

# NEWFINAL

Lizbeth, Nicole, Jenny, Anabel

2025-11-04

```
#set working directory
setwd("/cloud/project")
###Importing my original data set to RStudio and calling it MCE
MCE<-read.csv("Mall_Customers_extended.csv", header=TRUE)

#to view variable names
names(MCE)

## [1] "Unnamed..0"          "CustomerID"          "Genre"
## [4] "Age"                 "Annual.Income..k.."  "Spending.Score..1.100."
## [7] "IncomePerAge"        "SpendingEfficiency"  "AgeCategory"
## [10] "HighSpender"         "IncomeTier"          "OnlineShopFreq"
## [13] "LoyaltyScore"        "Satisfaction"        "CreditUtilization"

##### Step 1: Perform a K-Means Cluster analysis #####
# I created MCEcluster1 data frame with the only 3 variables needed
# for my cluster: annual income, spending score, age

MCEcluster<-data.frame(MCE)

MCEcluster1<-MCEcluster[,c("Annual.Income..k..", "Spending.Score..1.100.", "Age")]

#Now, I'm going to create a new file called MCEsv with my 3 values standardized
#so that I can begin running my clusters
MCEcluster2<-data.frame(MCEcluster1)

MCEsv <- scale(MCEcluster2)

#View first 5 observations
head(MCEsv,n=5)

##      Annual.Income..k.. Spending.Score..1.100.      Age
## [1,]          -1.734646          -0.4337131 -1.4210029
## [2,]          -1.734646           1.1927111 -1.2778288
## [3,]          -1.696572          -1.7116178 -1.3494159
## [4,]          -1.696572           1.0378135 -1.1346547
## [5,]          -1.658498          -0.3949887 -0.5619583

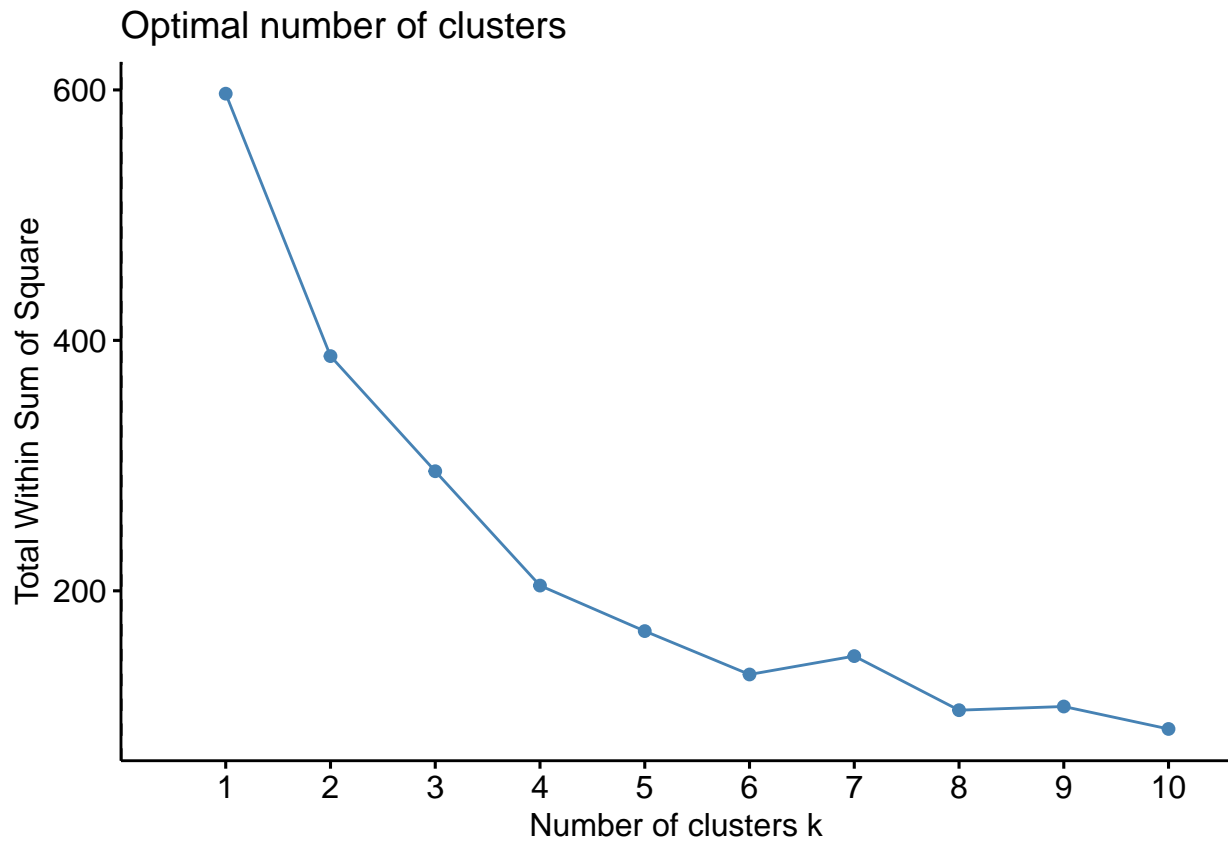
##### K-Means Algorithm #####
#The factoextra package creates clusters in R studio

install.packages("factoextra")
library(factoextra)
```

```
install.packages("rstatix")  
library(rstatix)
```

