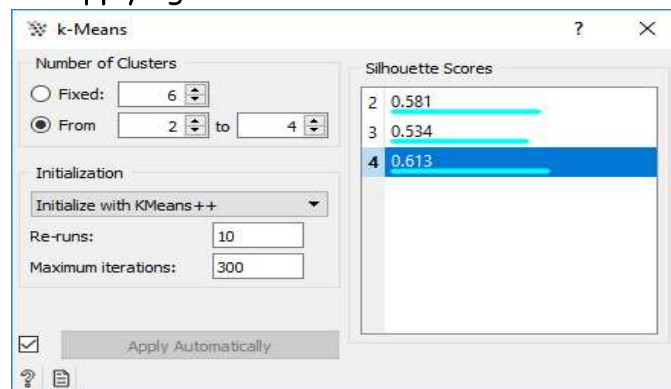
Name: Vishal Rajesh Mahajan

BI Lab EXP 07

Class: TE IT A

Roll No: 56

### Applying K-Means on the Painted Data:



### Output of Scatter Plot with Clusters:

