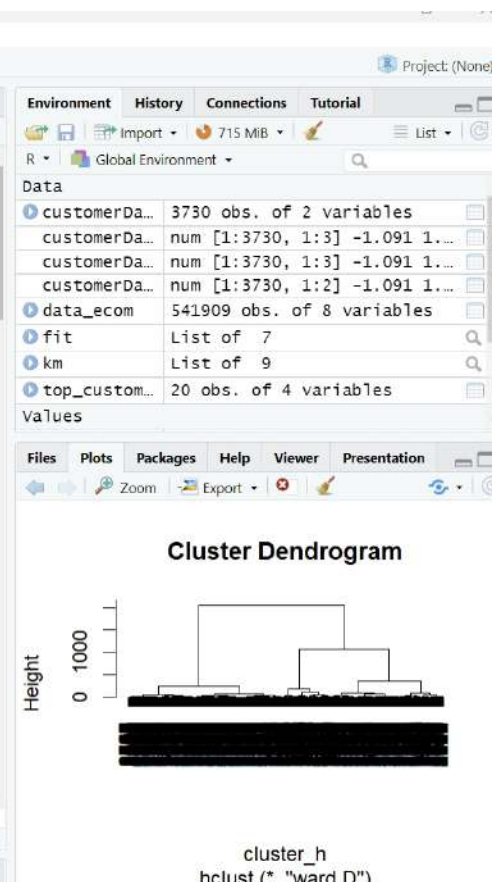


```
File Edit Code View Plots Session Build Debug Profile Tools Help
Go to file/function Addins
Ecommerce Clustering Project.R x customerData_Hier x customerData_km x customerData_sc x top_customers x custo...
Source on Save Run Source
1 #Selecting working Directory in my pc
2 setwd(choose.dir())
3 #Confirming the working directory path where Ecommrce Data Lies
4 getwd()
5
6 #install.packages("factoextra")--->example of installing necessary Libraries)
7 #Loading necessary Libraries
8 library(dplyr)
9 library(ggplot2)
10 library(cluster)
11 library(factoextra)
12 library(gridExtra)
13 library(purrr)
14
15 #Importing dataset from current working directory
16 data_ecom = read.csv("Ecommerce.csv")
17
18 #Viewing Dataset
19 View(data_ecom)
20 #Viewing first 6 rows of dataset in the console
21 head(data_ecom)
22
23 #Structure of the dataset
24 str(data_ecom)
25 #There is unnecessary column X with NA Values, Lets drop it
26 data_ecom <- subset(data_ecom,select=-X)
27
28 #We also saw date data in the character format,lets change it to Date format
29 data_ecom$InvoiceDate <- as.Date(data_ecom$InvoiceDate, "%d-%B-%y")
30
31 #Lets recheck the columns and their structure, and view it again
32 str(data_ecom)
33 View(data_ecom)
34
132:1 (Top Level) R Script
Console
```



File Edit Code View Plots Session Build Debug Profile Tools Help

Go to file/function Addins

```
Ecommerce Clustering Project R customerData_Hier customerData_km customerData_sc top_customers custoi
Source on Save Run Source
34 # We can still see the NA values in the data, Lets see its weight
35 na_count <- colSums(is.na(data_ecom))
36 na_count
37
38 # We have 135,080 missing Customer ID which is unique identifier.
39 # At this point, it's better to omit the data as we can't tag this transaction to right specific c
40 na.omit(data_ecom)
41
42 #Summary statistics of the dataset
43 summary(data_ecom)
44
45 #TotalSales & TotalUnitSales purchase by each customer
46 customerData <- data_ecom %>%
47   select(CustomerID, Country, Quantity, UnitPrice) %>%
48   group_by(CustomerID, Country) %>%
49   summarise(TotalRevenue = sum(Quantity * UnitPrice),
50             TotalItemsSold = sum(Quantity),
51             .groups = 'drop')
52 View(customerData)
53
54 ##TOPHighly Valued Customers based on Revenue Generated
55 top_customers <- customerData %>%
56   arrange(desc(TotalRevenue)) %>%
57   slice(1:20)
58 View(top_customers)
59
60
61 #We remove CustomerID and Country Columns as these features are of Less Importance/ can result mu
62 customerData <- customerData[,-c(1:2)]
63
64 #Finding and Sorting Outliers
65 boxplot(customerData)
66
67
132:1 (Top Level) R Script
```

Console

Project: (None)

Environment History Connections Tutorial

Global Environment 715 MiB

Data

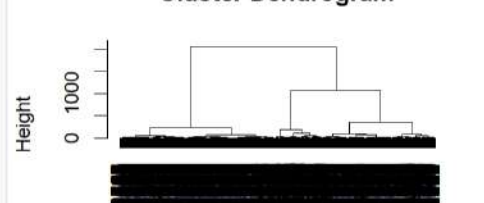
customerDa...	3730 obs. of 2 variables
customerDa...	num [1:3730, 1:3] -1.091 1...
customerDa...	num [1:3730, 1:3] -1.091 1...
customerDa...	num [1:3730, 1:2] -1.091 1...
data_ecom	541909 obs. of 8 variables
fit	List of 7
km	List of 9
top_custom...	20 obs. of 4 variables