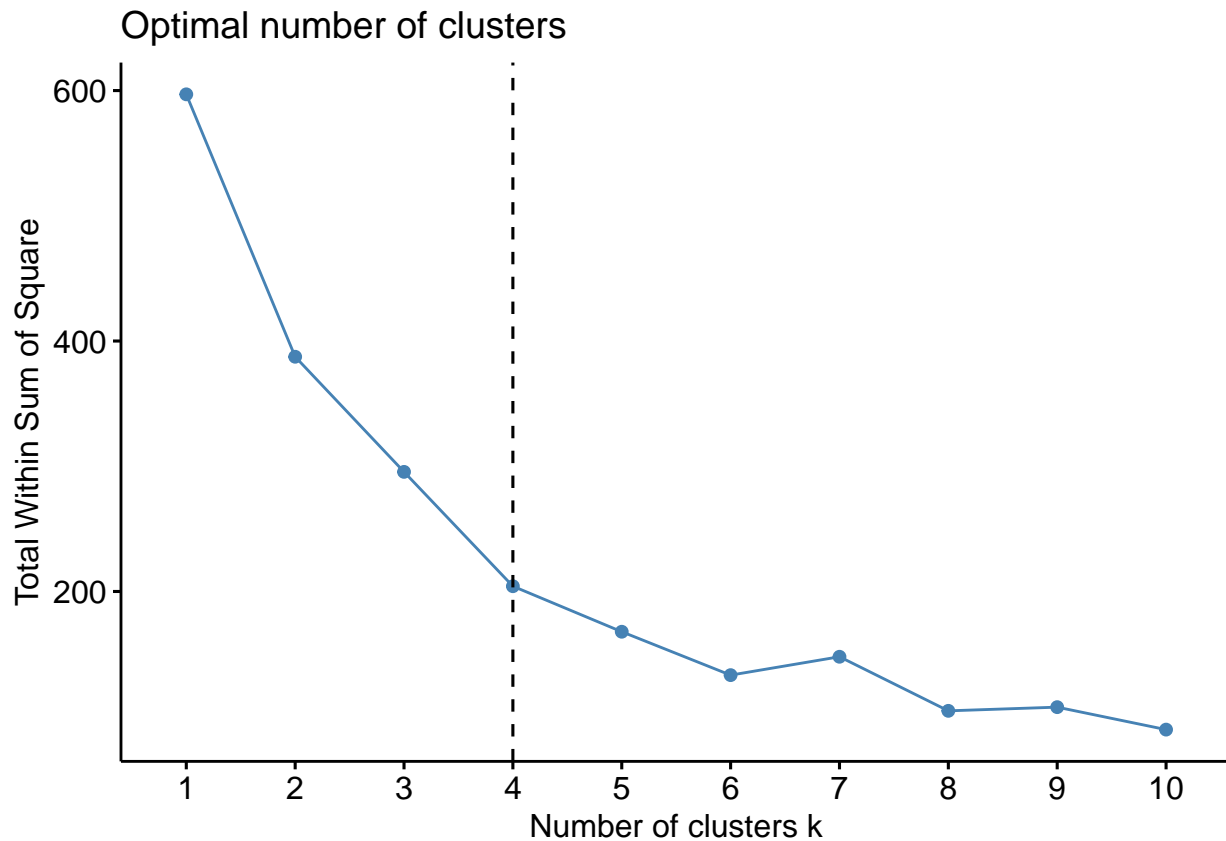
*#To find the number of clusters needed we use the fviz\_nbclust function*

```
fviz_nbclust(MCEsv, kmeans, method="wss") + geom_vline(xintercept = 0, linetype = 2)
```



