

## Cluster Analysis Document

Vivek I(IMT2013017)

Aravind K(IMT2013019)

Nikhilesh P(IMT2013031)

Arun Prasad(IMT2013037)

As part of clustering for the data we have chosen for the project, we have performed the following activities:

(1)K-means clustering

- Decided the number of clusters(k) to be chosen using elbow graph
- Plot of different clusterings.
- Silhouette plot to evaluate the clusterings we got.

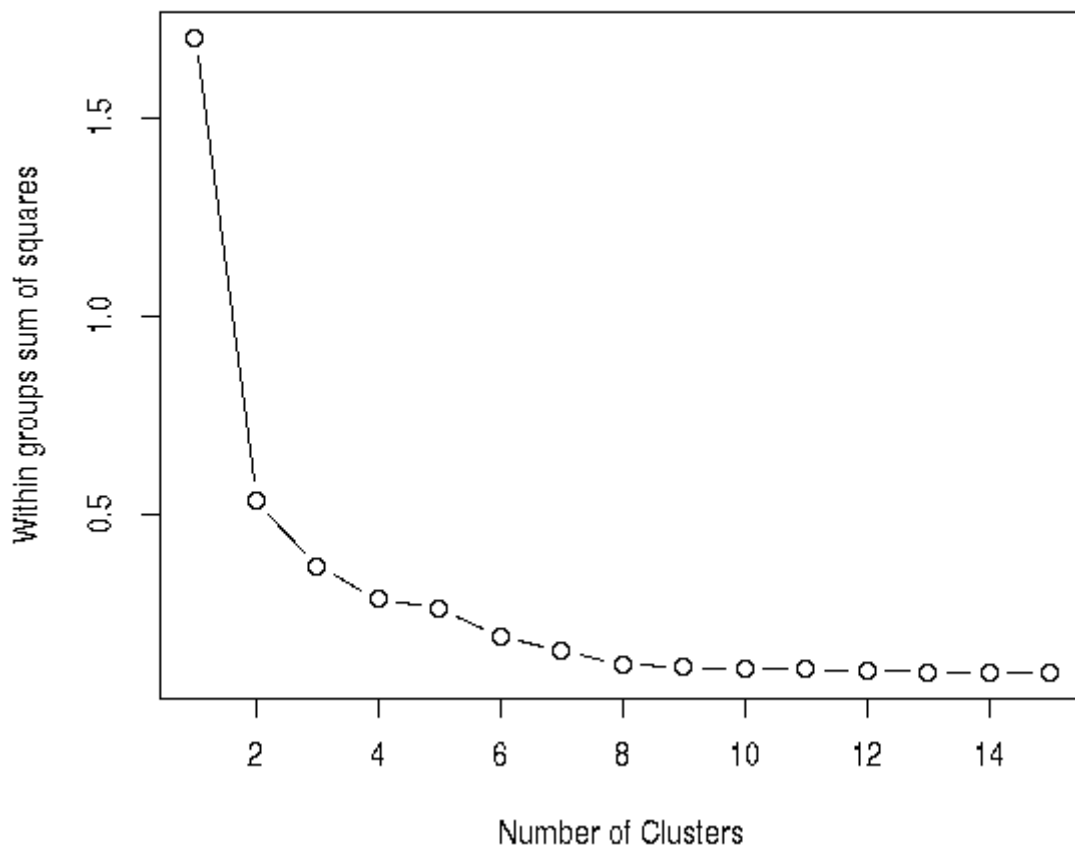
(2)Hirearcheal clustering

- Plot of different clusterings
- Silhouette plot to evaluate the clusterings we got.

Before performing the clustering algorithms, we have normalized our data and also converted categorical to numerical attributes to avoid the effect of outliers and to bring all the attribute's values to same scale(or range i.e b/w 0 and 1).

