# Unsupervised Algorithms in machine learning [*]

**Cesar Conejo Villalobos**    *Data Scientist*

---

This document provides some examples of unsupervised algorithms in machine learning. In these techniques, we need to infer the properties of the observations without the help of an output variable or *supervisor*. We review two methods: k-means and hierarchical clustering. Then we use some data from Kaggle for applying these techniques to produce a customer segmentation. The platform that we use is R. Because of the number of observations, we are going to use a parallel process for improving the execution times.

*Keywords*: Unsupervised, algorithms, k-means, hierarchical classification, kaggle, R, parallel

---

**Introduction**

In the book *The Elements of Statistical Learning* Hastie T. (2008) explains that in the case of unsupervised learning, data usually has a set of $N$ observations $(x_1, ..., x_N)$ of a random vector $X$ having joint density $Pr(X)$. The goal is to infer the properties of this probability density.

The techniques that statistics and machine learning offer us for unsupervised learning are the following:

1) Principal components, multidimensional scaling.
2) Cluster analysis.
3) Mixture modeling.
4) Association rules.

In this exercise, we are going to focus on cluster analysis. The basis of our model will be the Kaggle Credit Card dataset for Clustering. The data are an 8950 x 18 matrix. One variable is categorical and represents the customer ID, the next seventeen are real numbers each representing the behavior of credit cardholders. The goal is to define a marketing strategy based on customer segmentation.

The first aspect we need to solve is to find the number of clusters that we need. Hastie T. (2008) says that we can have two scenarios:

1) For data segmentation, the number of clusters is defined as part of the problem, and it is base on the capacity and resources of the company. The goal is to find observations that belong to each proposed group.

2) Determine how the observations belong to natural distinct groupings. In this case, the number of clusters is unknown.

For this exercise, we are going to use scenario 2 and trying to find the number of clusters and the characteristics of each group using the following techniques:

1) k-means
2) Hierarchical clustering.

---

[*]Template taken from (http://github.com/svmiller). **Corresponding author**: svmille@clemson.edu.

**Techniques**

As mentioned before, the goal of unsupervised algorithms is to get the classes as homogeneous as possible and such that they are sufficiently separated. This goal can be specified numerically from the following property:

Suppose that exist a partition $P = (C_1, ..., C_K)$ of $\Omega$, where $g_1, ..., g_K$ are the cluster center of the classes:

$$g_k = \frac{1}{|C_k|} \sum_{i \in C_k} x_i$$

Also $g$ is the global center $\frac{1}{N} \sum_{i=1}^{N} x_i$. We also define:

- Total point scatter: $T = \frac{1}{N} \sum_{i=1}^{N} ||x_i - g||^2$.

- Within-cluster point scatter: $W(C) = \frac{1}{N} \sum_{k=1}^{K} \sum_{i \in C_k} ||x_i - g_k||^2$.

- Between-cluster point scatter: $B(C) = \sum_{k=1}^{K} \frac{|C_k|}{N} ||g_k - g||^2$

In this case, the algorithm requires $B(C)$ to be maximum and $W(C)$ to be minimum. Since the total point scatter $T$ is fixed, then maximizing $B(C)$ automatically minimizes $W(C)$. Therefore, the two goals (homogeneity within classes and separation between classes) are achieved at the same time by minimizing $W(C)$.

Thus, the goal in the *K-means* method is to find a partition $C$ of $\Omega$. Also, we find some representatives of the classes, such that $W(C)$ is minimal. For determining how many clusters a dataset has, we can use the elbow method.

Furthermore, k-means depend on the choice of the number of clusters. On the other hand, hierarchical clustering methods do not expect such designations. Instead, Hastie T. (2008) claims that this method demands the user to specify a measure of dissimilarity between (disjoint) groups of observations, based on the pairwise dissimilarities among the observations.

