# EXPERIMENT-6

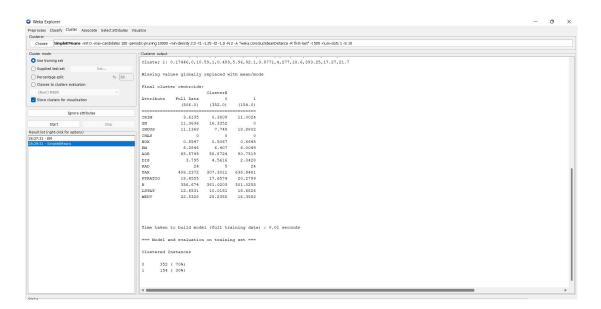## DATA SEGMENTATION BY K- MEANS CLUSTER
## USING WEKA AND R-TOOL

Name: S.G.DEVSACHIN
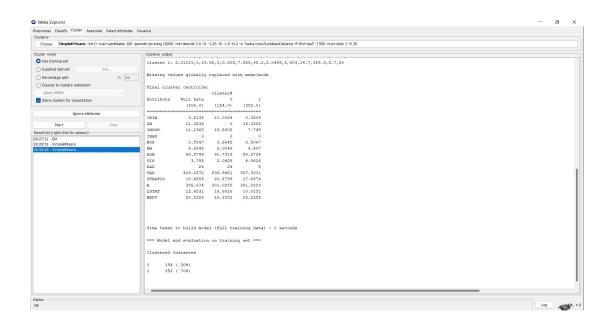Reg.No: 192111088
Subject: CSA1672, Data warehouse and data mining
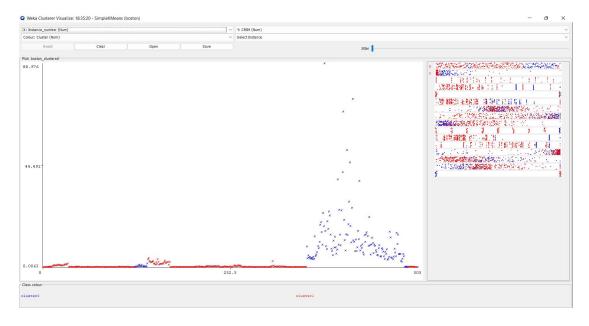
## OUTPUT:

1)**Choose a set of attributes for clustering and give a motivation.**



2)**Experiment with atleast 2 different number of clusters but with same seed values:**

**USING R-TOOL**



```
>
> clusters <- kmeans(citycrimes[,2:3], 5)
>
> citycrimes$Borough <- as.factor(clusters$cluster)
> str(clusters)
> str(clusters)
List of 9
 $ cluster     : int [1:24] 4 1 3 4 2 4 1 4 5 1 ...
 $ centers     : num [1:5, 1:2] 2079 8513 2908 1391 4509 ...
 ..- attr(*, "dimnames")=List of 2
 .. ..$ : chr [1:5] "1" "2" "3" "4" ...
 .. ..$ : chr [1:2] "Pop" "WC"
 $ totss       : num 1.39e+08
 $ withinss    : num [1:5] 315574 13643048 67068 52387 191844
 $ tot.withinss: num 14269920
 $ betweenss   : num 1.25e+08
 $ size        : int [1:5] 7 3 3 9 2
 $ iter        : int 3
 $ ifault      : int 0
 - attr(*, "class")= chr "kmeans"
> library(ggmap)
```

```
<environment: namespace.stats>
> ggmap(Map) + geom_point(aes(x = Pop[], y = WC[], colour = as.factor(Borough)),data = c
itycrimes)
Warning message:
Removed 24 rows containing missing values (geom_point).
> ggtitle("Map Boroughs using KMean")
$title
[1] "Map Boroughs using KMean"

$subtitle
NULL

attr(,"class")
[1] "labels"
>
> |
```


Map Boroughs using KMean