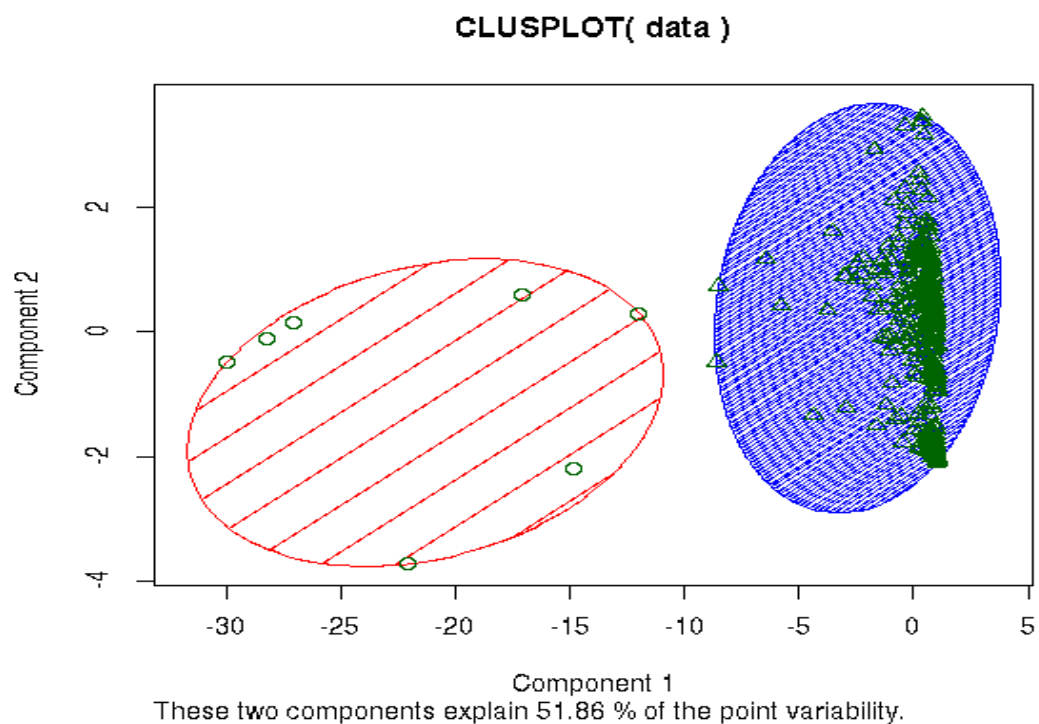
**K-means clustering :**

Since for kmeans function, number of clusters(k) needs to be given as input, we have first plotted the elbow graph for deciding this.



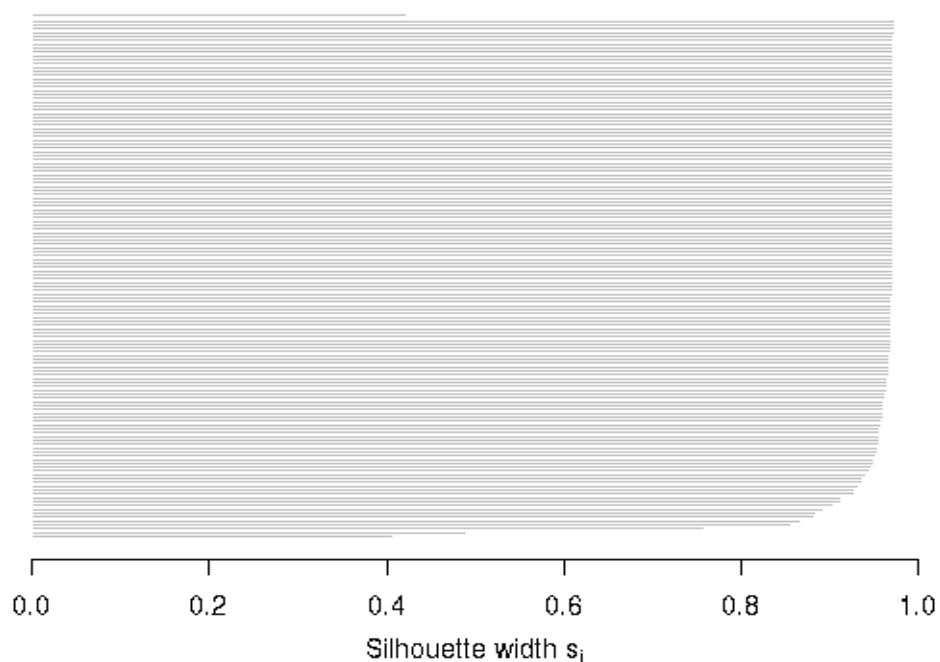
By seeing this graph, we have applied kmeans algorithm for k=2,3,4. The results of these are as follows:

K=2:



