

Report About

Clustering Techniques

Content

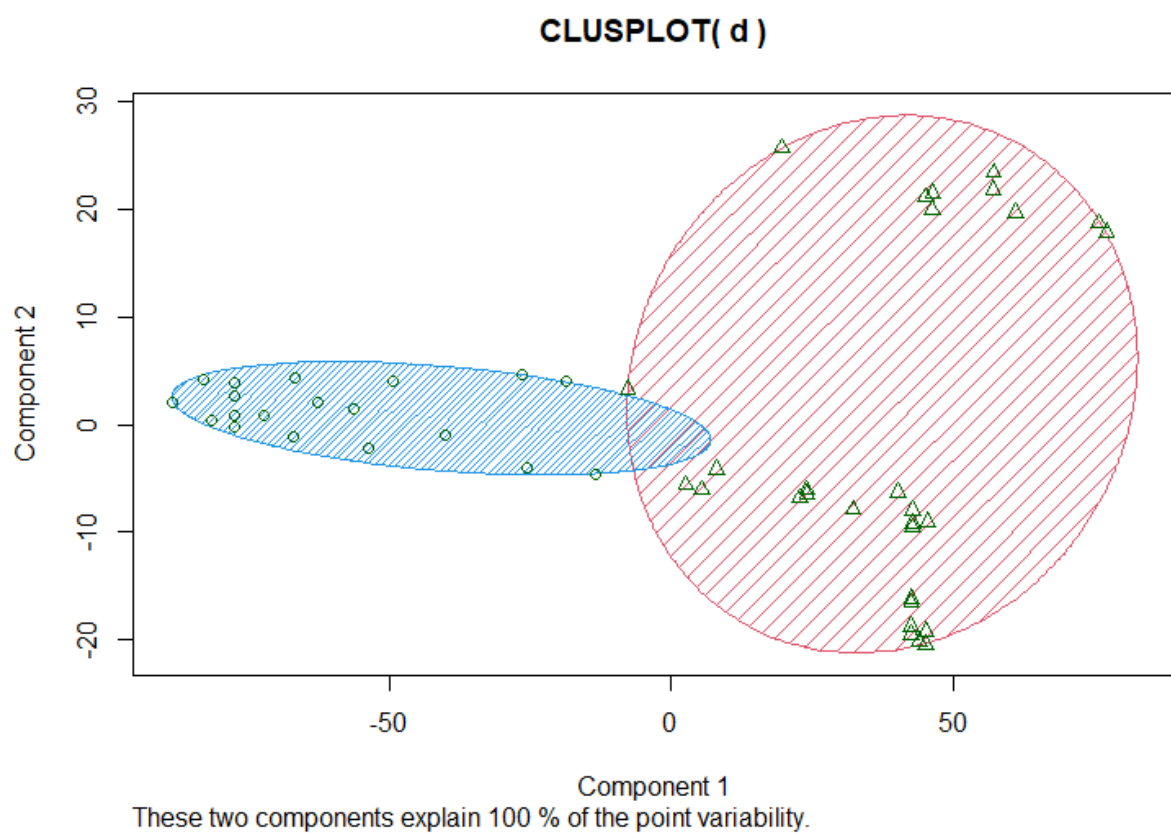
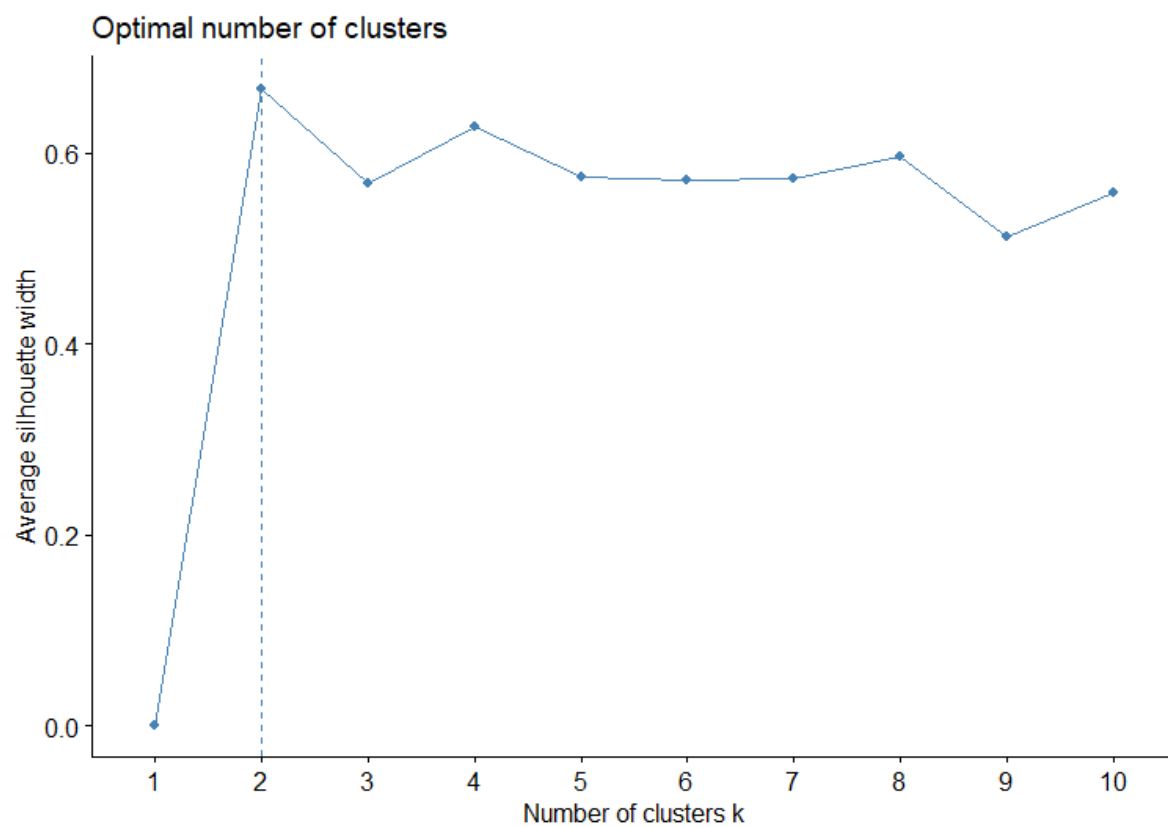
- Introduction.
- K-means clustering.
- Hierarchical clustering using Agglomerative.
- Hierarchical clustering using Divisive.
- Density Based clustering.