Values

Files Plots Packages Help Viewer Presentation

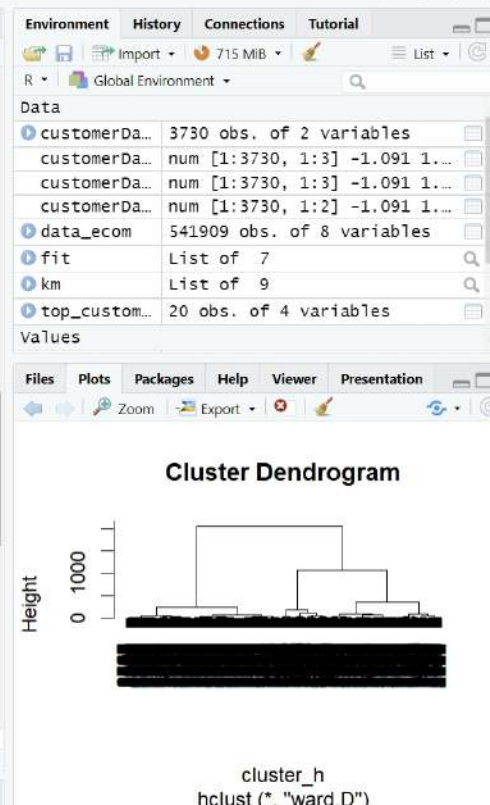
Zoom Export

Cluster Dendrogram



cluster_h
hclust (*, "ward.D")

```
File Edit Code View Plots Session Build Debug Profile Tools Help
Go to file/function Addins
Ecommerce Clustering Project.R x customerData_Hier x customerData_km x customerData_sc x top_customers x custo...
Source on Save Run Source
66
67 #Eliminating the outliers
68 library(dplyr)
69
70 customerData <- customerData %>%
71   filter(between(TotalRevenue, quantile(TotalRevenue, 0.25) - 1.5 * IQR(TotalRevenue),
72     quantile(TotalRevenue, 0.75) + 1.5 * IQR(TotalRevenue))) %>%
73   filter(between(TotalItemsSold, quantile(TotalItemsSold, 0.25) - 1.5 * IQR(TotalItemsSold),
74     quantile(TotalItemsSold, 0.75) + 1.5 * IQR(TotalItemsSold)))
75
76 #Lets Normalize the data as Sales Amount and Sales Units value has huge diff to compare without sc
77 customerData_sc <- scale(customerData)
78 View(customerData_sc)
79
80 ##KMeans Clustering
81 set.seed(42)
82
83 #Let's group the data into 5 cluster with k-means
84 km <- kmeans(customerData_sc, centers=5)
85 km
86 km$size
87
88 #Determining optimal clusters with Elbow Mthod
89 set.seed(42)
90
91 # Compute total within-cluster sum of squares for k=1 to 10
92 wss_values <- map_dbl(1:10, ~kmeans(customerData_sc, ., nstart = 10)$tot.withinss)
93
94 # Plot WSS values vs. number of clusters
95 plot(1:10, wss_values, type = "b", pch = 19, frame = FALSE,
96   xlab = "Number of clusters (k)", ylab = "Total within-cluster sum of squares (WSS)")
97
98 ###Looking at the plot we can say as per our K means graph 3 is the optimal number of clusters.
99
```



File Edit Code View Plots Session Build Debug Profile Tools Help

Go to file/function Addins

Project: (None)

Environm R Global

Data

C... 3...
C... n...
C... n...
C... n...
d... 5...
f... L...
km L...
t... 2...

Values

Files P

85:3 (Top Level) R Script

```
99 km <-kmeans(customerData_sc,centers=3)
100 km
101 km$size
102 str(km)
103
104 #Creating new column to specify the cluster
105 customerData_km <- cbind(customerData_sc,ClusterNum = km$cluster)
106 View(customerData_km)
107 #So, with k =3 , KMeans model is 81.6% confident to group new data to into distinct cluster.
108
109 ##Hierarchical Clustering
110 cluster_h <- dist(customerData_sc,method = "euclidian")
111 fit <- hclust(cluster_h,method = "ward.D")
112 groups <- cutree(fit, k = 3)
113 groups_factor <- as.factor(groups)
114 customerData_Hier <- cbind(customerData_sc, ClusterNum = groups_factor)
115 View(customerData_Hier)
116
117 # Dendrogram
118 plot(fit)
119
120 # Silhouette analysis
121 silhouette <- silhouette(groups_factor, dist(customerData_Hier))
122 silhouette_df <- summary(silhouette)
123 silhouette_df$avg.width
124
125
126 ###Identify the clustering algorithm that gives maximum accuracy and explains robust clusters.
127 #silhouette_df$avg.width is '0.499' which is positive and close to 1.
128 #Hierarchical model could be furthered improved by changing the value of K.
129 #KMeans model with value of k=3 gives maximum accuracy, and is 81.6% confident to cluster new data to the similar clusters.
130 #Since, Kmeans of 3 provides the better customer clustering and should be considered to roll out the loyalty program.
131 View(customerData_km)
132
```

e margins too