This method of classification uses a notion of proximity between groups of elements to measure the separation between the classes sought. To do this, the concept of aggregation is introduced, which is nothing more than a dissimilarity between groups of individuals: be $A, B \subseteq \Omega$ then the aggregation between $A$ and $B$ is $\delta(A, B)$. Then we have the following agglomerative clustering methods:

- Single linkage: $\delta_{SL}(A, B) = min\{d(x_i, d_j) | x_i \in A, x_j \in B\}$

- Complete linkage: $\delta_{CL}(A, B) = max\{d(x_i, d_j) | x_i \in A, x_j \in B\}$

- Average linkage: $\delta_{AL}(A, B) = \frac{1}{|A||B|} \sum_{x_i \in A, x_j \in B} d(x_i, d_j)$

- Ward linkage: $\delta_{Ward}(A, B) = \frac{|A||B|}{|A|+|B|} ||g_A - g_B||^2$

**Analysis**

We are going to define the marketing strategy using k-means and hierarchal clustering. But first, we will see the distribution of the data.

*Exploratory Analysis*

We create the function *cc_stats()* for analyzing some of the characteristics of the dataset such as:

- Number of complete observations.

- Number of NA values.

- Mean of complete observations.

- Standard desviation of complete observations.

- Number of outliers observations. $(Q3 + 1.5IQR)$

- Minimun value of complete observations.

- Maximun value of complete observations.

- 95 quantile.

- Upper limit for the value. (mean + 3 sd)

```r
---
# Basic statistics
# Input: x vector
# Output: Summary of statistics of the input

cc_stats <- function(x){

  #NA Values
  nas = sum(is.na(x))

  # Vector with complete values
  a = x[!is.na(x)]

  # Properties

  m   = mean(a)
  min = min(a)
  max = max(a)
  s   = sd(a)

  # Stats
  stats <- boxplot.stats(a)
  n     <- stats$n
  out   <- length(stats$out)

  Q95 = quantile(a, 0.95)
  UL = m + 3*s

  return(c(n      = n,
```

```
                nas    = nas,
                Mean   = m,
                StDev  = s,
                Q_out  = out,
                Min    = min,
                Max    = max,
                Q      = Q95,
                Upper_Limit = UL))
}
```

Using the function *apply()*, we see the statistical characteristics for each of the variables:

```
# Vector with the name of variables
vars <- c("BALANCE",
          "BALANCE_FREQUENCY",
          "PURCHASES",
          "ONEOFF_PURCHASES",
          "INSTALLMENTS_PURCHASES",
          "CASH_ADVANCE",
          "PURCHASES_FREQUENCY",
          "ONEOFF_PURCHASES_FREQUENCY",
          "PURCHASES_INSTALLMENTS_FREQUENCY",
          "CASH_ADVANCE_FREQUENCY",
          "CASH_ADVANCE_TRX",
          "PURCHASES_TRX",
          "CREDIT_LIMIT",
          "PAYMENTS",
          "MINIMUM_PAYMENTS",
          "PRC_FULL_PAYMENT",
          "TENURE")
```

```
# Apply the function for each variable
describe_stats <- t(data.frame(apply(cc_general[vars], 2, cc_stats)))
describe_stats
```

```
##                                      n nas    Mean   StDev Q_out    Min     Max
## BALANCE                           8950   0 1564.47 2081.53   695  0.000 19043.1
## BALANCE_FREQUENCY                 8950   0    0.88    0.24  1493  0.000     1.0
## PURCHASES                         8950   0 1003.20 2136.63   808  0.000 49039.6
## ONEOFF_PURCHASES                  8950   0  592.44 1659.89  1013  0.000 40761.2
## INSTALLMENTS_PURCHASES            8950   0  411.07  904.34   867  0.000 22500.0
## CASH_ADVANCE                      8950   0  978.87 2097.16  1030  0.000 47137.2
## PURCHASES_FREQUENCY               8950   0    0.49    0.40     0  0.000     1.0
## ONEOFF_PURCHASES_FREQUENCY        8950   0    0.20    0.30   782  0.000     1.0
## PURCHASES_INSTALLMENTS_FREQUENCY  8950   0    0.36    0.40     0  0.000     1.0
## CASH_ADVANCE_FREQUENCY            8950   0    0.14    0.20   525  0.000     1.5
## CASH_ADVANCE_TRX                  8950   0    3.25    6.82   804  0.000   123.0
## PURCHASES_TRX                     8950   0   14.71   24.86   766  0.000   358.0
```

```
## CREDIT_LIMIT                         8949    1 4494.45 3638.82    248 50.000 30000.0
## PAYMENTS                             8950    0 1733.14 2895.06    808  0.000 50721.5
## MINIMUM_PAYMENTS                     8637  313  864.21 2372.45    841  0.019 76406.2
## PRC_FULL_PAYMENT                     8950    0    0.15    0.29   1474  0.000     1.0
## TENURE                              8950    0   11.52    1.34   1366  6.000    12.0
##                                      Q.95% Upper_Limit
## BALANCE                             5909.11      7809.07
## BALANCE_FREQUENCY                      1.00         1.59
## PURCHASES                           3998.62      7413.11
## ONEOFF_PURCHASES                    2671.09      5572.10
## INSTALLMENTS_PURCHASES              1750.09      3124.08
## CASH_ADVANCE                        4647.17      7270.36
## PURCHASES_FREQUENCY                    1.00         1.69
## ONEOFF_PURCHASES_FREQUENCY             1.00         1.10
## PURCHASES_INSTALLMENTS_FREQUENCY       1.00         1.56
## CASH_ADVANCE_FREQUENCY                 0.58         0.74
## CASH_ADVANCE_TRX                      15.00        23.72
## PURCHASES_TRX                         57.00        89.28
## CREDIT_LIMIT                       12000.00     15410.90
## PAYMENTS                            6082.09     10418.34
## MINIMUM_PAYMENTS                    2766.56      7981.55
## PRC_FULL_PAYMENT                       1.00         1.03
## TENURE                                12.00        15.53
```