Introduction

Cluster analysis or clustering is the task of grouping a set of objects in such a way that objects in the same group (called a cluster) are more similar (in some sense) to each other than to those in other groups (clusters). It is a main task of exploratory data analysis, and a common technique for statistical data analysis, used in many fields, including pattern recognition, image analysis, information retrieval, bioinformatics, data compression, computer graphics and machine learning.

K means clustering

```
5 # Loading package
6 library(factoextra)
7 library(ggfortify)
8 library(ggplot2)
9 library(dplyr)
10 library(cluster)
11 df<-read.csv("data 2.enc") #read the data
12 df #View the data
13 d= select(df,c(2,3))
14 d
15 #function to compute average silhouette for k clusters
16 fviz_nbclust(d, kmeans, method = "silhouette")
17
18 # Fitting K-Means clustering Model
19 kms=kmeans(d,2)
20 kms
21
22 #the centers of clusters
23 kms$centers
24 # Cluster identification for each observation
25 kms$cluster
26
27 # Visualizing clusters
28
29 autoplot(kms,d,frame=T)
30 clusplot(d,kms$cluster ,color=T,shade=T)
31 #plot results of final k-means model
32 fviz_cluster(kms, data = d)
```



Advantages:

- If variables are huge, then K-Means most of the times computationally faster than hierarchical clustering, if we keep k small.
- K-Means produce tighter clusters than hierarchical clustering, especially if the clusters are globular.
- Relatively simple to implement.
- Scales to large data sets.
- Can warm-start the positions of centroids.
- Easily adapts to new examples.

Disadvantages:

- Clustering outliers.
- Centroids can be dragged by outliers, or outliers might get their own cluster instead of being ignored. Consider removing or clipping outliers before clustering.
- Scaling with number of dimensions.
- As the number of dimensions increases, a distance-based similarity measure converges to a constant value between any given examples. Reduce dimensionality either by using PCA on the feature data, or by using “spectral clustering” to modify the clustering algorithm as explained below.
- Different initial partitions can result in different final clusters.
- It does not work well with clusters (in the original data) of Different size and Different density.

when k-means preferable to be used?

k-means becomes a great solution for pre-clustering, reducing the space into disjoint smaller sub-spaces where other clustering algorithms can be applied.

The K-means clustering algorithm is used to find groups which have not been explicitly labeled in the data. This can be used to confirm business assumptions about what types of groups exist or to identify unknown groups in complex data sets.

Can warm-start the positions of centroids. Easily adapts to new examples. Generalizes to clusters of different shapes and sizes, such as elliptical clusters.

if variables are huge, then K-Means most of the times computationally faster than hierarchical clustering.

How can k-means clustering be improved?

K-means clustering algorithm can be significantly improved by using a better initialization technique, and by repeating (re-starting) the algorithm.

When the data has well separated clusters, the performance of k-means depends completely on the goodness of the initialization.

Initialization using simple furthest point heuristic (Maxmin) reduces the clustering error of k-means from 15% to 6%, on average.

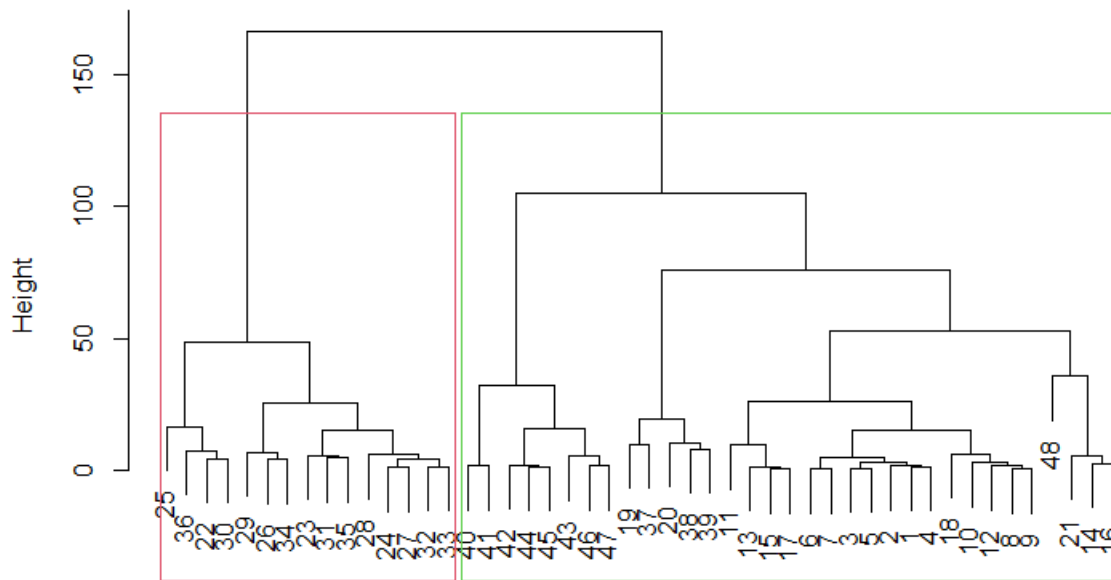
Hierarchical clustering using Agglomerative

```
1 # importing libraries
2 library(ggplot2)
3 library(cluster)
4 library(ggfortify)
5 library(stats)
6 library(dplyr)
7 library(factoextra)
8 #import data
9 data = read.csv("data 2.enc")
10 data = data[,c(2,3)]
11
12 #Agglomerative using hclust
13 d = dist(data, method = "euclidean")
14 #complete method
15 hcluster = hclust(d, method = "complete")
16 plot(hcluster)
17 group1 <- cutree(hcluster, k=3) # cut tree into 3 clusters
18 # draw dendrogram with red borders around the 5 clusters
19 rect.hclust(hcluster, k=3, border=2:5)
20 table(group1)
21 #single method
22 hcluster = hclust(d, method = "single")
23 plot(hcluster)
24 #select best number of clusters(k)
25 fviz_nbclust(data, FUN = hcut, method = "silhouette")#3 clusters
26 #
```

```
28 #Agglomerative using agnes
29 agh=agnes(data,method="complete")
30 agh$ac #look at the agglomerative coefficient
31 group2 = cutree(as.hclust(agh), k = 3)
32 table(group2)
33 pltree(agh)
34 #Single method
35 agh=agnes(data,method="single")
36 agh$ac
37 pltree(agh)
```