**Silhouette plot of (x = km\$cluster, dlist = dls)**

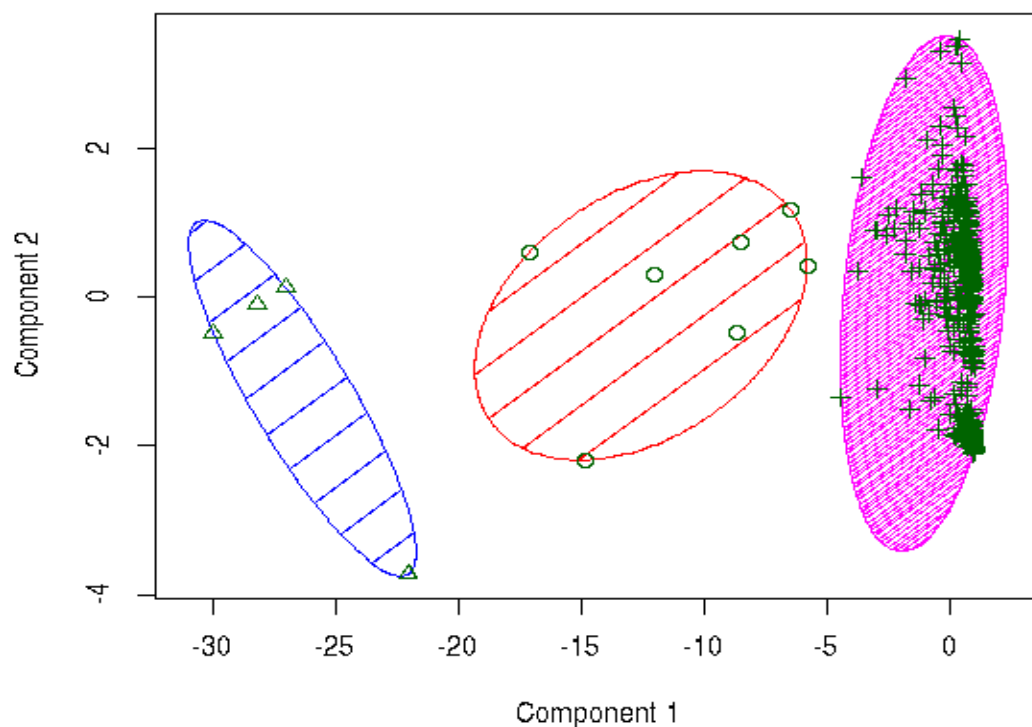
n = 180



Average silhouette width : 0.95

K=3 :

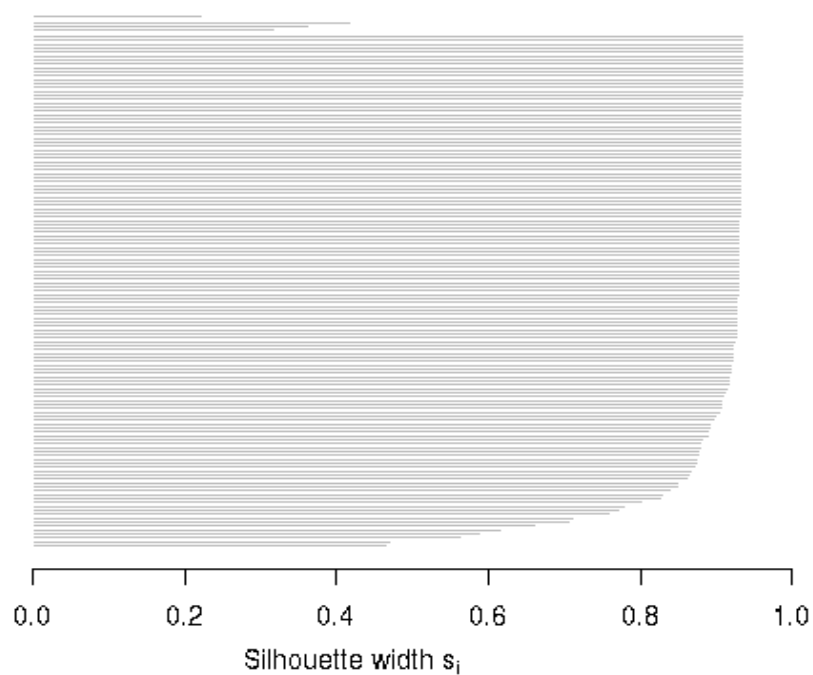
**CLUSPLOT( data )**



These two components explain 51.86 % of the point variability.