First of all, there is only a few values with `NA`. If we want to see if they both happen at the same time, we can do:

```
sum(is.na(cc_general$CREDIT_LIMIT) & is.na(cc_general$MINIMUM_PAYMENTS))
```

```
## [1] 0
```

As a result, the `NA` values do not occur in the same row. For fixing these unknown values, we can follow three alternatives:

- Remove the cases.

- Fill in the unknowns using some strategy.

- Use tools that handle these types of values.

In this case, the unknown values only represent 3.51% of the data, so we decide to delete the observations.

```
cc_general <- cc_general[-which(is.na(cc_general$CREDIT_LIMIT)
                        | is.na(cc_general$MINIMUM_PAYMENTS)),]
```

Another aspect we see is that the variables are measure in different scales. For example *BALANCE FREQUENCY*, *PURCHASES FREQUENCY*, *ONE OFF PURCHASES FREQUENCY* and *PURCHASES INSTALLMENTS FREQUENCY* are measure with a score between 0 and 1. Other values are measure in money units and others in the number of transactions. Because there are different units then we should scaling variables. We do that with the function ***normalize()***:

```
---
## Normalize
## Input: Numeric vector
## Output: Vector normalized.

normalize <- function(x){

  min_x <- min(x)
  max_x <- max(x)

  return((x - min_x)/(max_x - min_x))

}
```

Then, we apply the function to each variable:

```
cc_general_norm <- data.frame(apply(cc_general[vars], 2, normalize))
```