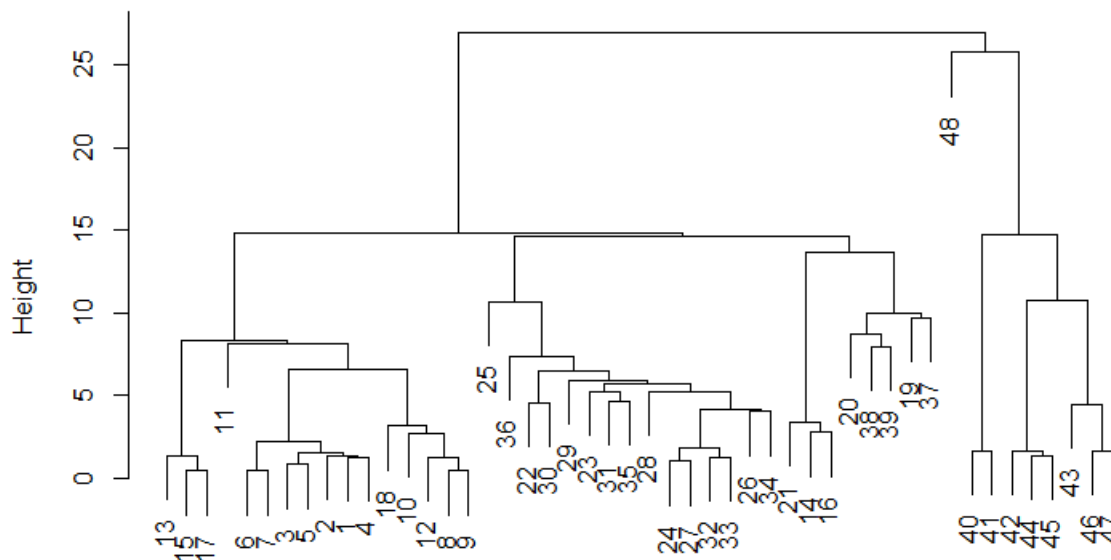

Result

Cluster Dendrogram



```
hclust (*, "complete")
```

Cluster Dendrogram

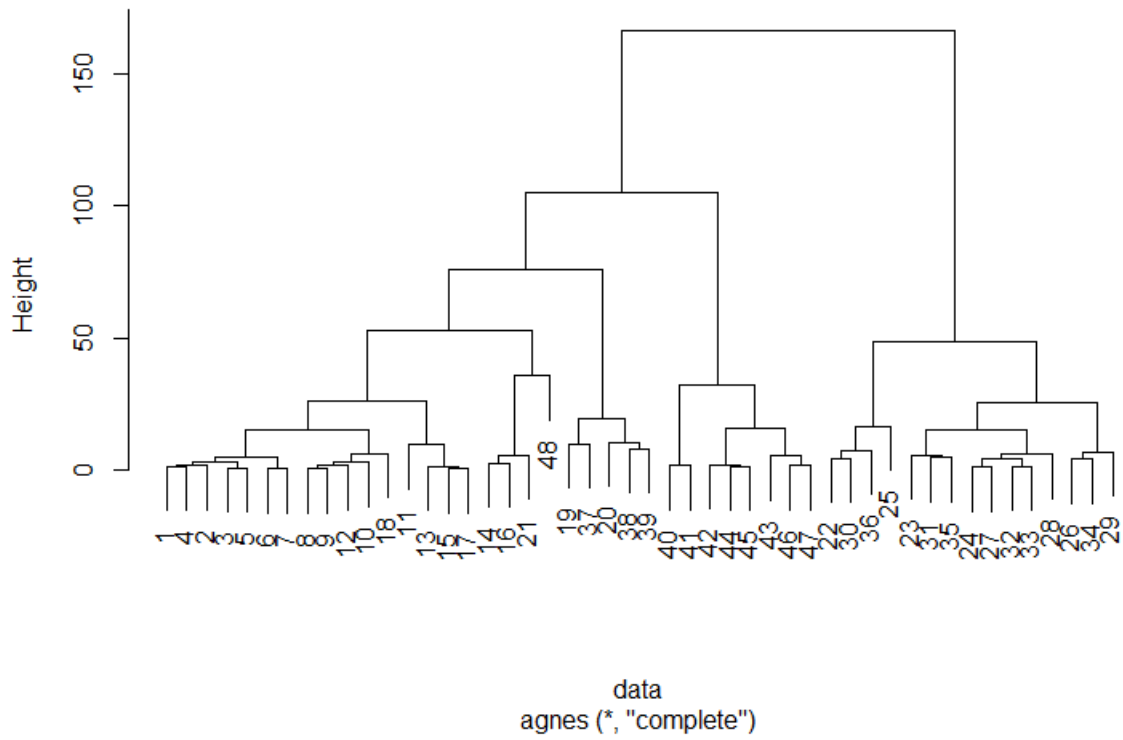


```

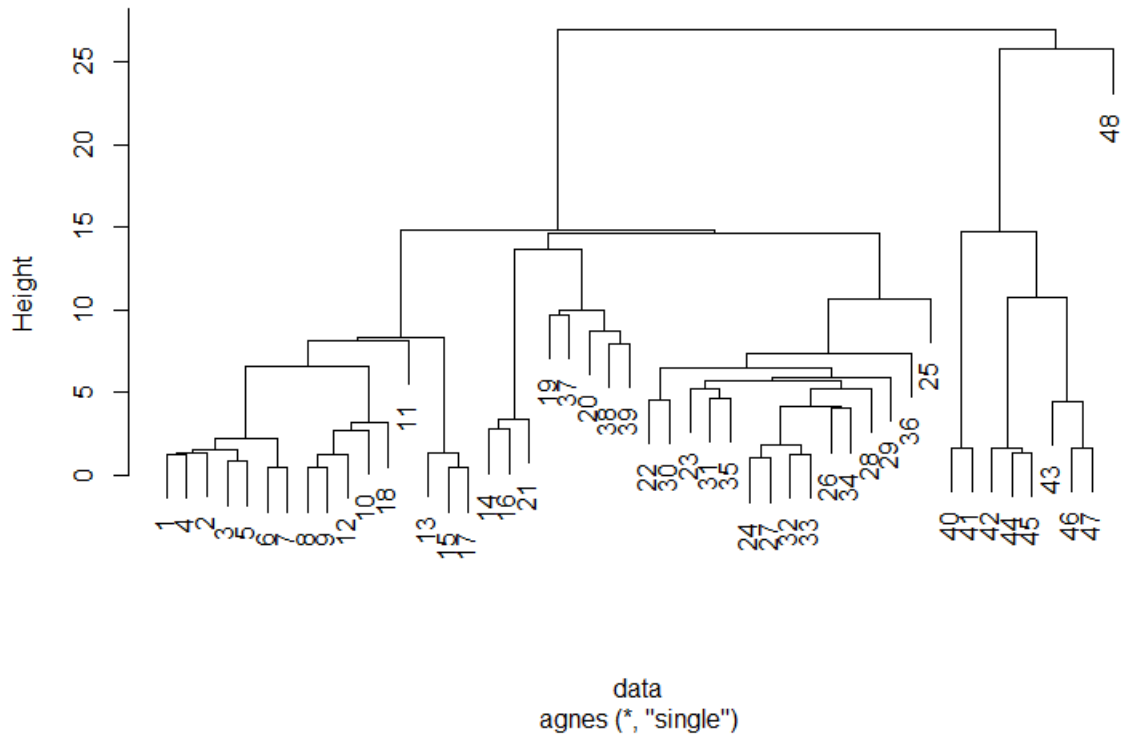
d
hclust (*, "single")

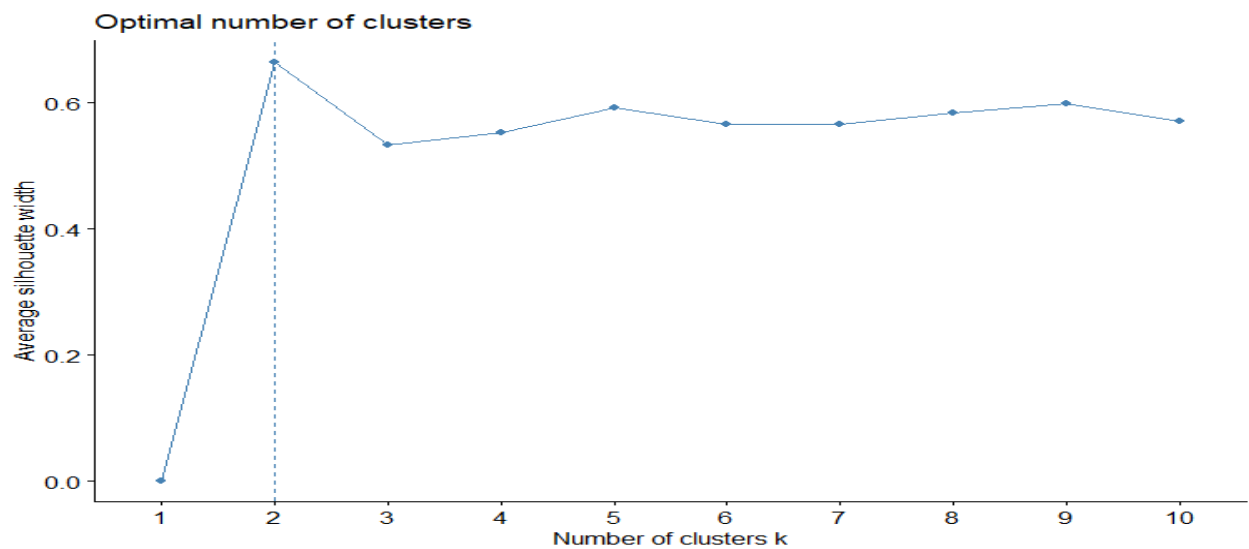
```