*## I'm using 4 clusters based on the graph in my plots box*

```
fviz_nbclust(MCEsv, kmeans, method="wss") + geom_vline(xintercept = 4, linetype = 2)
```



```
#to obtain descriptive stats on 4 clusters
set.seed(123)

km.res <- kmeans(MCEsv, 4, nstart=25)
##Per homework instructions, do not run line 52
print(km.res)

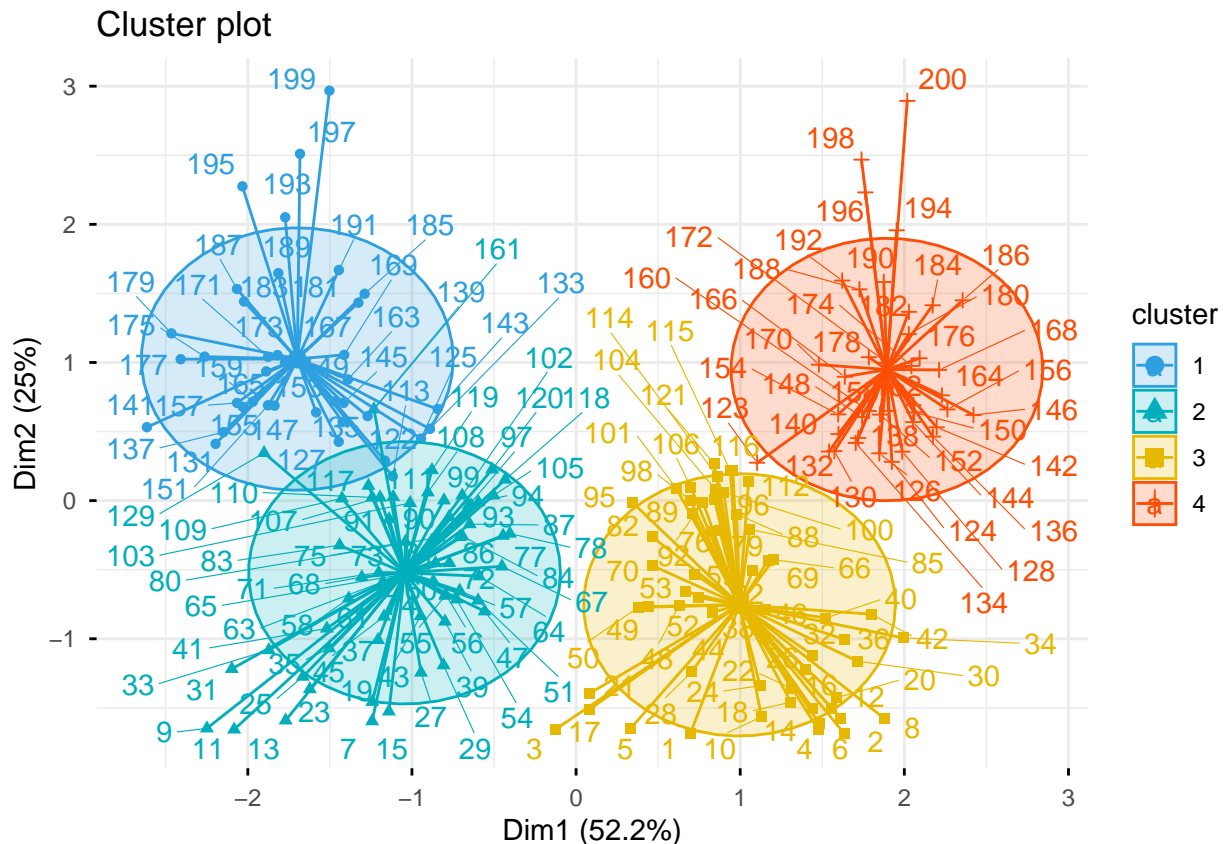
## K-means clustering with 4 clusters of sizes 38, 65, 57, 40
##
## Cluster means:
##   Annual.Income..k.. Spending.Score..1.100.      Age
## 1      0.9876366      -1.1857814  0.03711223
## 2      -0.4893373      -0.3961802  1.08344244
## 3      -0.7827991       0.3910484 -0.96008279
## 4       0.9724070       1.2130414 -0.42773261
##
## Clustering vector:
##   [1] 3 3 3 3 3 3 2 3 2 3 2 3 2 3 3 3 2 3 3 3 2 3 2 3 2 3 2 3 2 3 2 3 2 3 2
##  [38] 3 2 3 2 3 2 3 2 3 2 3 3 3 2 3 3 2 2 2 2 2 3 2 2 3 2 2 3 2 2 3 3 2 2 2 2
##  [75] 2 3 2 2 3 2 2 3 2 2 3 2 2 3 3 2 2 3 2 2 3 3 2 3 3 2 2 3 2 3 2 2 2 2 2 2
## [112] 3 1 3 3 3 2 2 2 2 3 1 4 4 1 4 1 4 2 4 1 4 1 4 1 4 1 4 1 4 1 4 1 4 1 4
## [149] 1 4 1 4 1 4 1 4 1 4 1 4 2 4 1 4 1 4 1 4 1 4 1 4 1 4 1 4 1 4 1 4 1 4 1
## [186] 4 1 4 1 4 1 4 1 4 1 4 1 4 1 4 1 4
##
## Within cluster sum of squares by cluster:
## [1] 44.01863 74.83280 61.43215 23.91544
## (between_SS / total_SS = 65.8 %)
```

```
##
## Available components:
##
## [1] "cluster"      "centers"      "totss"        "withinss"     "tot.withinss"
## [6] "betweenss"    "size"         "iter"         "ifault"
##### Step 2: Create the Clustering Visual #####

#reating a new dataset using the cbind() function to merge the subset data
#created in step 1, and the km.res$clsuter

dd <- cbind(MCEsv, cluster=km.res$cluster)

fviz_cluster(km.res, data=dd,
  palette = c("#2E9FDF", "#00AFBB", "#E7B800", "#FC4E07"),
  ellipse.type = "euclid", #Concentration ellipse
  star.plot = TRUE, #Add segments from centroids to items
  repel = TRUE, #Avoid label overplotting
  ggtheme = theme_minimal())
```



```
#boxplot
# Step 1: Select columns for clustering
data_for_clustering <- MCE[, c("Age", "Annual.Income..k..", "Spending.Score..1.100.")]

# Step 2: Run k-means clustering

set.seed(123) # for reproducibility
km.res <- kmeans(data_for_clustering, centers = 3, nstart = 25)
```

