Name:-Shantanu S. Mangale
B.E B 41242 DMW
LP-II

# Assignment No. 2

**Title**:- Consider a suitable dataset. for a clustering of data instances in different groups, apply different clustering techniques (minimum 2). ~~Virt~~ Visualize the clusters using suitable tool.

## Problem Definition:-

Visualize the cluster using suitable tool.

## Pre-requisite:-

i) Basic concepts of ETL.
ii) Knowledge about R tool.

## S/W & H/W Requirements:-

R-tool, PIV, 2GB RAB, 500 GB HDD.

## Learning Objectives:-

Use R functions to create K-means Clustering models & Heirarchical clustering models.

## Outcomes:-

Visualize the effectiveness of the K-means clustering algorithm & hierarchical clustering using graphic
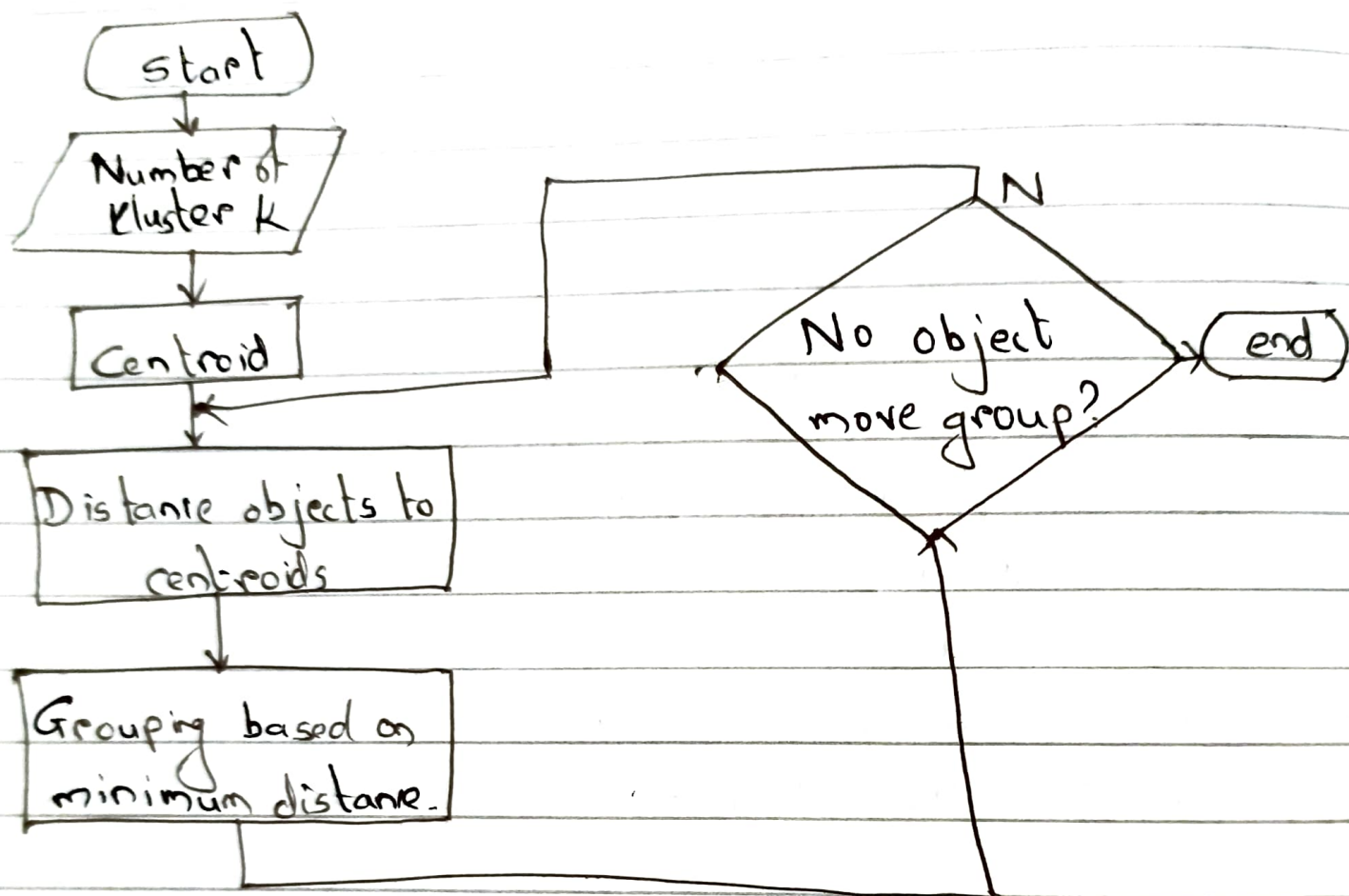
capabilities in R.