Dendrogram of agnes(x = data, method = "complete")



Dendrogram of agnes(x = data, method = "single")





Hclust

```
> table(group1)
group1
 1  2
33 15
```

Agnes

```
> table(group2)
group2
 1  2
33 15
```

Comments on result

Both hclust and agnes have the same result and behave very similarly in single and complete methods. As we see in group 1 and group 2, Agnes has more methods than hclust.

The proximity of two observations can be drawn based on the height where branches containing those two observations first are fused (like height of 14&16 in complete). But we can't do that with the horizontal axis as a criterion of their similarity. In complete method it tends to produce more compact clusters. But in single it tends to produce long "loose" clusters.

We see that the best number of clusters here is 2 by silhouette method.

Agnes yields the agglomerative coefficient which measures the amount of clustering structure found (0.9734672 in complete and 0.8583913 in single), the closer to 1 the strongest clustering structure.

Advantage of Agglomerative

- Agnes function can get the agglomerative coefficient, which measures the amount of clustering structure found (values closer to 1 suggest strong clustering structure).
- Agnes, it allows us to find certain hierarchical clustering methods that can identify stronger clustering structures.
- Agnes contains six agglomerative algorithms, some not included in hclust. Has it's own plot method.
- Hclust allow us to use rect.hclust method - which divides dendrogram into given number of groups (argument k) and draws rectangles around samples in these groups.
- Hclust allow us to use cutree method - cuts the tree (dendrogram) into given number of clusters (argument k) or according to given level of similarity
- hclust function implements both Ward's algorithms (the genuine one, named ward.D2, as well as the second one, called ward.D) unlike agnes which use the genuine one.
- agglomerative clustering is Easy to implement and good at identifying small clusters
- It doesn't require the number of clusters to be specified.

Disadvantage of Agglomerative

- Generally, has long runtimes.
- Time complexity is higher at least $O(n^2 \log n)$.
- Can never undo what was done previously, which means if the objects may have been incorrectly grouped at an earlier stage, and the same result should be close to ensure it.
- the usage of various distance metrics for measuring distances between the clusters may produce different results.
- Not suitable for large datasets.
- Initial seeds have a strong impact on the results.
- The order of the data has an impact on the results
- Very sensitive to outliers

Some cases in which Agglomerative is preferable to be used

when we want to partition our data into groups at different levels such as in a hierarchy and use an algorithm which is Easy to implement and good at identifying small clusters and when it doesn't require the number of clusters to be specified.

How can we improve Agglomerative?