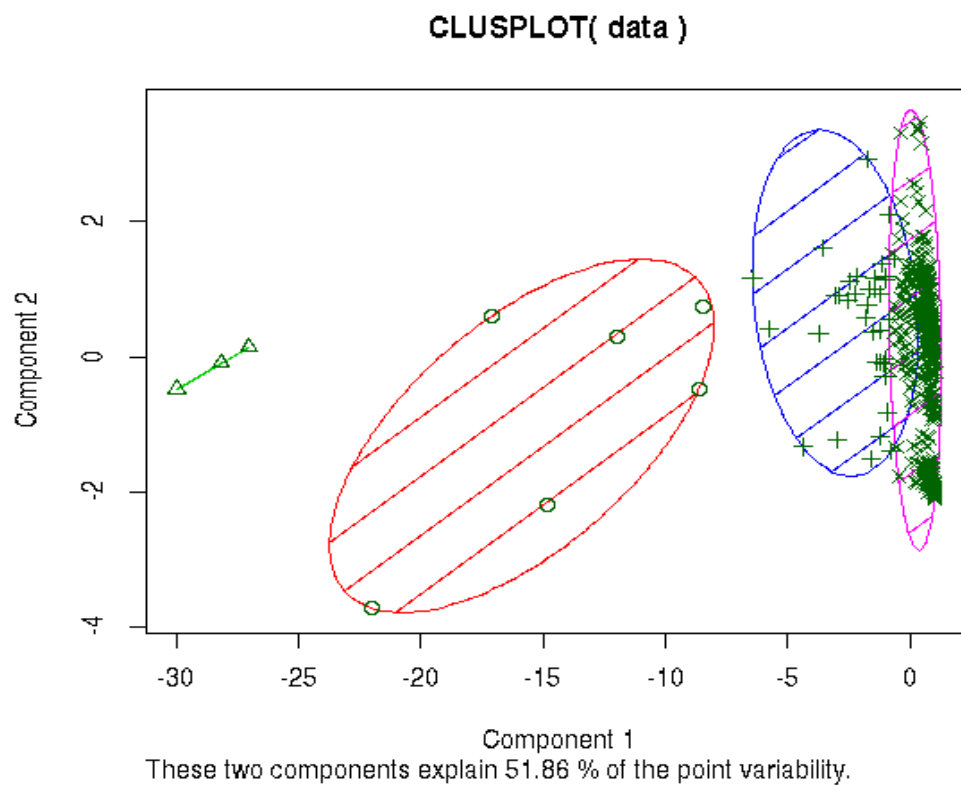
**Silhouette plot of (x = km\$cluster, dlist = dls)**

n = 180



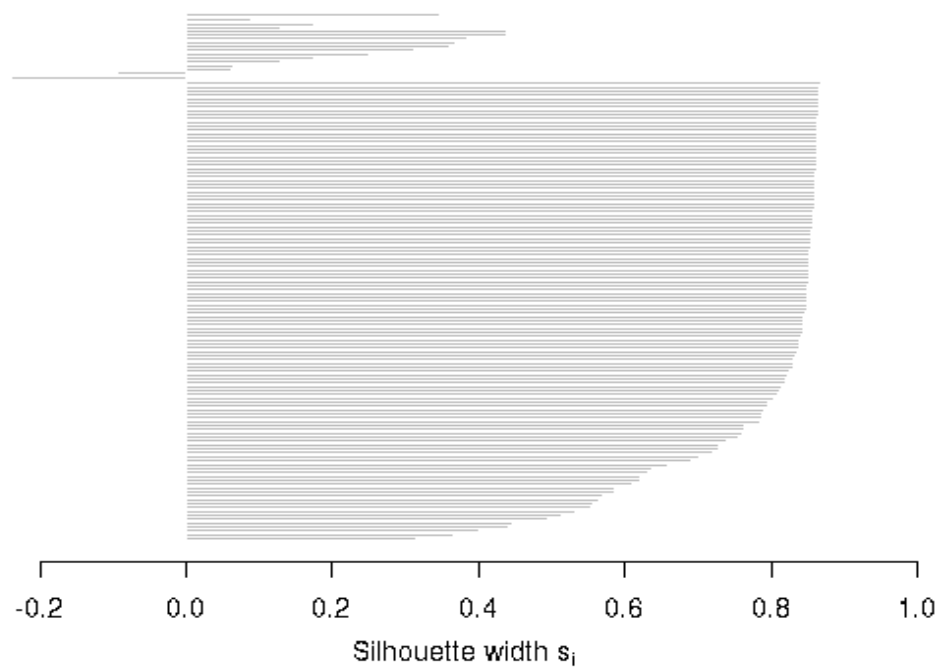
Average silhouette width : 0.87

K=4 :