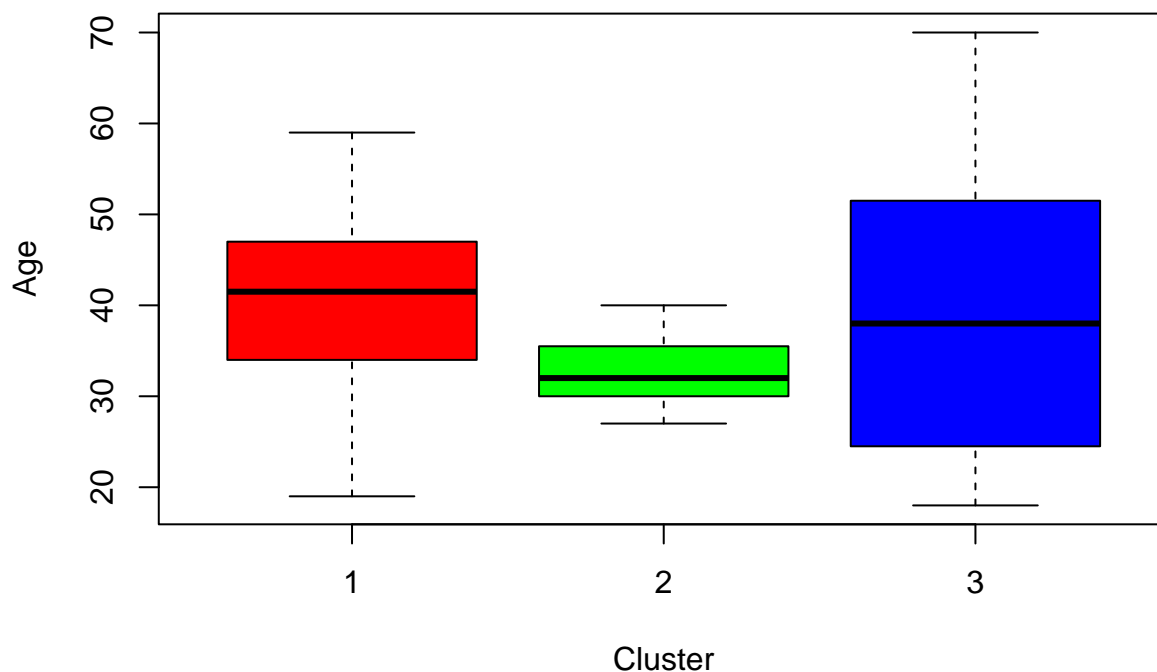
```
km.res$cluster
```

```
# Merge cluster with original data
```

```
# Boxplot for Age across clusters
```