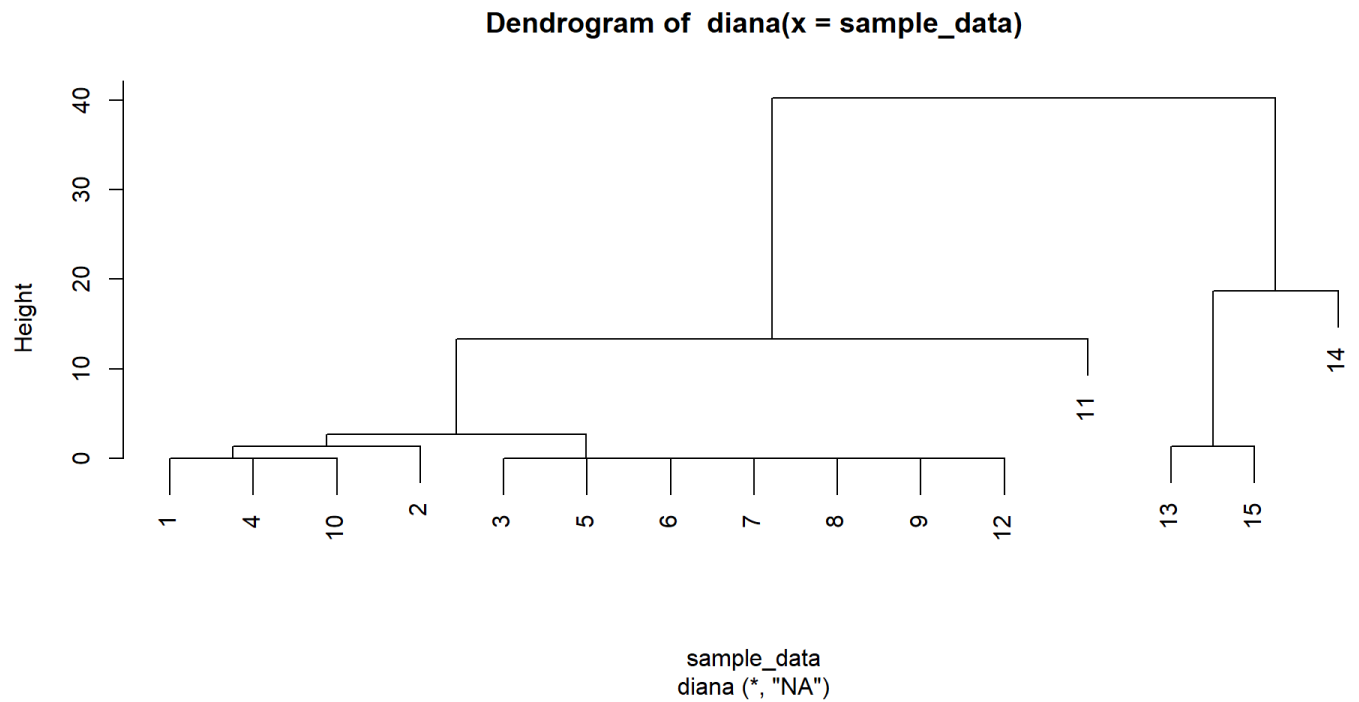
A main drawback of the Agglomerative is that it is very sensitive to noise and thus results in its lack of robustness, which limits its use in real-world applications.

The kernel method in the machine learning theory refers to increasing the computational power of linear methods by mapping the data into high-dimensional feature space and has shown its great power in several kernel-based learning machines, e.g., support vector machines (SVMs) and kernel principal component analysis (KPCA). So, a robust batch clustering algorithm has been developed for clustering incomplete data (nonstationary data) using the kernel method

Hierarchical clustering using Divisive

```
1 install.packages("ggplot2")
2 install.packages("ggfortify")
3 install.packages("factoextra")
4 library(tidyverse)
5 library(cluster)
6 library(ggplot2)
7 library(factoextra)
8 library(mclust)
9 library(ggfortify)
10 library(stats)
11 library(dplyr)
12 View(data.2)
13 data = select(data.2,c(2,3))
14 sample_data=(data[1:15,2,3])
15 dh=diana(sample_data)
16 dh$dc
17 pltree(dh)
18
```

Result:



```
> data = select(data.2,c(2,3))
> sample_data=(data[1:15,2,3])
> dh=diana(sample_data)
> dh$dc
[1] 0.9404643
> pltree(dh)
```

Advantages of divisive:

Divisive clustering is more efficient if we do not generate a complete hierarchy all the way down to individual data leaves. Time complexity of a naive agglomerative clustering is $O(n^3)$ because we exhaustively scan the $N \times N$ matrix `dist_mat` for the lowest distance in each of $N-1$ iterations

Increased performance: Multiple machines provide greater processing power. Greater scalability: As your user base grows and report complexity increases, your resources can grow. Simplified management: Clustering simplifies the management of large or rapidly growing systems.

divisive algorithm does not require to prespecify the number of clusters

Disadvantages:

A challenge with divisive methods is how to partition a large cluster into several smaller ones. For example, there are $2^n - 1$ possible ways to partition a set of n objects into two exclusive subsets, where n is the number of objects. When n is large, it is computationally prohibitive to examine all possibilities. Consequently, a divisive method typically uses heuristics in partitioning, which can lead to inaccurate results. For the sake of efficiency, divisive methods typically do not backtrack on partitioning decisions that have been made. Once a cluster is partitioned, any alternative partitioning of this cluster will not be considered again. Due to the challenges in divisive methods, there are many more agglomerative methods than divisive methods.

Hierarchical Agglomerative vs Divisive clustering:

Divisive clustering is more

complex as compared to agglomerative clustering, as in the case of divisive clustering we need a flat clustering method as “subroutine” to split each cluster until we have each data having its own singleton cluster.

Divisive clustering is more efficient if we do not generate a complete hierarchy all the way down to individual data leaves. The time complexity of a naive agglomerative clustering is $O(n^3)$ because we exhaustively scan the $N \times N$ matrix `dist_mat` for the lowest distance in each of $N-1$ iterations. Using priority queue data structure we can reduce this complexity to $O(n^2 \log n)$. By using some more optimizations it can be brought down to $O(n^2)$. Whereas for divisive clustering given a fixed number of top levels, using an efficient flat algorithm like K-Means, divisive algorithms are linear in the number of patterns and clusters.

A divisive algorithm is also more accurate. Agglomerative clustering makes decisions by considering the local patterns or neighbor points without initially taking into account the global distribution of data. These early decisions cannot be undone. whereas divisive clustering takes into consideration the global distribution of data when making top-level partitioning decisions.

Advantages of Hierarchical clustering

- Hierarchical clustering outputs a hierarchy, is a structure that is more informative than the unstructured set of flat clusters returned by k-means. Therefore, it is easier to decide on the number of clusters by looking at the dendrogram.
- Easy to implement.
- Unlike, K-means clustering where we set the number of clusters and then analyze the data and again change number of clusters to be formed for better results, here the algorithm is much easy for applying and interpretation.
- even after hierarchical clustering's more time complexity, it is worth to use it over other clustering technique

Disadvantages

- the time taken to apply this clustering is more than K-means. So, hierarchical is advised not to be used for large datasets, as the time complexity is $O(n^2)$, that is with larger dataset, time taken would take squared more time.
- It does not always provide the best solution: When you cluster multi-dimensional retail data that cannot always be visualized on a plot, poor solutions may be tricky to spot and resolve.
- The algorithm cannot run if there is missing data: You will need to remove these lines or estimate values to ensure the algorithm is able to run.
- The algorithm cannot run with different data types: When you use many different data types, it becomes difficult to compute a distance matrix. There is no simple formula that can work with both qualitative and numerical data at the same time.
- The dendrogram can be misinterpreted: The descriptors and composition of clusters may be difficult to interpret for all your business stakeholders involved with clustering.