Theory:-

• What is K-means clustering?

K-means clustering is a type of unsupervised learning which is used when you have unlabeled data (i.e., data without defined categories or groups). The goal of this algorithm is to find groups in the data, with the number of groups represented by the variable K. Data points are clustered based on feature similarity. The results of the K-means clustering algorithm are:

i) The centroids of the K-clusters, which can be used to label new data.

ii) Labels for the training data (each data point is assigned to a single cluster).

• Steps to perform K-means clustering.

```
┌──────────┐
│  start   │
└──────────┘
      │
      ▼
 ╱──────────╲
│ Number of  │
│ Kluster K  │
 ╲──────────╱
      │
      ▼
┌──────────┐                                          ┌─────╲
│ Centroid │                              N           │      ╲    ┌──────┐
└──────────┘                            ◇─────────────┤  No object  ├──────┤ end  │
      │                                 │ move group? │      ╱    └──────┘
      ▼                                  ╲──────────╱
┌──────────────────┐
│ Distance objects to│
│   centroids        │
└──────────────────┘
      │
      ▼
┌──────────────────┐
│ Grouping based on │
│ minimum distance. │
└──────────────────┘
```

- **R Implementations:**

The k-means function, provided by the cluster package, is used as follows:

```
K-means (x, centers, iter.max' = 10, nstart =1,
          algorithm = c("Hartigan-Wong", "Lloyd",
          "Forgy", "MacQueen"))
```

where the arguments are:

$x$ → A numeric matrix of data or an object that can be co-erced to such a matrix (such as a numeric

vector or a data frame with all numeric columns).

centers → Either the number of clusters or a set of initial (distinct) cluster centers. If a number a random set of (distinct) rows in $x$ is chosen as the initial centers

iter.max → The maximum number of iterations allowed.

nstart → If centers is a number, nstart gives the number of random sets that should be chosen.

algorithm → The algorithm to be used. It should be one of the following.
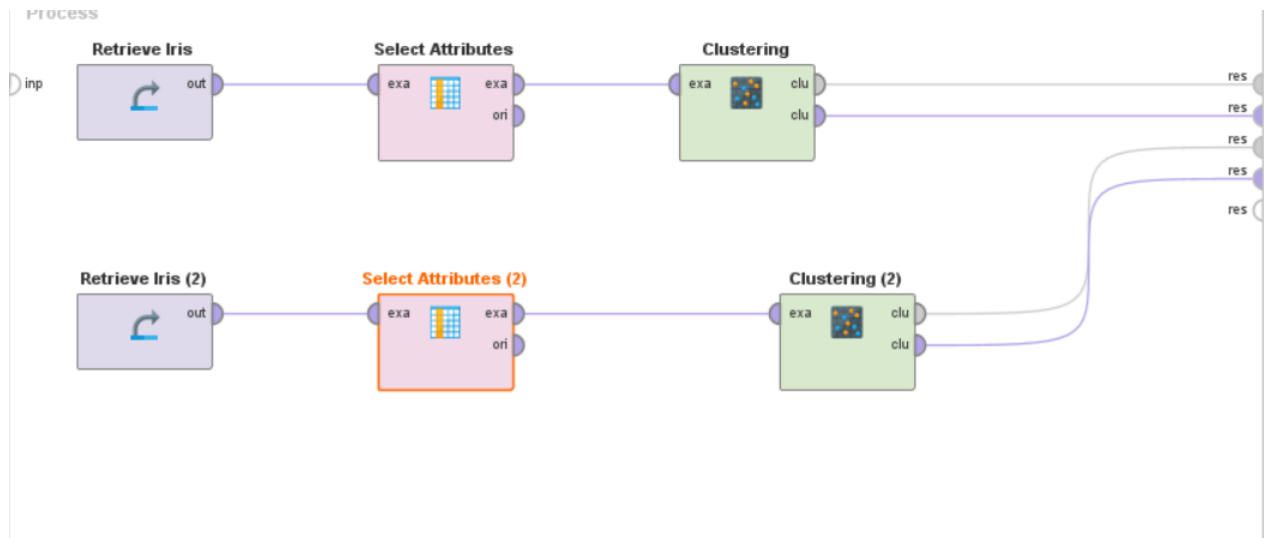
• IRIS dataset:

This is perhaps the best known database to be found in the pattern recognition literature. The data set contains 3 classes of 50 instances each, where each class refes to a type of its iris plant.

Conclusion:-

Visualized the effectiveness of the k-means clustering algorithm & hierarchical clustering using graphic capabilities in R.

**Output:**

**Clustering using K-means and DB Scan**



**Visualization:**

ExampleSet

● Iris-setosa ● Iris-versicolor ● Iris-virginica