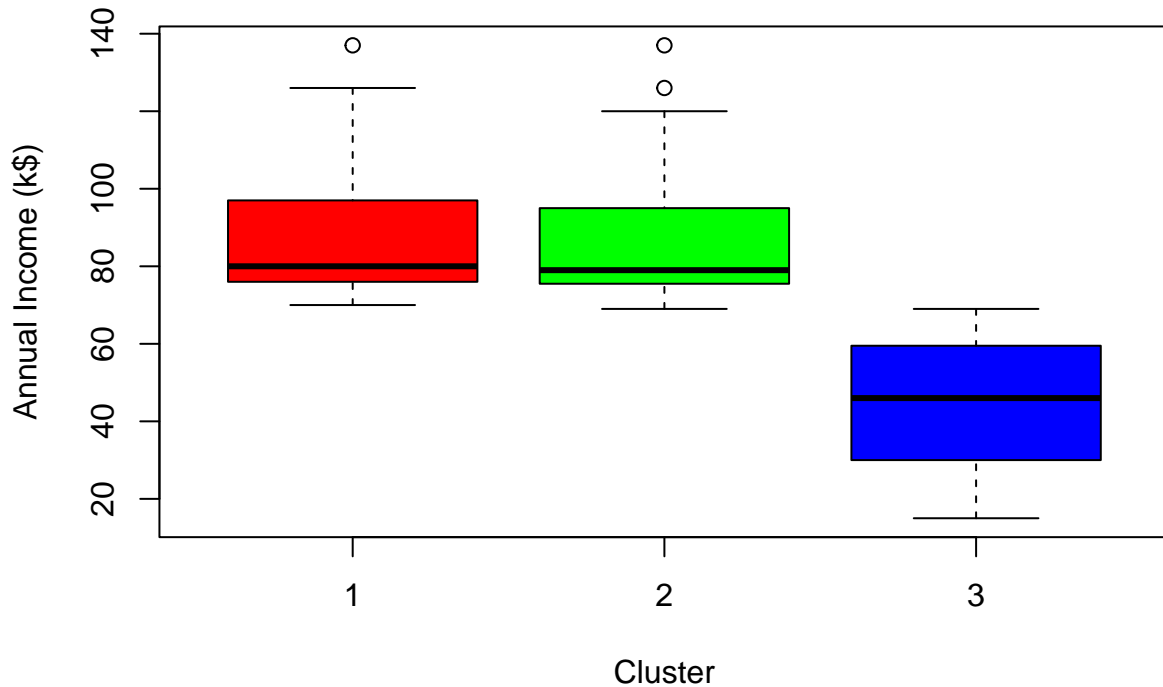
```
boxplot(Age ~ Cluster, data = MCE_clustered,
        main = "Distribution of Age Across Clusters",
        xlab = "Cluster", ylab = "Age",
        col = c("red", "green", "blue"))
```