Finally, we apply *cc_stats()* again for seeing the changes in our data.

```
describe_stats_norm <- t(data.frame(apply(cc_general_norm[vars], 2, cc_stats)))
describe_stats_norm
```

```
##                                      n nas  Mean StDev Q_out Min Max Q.95%
## BALANCE                           8636   0 0.084 0.110   666   0   1 0.312
## BALANCE_FREQUENCY                 8636   0 0.895 0.208  1511   0   1 1.000
## PURCHASES                         8636   0 0.021 0.044   767   0   1 0.083
## ONEOFF_PURCHASES                  8636   0 0.015 0.041   961   0   1 0.067
## INSTALLMENTS_PURCHASES            8636   0 0.019 0.041   811   0   1 0.080
## CASH_ADVANCE                      8636   0 0.021 0.045   976   0   1 0.100
## PURCHASES_FREQUENCY               8636   0 0.496 0.401     0   0   1 1.000
## ONEOFF_PURCHASES_FREQUENCY        8636   0 0.206 0.300   749   0   1 1.000
## PURCHASES_INSTALLMENTS_FREQUENCY  8636   0 0.369 0.398     0   0   1 1.000
## CASH_ADVANCE_FREQUENCY            8636   0 0.092 0.135   346   0   1 0.389
## CASH_ADVANCE_TRX                  8636   0 0.027 0.056   794   0   1 0.122
## PURCHASES_TRX                     8636   0 0.042 0.070   716   0   1 0.165
## CREDIT_LIMIT                      8636   0 0.149 0.122   243   0   1 0.399
## PAYMENTS                          8636   0 0.035 0.057   785   0   1 0.121
## MINIMUM_PAYMENTS                  8636   0 0.011 0.031   841   0   1 0.036
## PRC_FULL_PAYMENT                  8636   0 0.159 0.296  1343   0   1 1.000
## TENURE                            8636   0 0.922 0.218  1290   0   1 1.000
##                                  Upper_Limit
## BALANCE                                 0.41
## BALANCE_FREQUENCY                       1.52
## PURCHASES                               0.15
## ONEOFF_PURCHASES                        0.14
## INSTALLMENTS_PURCHASES                  0.14
## CASH_ADVANCE                            0.16
```

```
## PURCHASES_FREQUENCY                      1.70
## ONEOFF_PURCHASES_FREQUENCY               1.11
## PURCHASES_INSTALLMENTS_FREQUENCY         1.56
## CASH_ADVANCE_FREQUENCY                    0.50
## CASH_ADVANCE_TRX                          0.20
## PURCHASES_TRX                             0.25
## CREDIT_LIMIT                              0.52
## PAYMENTS                                  0.21
## MINIMUM_PAYMENTS                          0.10
## PRC_FULL_PAYMENT                          1.05
## TENURE                                    1.58
```

Now, all the variables are on a scale from 0 to 1. Also, there is no change in the variance of the variables because we are only scaling and no standardizing. Finally, we save the new data set for being used as the source in the parallel executions:

```
#Save table
write.table(x = cc_general_norm,
            file = "output/data/cc_general_norm.csv",
            sep = ",",
            dec = ".")
```

### K-means

For applying *k-means*, we develop the following code. We set the seed 1234 for reproducibility purposes.

```
set.seed(1234)
```

First, we need to decide the number of clusters. We apply K-means clustering to the data using the following techniques:

- Hartigan-Wong

- MacQueen

- Lloyd

- Forgy

Also, we use the snow library for performing parallelizable operations.

```
library(snow)

cl <- makeCluster(4, type="SOCK")
# Read data in each cluster
ignore <- clusterEvalQ(cl,
                       { data <- read.csv("output/data/cc_general_norm.csv",
                                          header = TRUE,
```

```r
                                          sep = ",",
                                          dec = "."
                                          )
                          set.seed(1234)
                          NULL})

# Hartigan-Wong
results_HW <- clusterApply(cl,
                           seq(1,20),
                           function(x) kmeans(data,
                                              centers = x,
                                              algorithm = "Hartigan-Wong",
                                              nstart = 50))

withinss_HW <-  sapply(results_HW, function(results_HW) results_HW$tot.withinss)

# MacQueen
results_MQ <- clusterApply(cl,
                           seq(1,20),
                           function(x) kmeans(data,
                                              centers = x,
                                              algorithm = "MacQueen",
                                              nstart = 50))

withinss_MQ <-  sapply(results_MQ, function(results_MQ) results_MQ$tot.withinss)

# Lloyd
results_Ll <- clusterApply(cl,
                           seq(1,20),
                           function(x) kmeans(data,
                                              centers = x,
                                              algorithm = "Lloyd",
                                              nstart = 50))

withinss_Ll <-  sapply(results_Ll, function(results_Ll) results_Ll$tot.withinss)


# Forgy
results_FG <- clusterApply(cl,
                           seq(1,20),
                           function(x) kmeans(data,
                                              centers = x,
                                              algorithm = "Forgy",
                                              nstart = 50))

withinss_FG <-  sapply(results_FG, function(results_FG) results_FG$tot.withinss)
```

```
# end paralell
stopCluster(cl)
```