Density-based clustering

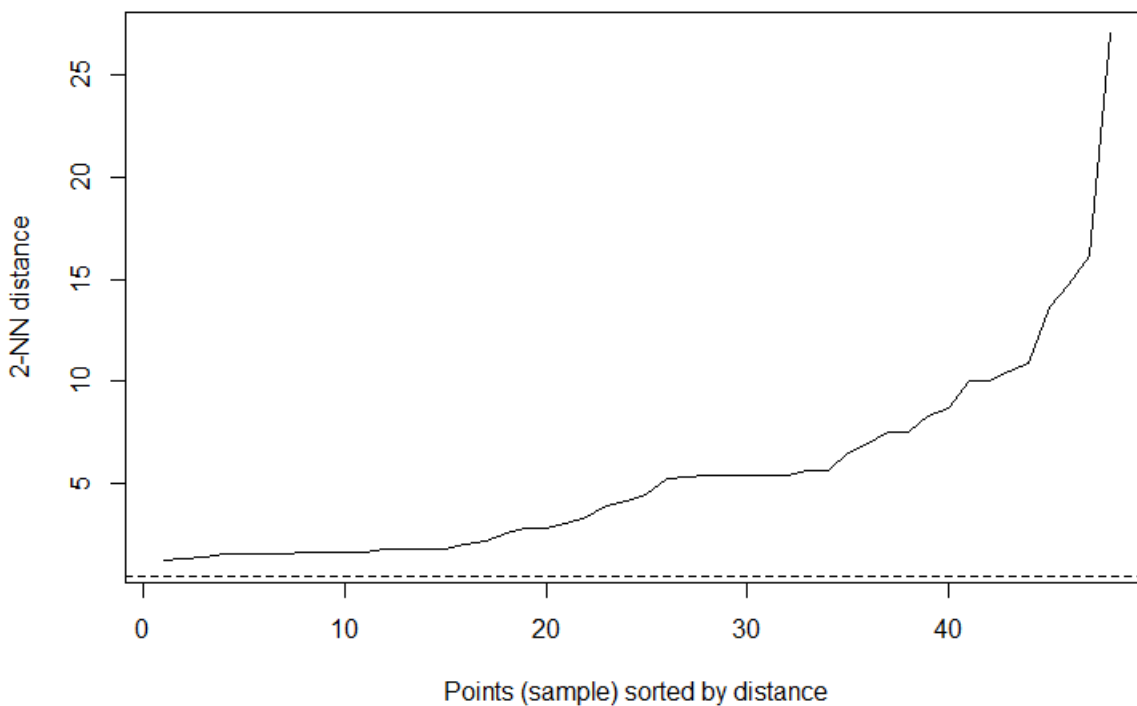
```
1 # Density-based clustering
2
3 #Importing libraries
4 library(dbscan)
5 library(fpc)
6 library(pch)
7
8 # data.2 data
9 data.2 ← read.csv("data 2.enc")
10 str(data.2)
11 D2← data.2[,c(2,3)]
12 str(D2)
13
14 # Installing packages
15 if(!require(devtools)) install.packages("devtools")
16 devtools::install_github("kassambara/factoextra")
17 library(factoextra)
18
19 # Obtaining optimal eps value
20 kMNdistsplot(D2, k=2)
21 abline(h = 0.45, lty=2)
22
23 # Density-based clustering with fpc & dbscan
24 set.seed(25)
25
26 f ← fpc::dbscan(D2, eps=15, MinPts =5)
27 f
28 d ← dbscan::dbscan(D2, 15, 5)
29 d
30
31 # Cluster visualization
32 fviz_cluster(d, D2, geom = "point")
33
34 # Plotting Cluster
35 plot(d, D2[,1:2], main = "DBScan")
36
```