## Distribution of Age Across Clusters



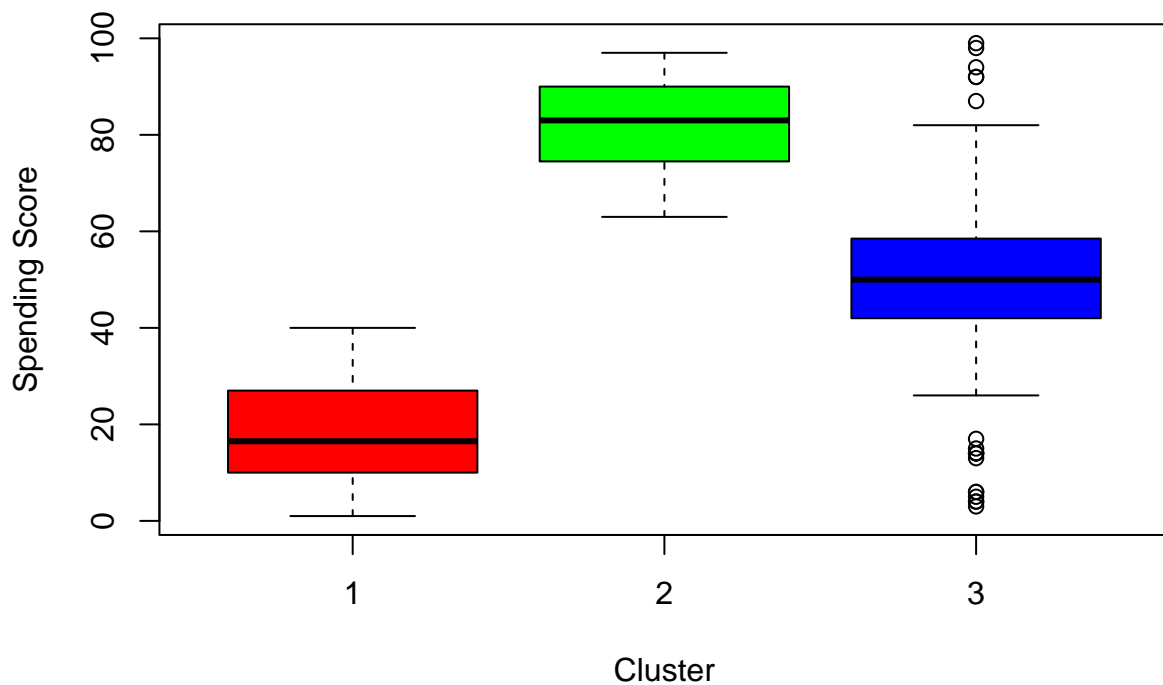
```
boxplot(Annual.Income..k.. ~ Cluster, data = MCE_clustered,
        main = "Distribution of Annual Income Across Clusters",
        xlab = "Cluster", ylab = "Annual Income (k$)",
        col = c("red", "green", "blue"))
```

## Distribution of Annual Income Across Clusters



```
# Boxplot for Spending Score across clusters
boxplot(Spending.Score..1.100. ~ Cluster, data = MCE_clustered,
        main = "Distribution of Spending Score Across Clusters",
        xlab = "Cluster", ylab = "Spending Score",
        col = c("red", "green", "blue"))
```

## Distribution of Spending Score Across Clusters



```
#####
# descriptive tables #
#####

install.packages("stargazer")
library(stargazer)

cluster1 <- subset(MCE_clustered, Cluster == 1)
cluster2 <- subset(MCE_clustered, Cluster == 2)
cluster3 <- subset(MCE_clustered, Cluster == 3)

variables <- c("Age", "Annual.Income..k..", "Spending.Score..1.100.")

stargazer(cluster1[, variables], type = "text", title = "Cluster 1 Descriptive Stats")
```

```
##
## Cluster 1 Descriptive Stats
## =====
## Statistic          N    Mean  St. Dev. Min Max
## -----
## Age                38 40.395  11.377  19  59
## Annual.Income..k.. 38 87.000  16.271  70 137
## Spending.Score..1.100. 38 18.632  10.916   1  40
## -----
```

```
stargazer(cluster2[, variables], type = "text", title = "Cluster 2 Descriptive Stats")
```

```
##
## Cluster 2 Descriptive Stats
## =====
## Statistic          N    Mean  St. Dev. Min Max
## -----
## Age                39 32.692   3.729  27  40
## Annual.Income..k.. 39 86.538  16.312  69 137
## Spending.Score..1.100. 39 82.128   9.364  63  97
## -----
```

```
stargazer(cluster3[, variables], type = "text", title = "Cluster 3 Descriptive Stats")
```

```
##
## Cluster 3 Descriptive Stats
## =====
## Statistic          N    Mean  St. Dev. Min Max
## -----
## Age                123 40.325  16.114  18  70
## Annual.Income..k.. 123 44.154  16.038  15  69
## Spending.Score..1.100. 123 49.829  19.694   3  99
## -----
```

```
#How to do multiple linear regression in R
```

```
MallData <- read.csv("Mall_Customers_extended.csv", header=TRUE)
```

```
names(MallData)
```

```
## [1] "Unnamed..0"          "CustomerID"          "Genre"
## [4] "Age"                  "Annual.Income..k.."  "Spending.Score..1.100."
## [7] "IncomePerAge"         "SpendingEfficiency"  "AgeCategory"
## [10] "HighSpender"         "IncomeTier"          "OnlineShopFreq"
## [13] "LoyaltyScore"        "Satisfaction"        "CreditUtilization"
```

```
y <- MallData$y.LoyaltyScore
```

```
x1 <- MallData$Age
```

```
x2 <- MallData$Annual.Income..k..
```

```
x3 <-MallData$Spending.Score..1.100.
```