Then, we can observe the total within-cluster sum of squares for K-means clustering for some clusters from 1 to 20.

```
plot(withinss_HW,
     col = "red",
     type = "b",
     xlab = "Number of cluster k",
     ylab = "Sum of squares",
     main = "Elbow Method: No. of clusters by algorithm")
points(withinss_MQ, col = "blue",    type = "b")
points(withinss_Ll, col = "green",   type = "b")
points(withinss_FG, col = "magenta", type = "b")
legend("topright",
       legend = c("Hartigan","MacQueen","Lloyd","Forgy"),
       col = c("red", "blue", "green", "magenta"),
       lty = 1,
       lwd = 1)
```



Figure 1: Total within-cluster sum of squares for K-means clustering

We can see that the kink occurs at $k = 5$, so this is the number of clusters that we propose for the marketing strategy. Now, the question that we need to answer is which method for k-means to use. Therefore, we code the following lines that show the results of the clustering in 5 groups.

9

```r
# Which method:
cl <- makeCluster(4, type="SOCK")
# Read data in each cluster
ignore <- clusterEvalQ(cl,
                       { data <- read.csv("output/data/cc_general_norm.csv",
                                          header = TRUE,
                                          sep = ",",
                                          dec = "."
                       )
                       set.seed(1234)
                       NULL})


# Evaluate all the posible algorithms
results_Algorithm <- clusterApply(cl,
                         c("Hartigan-Wong","MacQueen","Lloyd","Forgy"),
                         function(algorithm) kmeans(data,
                                            centers = 5,
                                            algorithm = algorithm,
                                            nstart = 50))
for(i in 1:4) print(results_Algorithm[[i]]$betweenss)
# end cluster
stopCluster(cl)
```

```
## [1] 3496
## [1] 3496
## [1] 3490
## [1] 3490
```

Because *Hartigan-Wong* reaches the greatest (Also MacQueen does) between-cluster point scatter value, we choose that method. As a result, we code the method in the following way:

```r
Cluster_HW <- kmeans(x = cc_general_norm,
                     centers =  5,
                     nstart = 250,
                     algorithm = "Hartigan-Wong")
```

Finally, we can see some properties of the cluster. For example, the number of observations by the cluster is:

```r
Cluster_HW$size
```

```
## [1] 1169  823 3510 2106 1028
```