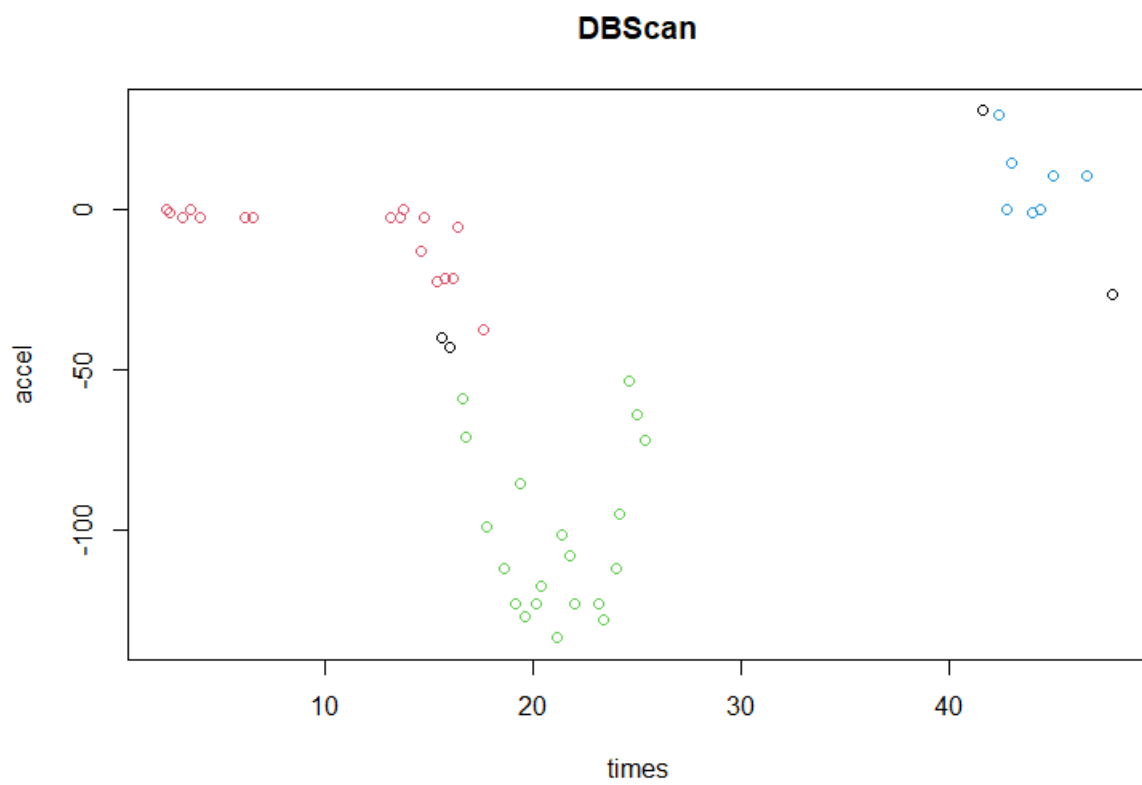
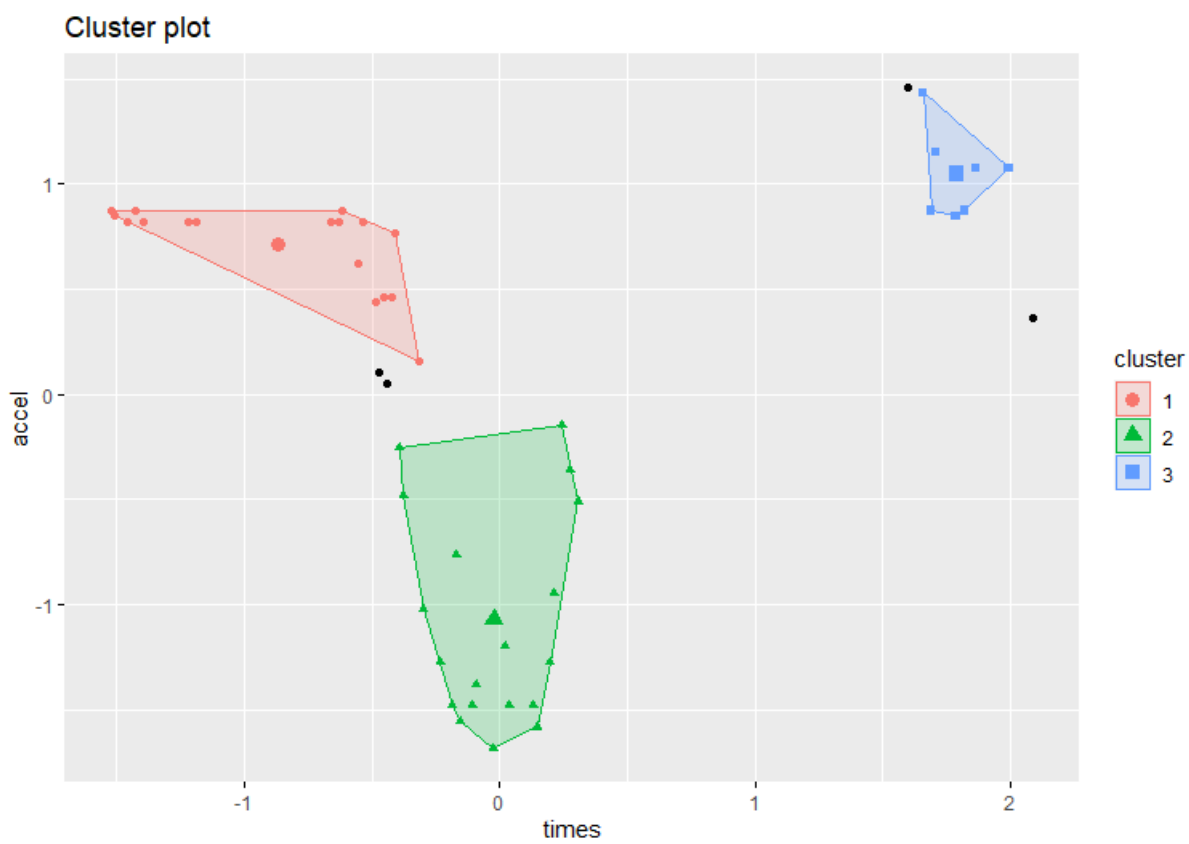
Results

```
DBSCAN clustering for 48 objects.  
Parameters: eps = 15, minPts = 5  
The clustering contains 3 cluster(s) and 4 noise points.
```

```
0  1  2  3  
4 17 20  7
```

```
Available fields: cluster, eps, minPts
```





Comment

To get the best result of clustering my data, I tried to increase and decrease the number of minPoints and EPS to see which solution is the suitable .

When

```
f <- fpc::dbscan(data.2, eps=20, MinPts = 5)
```

f

output:

dbscan Pts=48 MinPts=5 eps=20

0 1 2 3 4

border 5 0 3 1 2

seed 0 7 9 15 6

total 5 7 12 16 8

```
d <- dbscan::dbscan(data.2, 20, 5)
```

d

output:

DBSCAN clustering for 48 objects.