**Silhouette plot of (x = km\$cluster, dlist = dls)**

n = 180

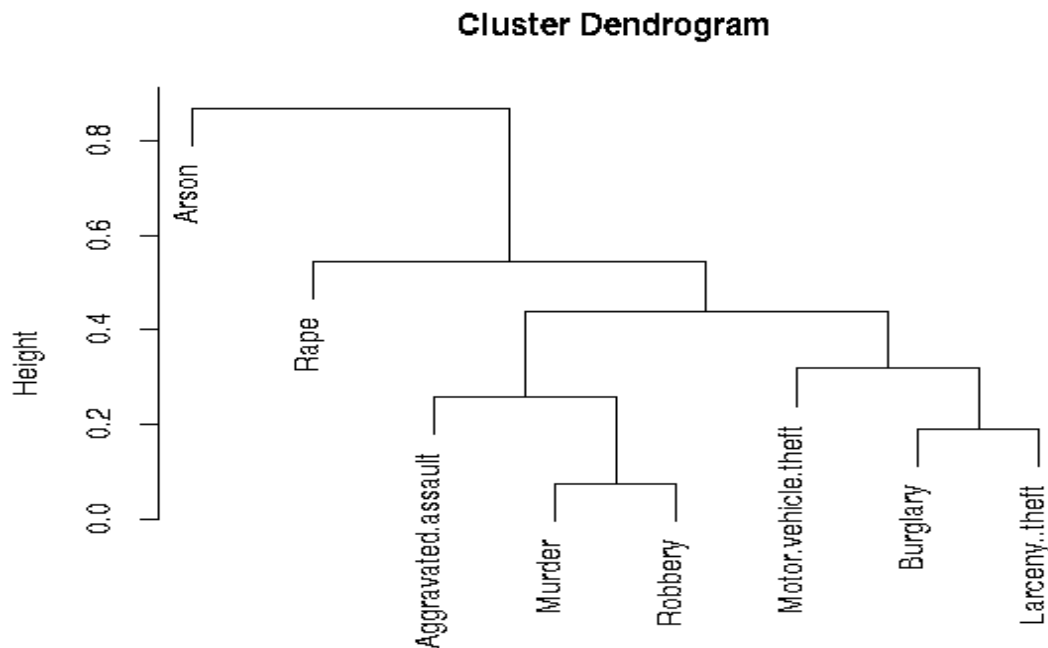


Average silhouette width : 0.71

So, from above silhouette graphs, it's clear that the average silhouette width is more for  $k=2$ . Also, it's clear from the graph that, more points fall in one cluster and very few fall in the remaining clusters, which implies that all the data points are close to each other except few data points (or) in other words, the number of crimes (in all categories) in all the cities almost remain the same except few cities which have more (or) less number of crimes and these fall into other clusters.

### Hierarchical clustering :

Below shown is the output of hierarchical clustering which is a cluster dendrogram:

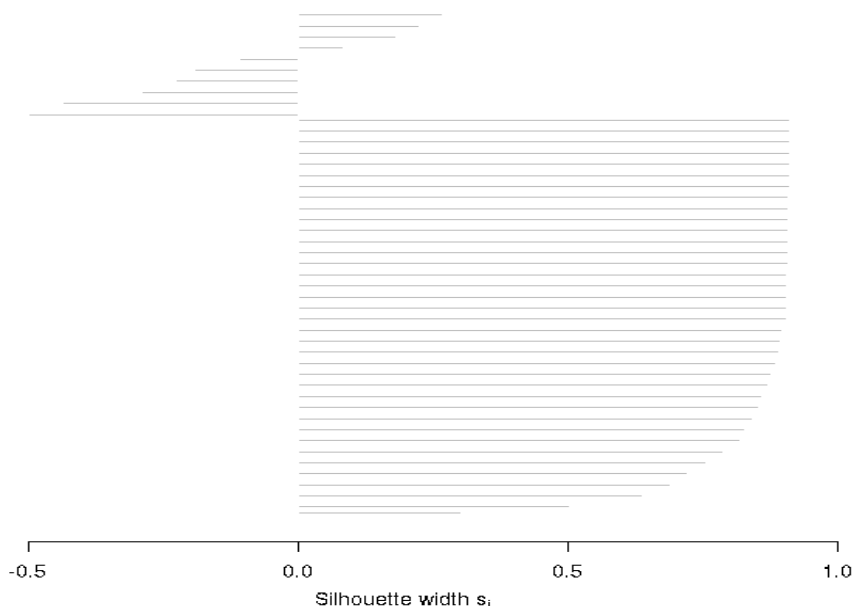


This implies data in the features which are in same subtree of any node are closer (or has similar spread among them) compared to data in the features which are in the other subtree.

Below is the silhouette plot for the above clustering:

**Silhouette plot of pam(x = dist(data[c(2:9)]), k = 3)**

n = 436



Average silhouette width : 0.63