In this case, we can see that the classes are well balanced. Also, we can see the center of each cluster and realize some interpretations.
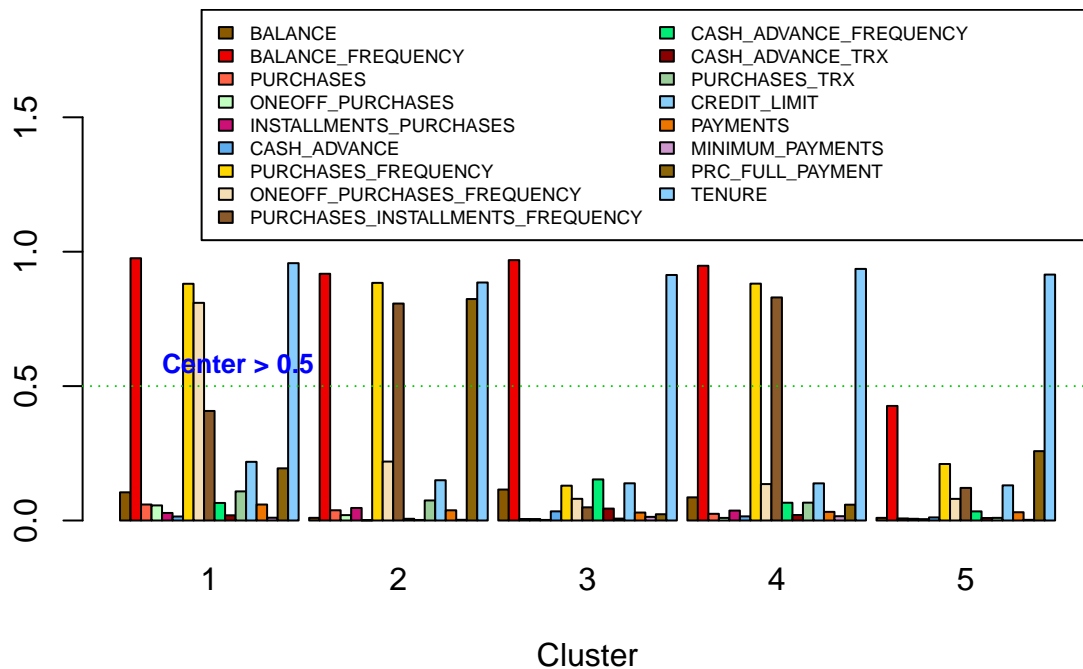
```
palette <- c("orange4", "red2", "tomato", "darkseagreen1",  "deeppink3", "steelblue2",
             "gold1", "wheat", "tan4", "springgreen2", "darkred", "darkseagreen3",
             "lightskyblue", "darkorange2", "plum3", "darkgoldenrod4", "skyblue1")

# Graph
barplot(t(Cluster_HW$centers),
        main = "Centers by cluster K-means clustering",
        xlab = "Cluster",
        beside = TRUE,
        col = palette,
        ylim = c(0, 1.9)
        )
abline(h = 0.5, col = 3, lty = 3)
text(5, 0.55, labels= "Center > 0.5", col = "blue", adj = c(0, 0), font=2, cex= 0.8 )
legend("topright",
       legend = colnames(Cluster_HW$centers),
       fill = palette, ncol = 2,
       cex = 0.6)
```

**Centers by cluster K−means clustering**



In cluster 1, these customers have the highest value in ONEOFF_PURCHASES_FREQUENCY. So, this group usually does the greatest purchase amount done in one-go.

Cluster 2 and 4 have the highest value in purchases_installments_frequency, so these customers used frequently their credit cards. But cluster 2 has the highest value in PRC_FULL_-PAYMENT, so this is a cluster of customers that used their card and also pay the products.

Cluster 3 has customers with less Purchase frequency, so this group does not make frequent purchases with their credit cards. Finally, cluster 5 has customers with less balance, balance frequency, and purchase installments frequency. This group is similar to group 3 because this is the group of customers that do not use frequently their credit cards, so their balance is not frequently updated.

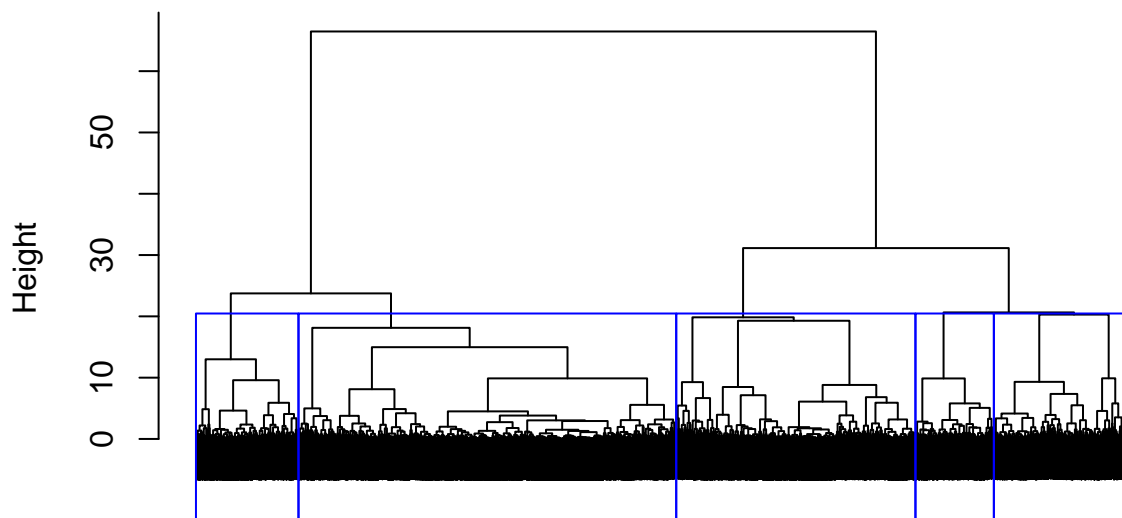In this way, we can segment the marketing strategy among these groups.