Parameters: eps = 20, minPts = 5

The clustering contains **4 cluster(s) and 5 noise points.**

0 1 2 3 4

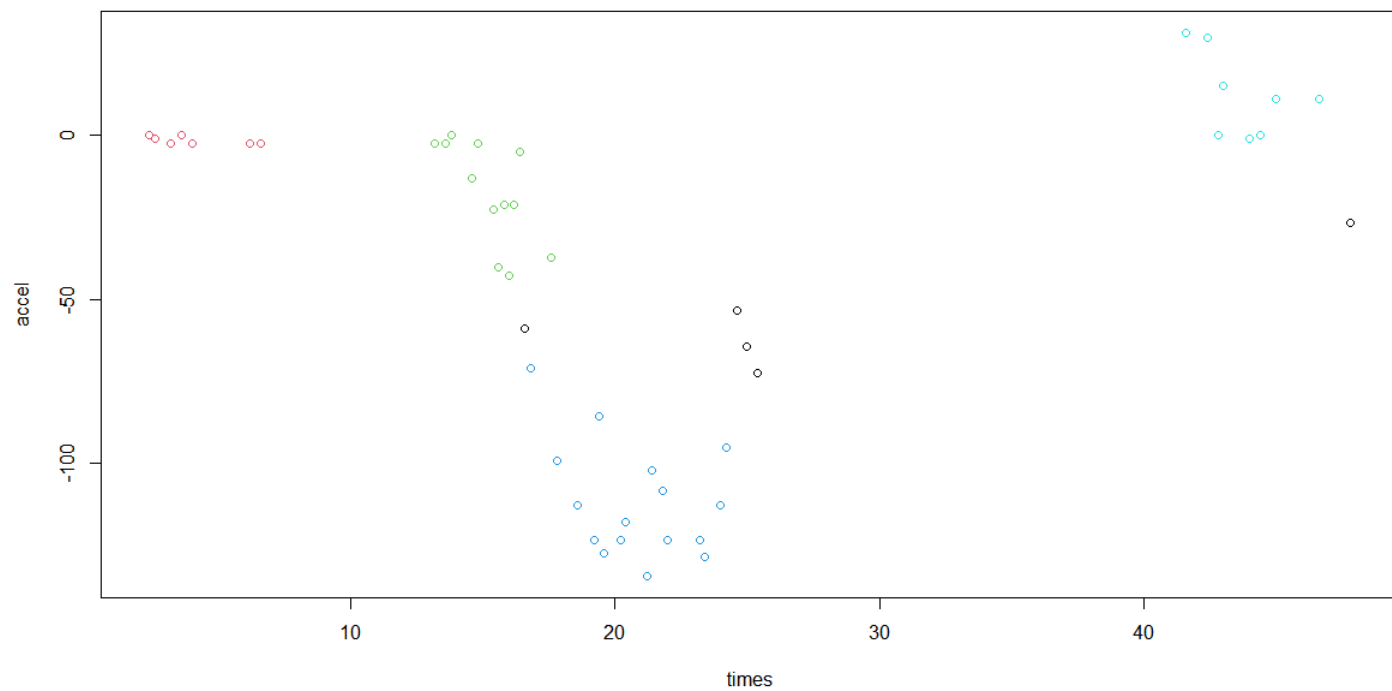
5 7 12 16 8

We notice that the number of noise points increased by 1 and the number of clusters increased to 4

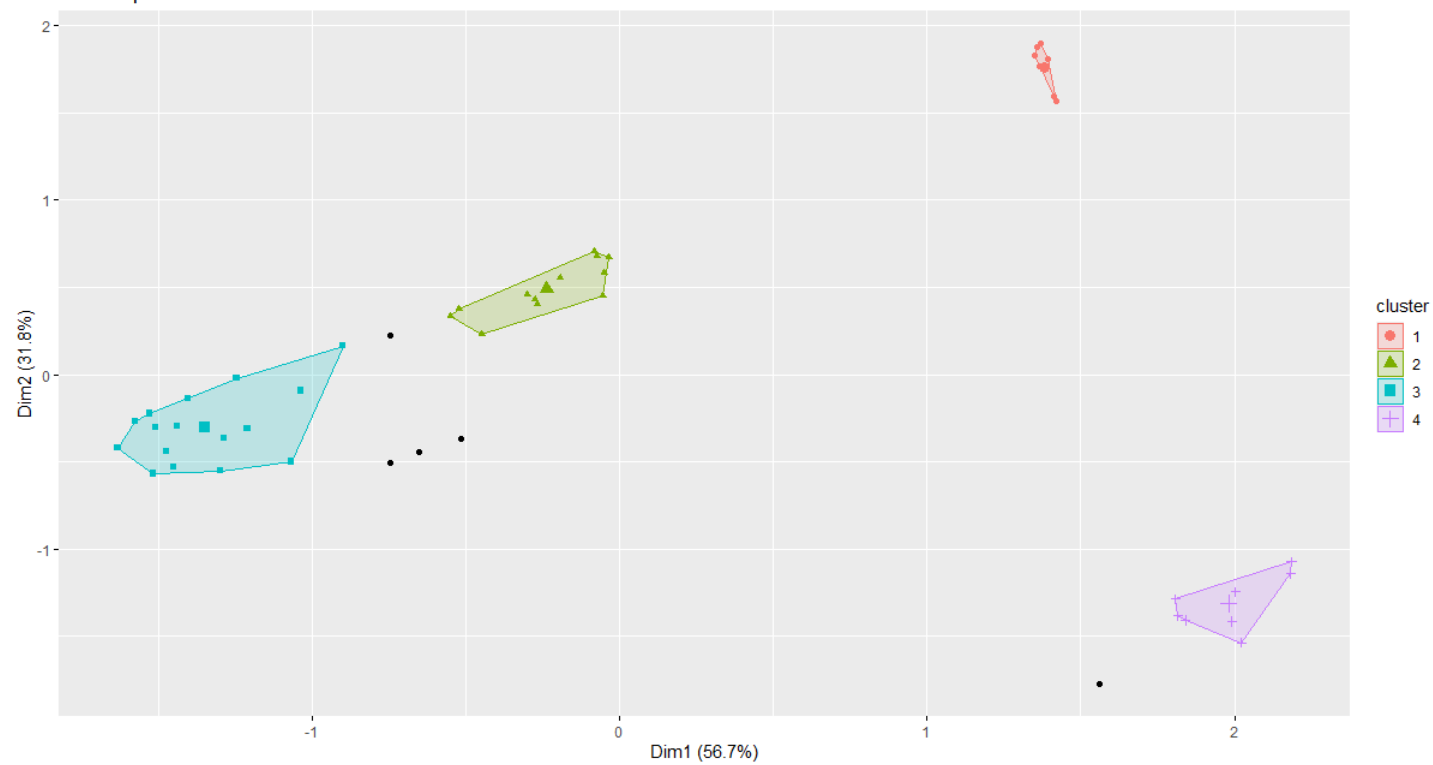
When using Eps =20 and minpoints =5

```
plot(d, D2[,1:2], main = "DBScan")
```

DBScan



Cluster plot



Advantages

- Can discover arbitrarily shaped clusters.
- Does not require a-priori specification of number of clusters.
- Able to identify noise data while clustering.
- Find cluster completely surrounded by different clusters.
- Require just two points which are very insensitive to the ordering of the points in the database.

Disadvantages

- Not partitionable for multiprocessor systems.
- Datasets with altering densities are tricky.
- Sensitive to clustering parameters minPoints and EPS.
- Fails to identify cluster if density varies and if the dataset is too sparse.
- Sampling affects density measures.

When this method is preferable >>

DBSCAN. DBSCAN is a clustering method that is used in machine learning to separate clusters of high density from clusters of low density. Given that DBSCAN is a density-based clustering algorithm, **it does a great job of seeking areas in the data that have a high density of observations, versus areas of the data that are not very dense with observations. DBSCAN can sort data into clusters of varying shapes as well, another strong advantage.**

Improve this method: when we increase Eps so that outliers decrease or the clusters increase and outliers increase , we get better normalization .