1) Customers with a high amount of purchases.
2) Customers that used their cards frequently.
3) Customers that no use their cards frequently.

*Hierachical cluster*

Using **Ward** aggregation, we code the method in the following way:

```
model_Ward <- hclust(dist(cc_general_norm), method = "ward.D2")

plot(model_Ward, labels = FALSE, xlab = "Hierachical cluster using Ward")
rect.hclust(model_Ward, k = 5, border = "blue")
```

## Cluster Dendrogram



Hierachical cluster using Ward
hclust (*, "ward.D2")

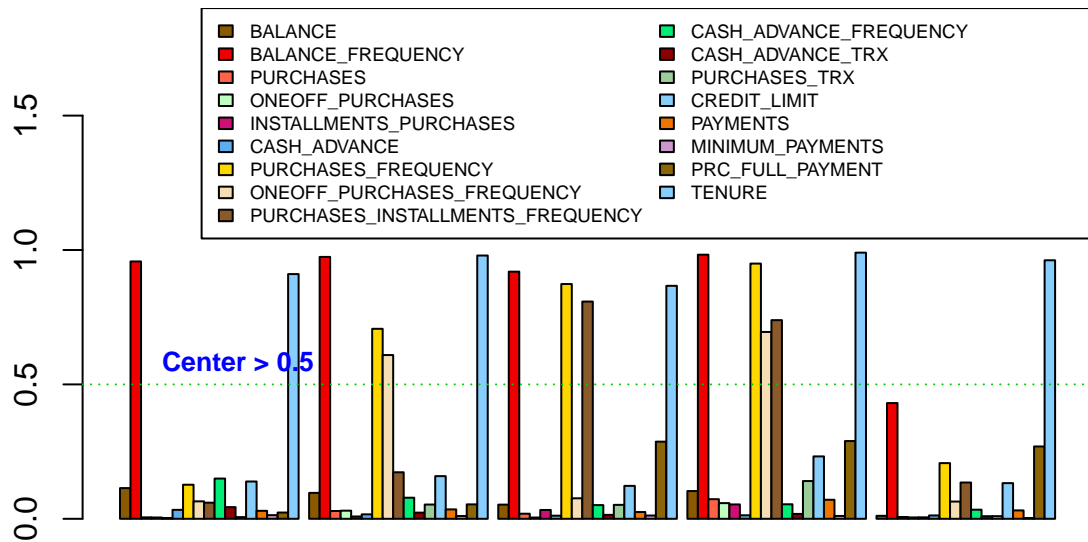```
cluster <- cutree(model_Ward, k = 5)
cc_general_norm_cluster <-cbind(cc_general_norm, cluster)
library(rattle)
```

```
cc_general_center <- centers.hclust(cc_general_norm,
                                    model_Ward,
                                    nclust = 5,
                                    use.median = FALSE)
```

As in *k-means*, we can see the properties of the centers cluster in the following graph:

```
barplot(t(cc_general_center),
        beside = TRUE,
        main = "Center by cluster Hierarchical Classification",
        col = palette,
        ylim = c(0, 1.9)
        )
abline(h = 0.5, col = 3, lty = 3)
text(5, 0.55, labels= "Center > 0.5", col = "blue", adj = c(0, 0), font=2, cex= 0.8 )
legend("topright",
       legend = colnames(cc_general_center),
       fill = palette, ncol = 2,
       cex = 0.6)
```



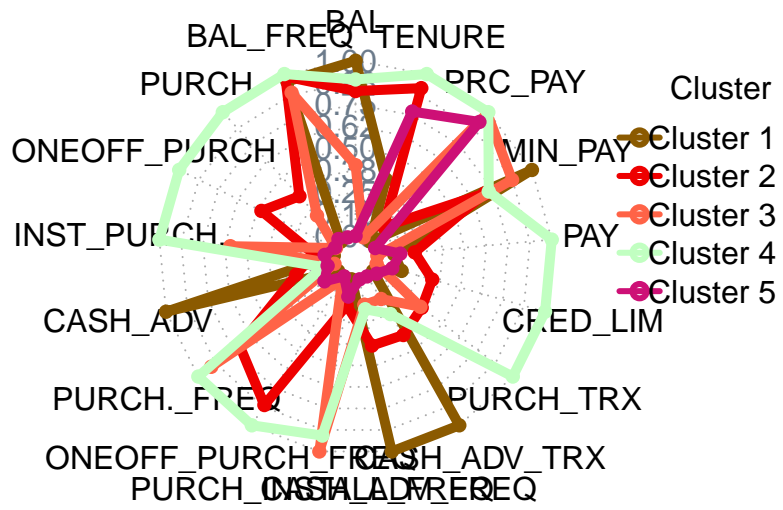For interpretation of the results, we can use the radar chart

- In cluster 1, the most important variables are balance and cash advance.
- In cluster 2 the most important variable is purchase frequency.

- In cluster 3 the most important variable is One-off purchase frequency.
- Cluster 4 has the highest influence on several of the variables.
- Cluster 5 has the influence of the variable PRC_Payments.

```
# radar graph
center   <- as.data.frame(cc_general_center)
maximos <- apply(center,2,max)
minimos <- apply(center,2,min)
center   <- rbind(minimos,center)
center   <- rbind(maximos,center)
colnames(center) <- c("BAL", "BAL_FREQ", "PURCH", "ONEOFF_PURCH","INST_PURCH.",
                      "CASH_ADV", "PURCH._FREQ", "ONEOFF_PURCH_FREQ",
                      "PURCH_INSTALL_FREQ","CASH_ADV_FREQ", "CASH_ADV_TRX",
                      "PURCH_TRX", "CRED_LIM", "PAY", "MIN_PAY", "PRC_PAY", "TENURE")
```

```
library(fmsb)
radarchart(center,
           maxmin = TRUE,
           axistype = 4,
           axislabcol = "slategray4",
           centerzero = FALSE,
           seg = 8,
           cglcol = "gray67",
           pcol= palette,
           plty = 1,
           plwd = 5,
           title = "Cluster comparation Hierarchical Classification")
det_radar   <-legend(1.5,1, legend=c("Cluster 1","Cluster 2","Cluster 3",
                                     "Cluster 4", "Cluster 5"),
               seg.len=-1.4,
               title="Cluster",
               pch=21,
               bty="n" ,lwd=3, y.intersp=1, horiz=FALSE,
               col= palette
               )
```

**Cluster comparation Hierarchical Classification**



Figure labels: BAL, BAL_FREQ, TENURE, PURCH, PRC_PAY, ONEOFF_PURCH, MIN_PAY, INST_PURCH, PAY, CASH_ADV, CRED_LIM, PURCH._FREQ, PURCH_TRX, ONEOFF_PURCH_FREQ, CASH_ADV_TRX, PURCH_INSTALL_FREQ, CASH_ADV_FREQ

Legend: Cluster — Cluster 1, Cluster 2, Cluster 3, Cluster 4, Cluster 5

Scale values: 0.13, 0.25, 0.38, 0.50, 0.62, 0.75, 0.87, 1.00

## References

Hastie T., Tibshirani R., Friedman J. 2008. *The elements of Statistical Learning*. Springer.