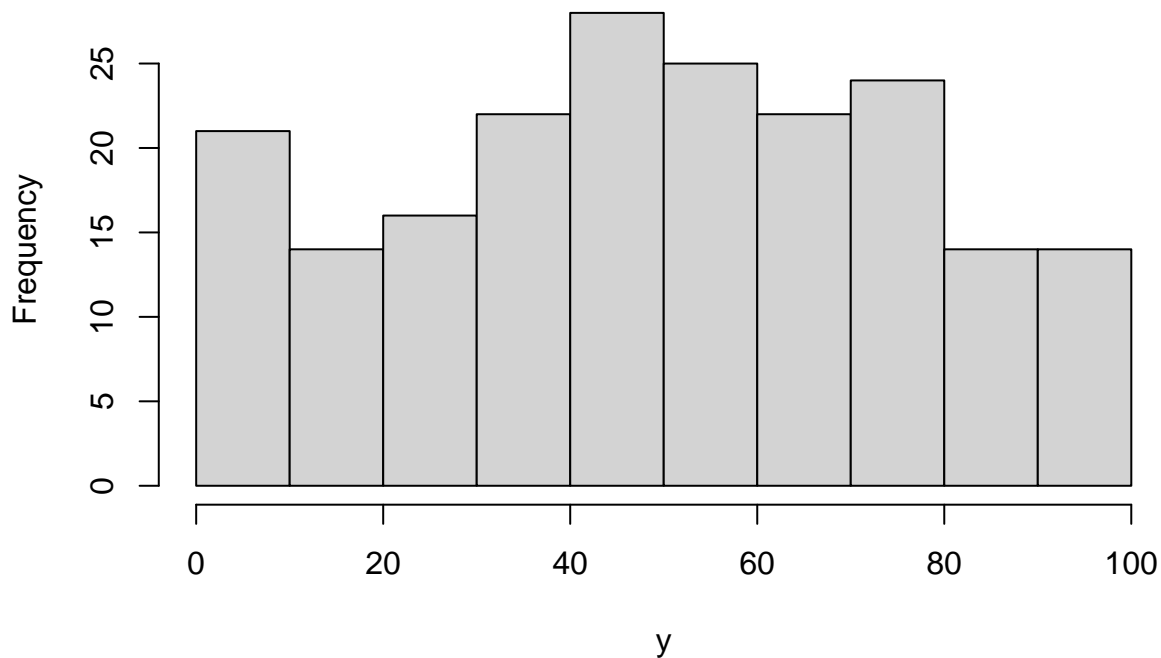
```
str(MallData$LoyaltyScore)
```

```
## num [1:200] 46 75 28 83.9 39.9 ...
```

```
y <- as.numeric(as.character(MallData$LoyaltyScore))
```

```
hist(y)
```

**Histogram of y**



```
model2 <- lm(y ~ x1 + x2 + x3, data=MallData)
```

```
model2
```

```
##
```

```
## Call:
```

```
## lm(formula = y ~ x1 + x2 + x3, data = MallData)
```

```
##
```

```
## Coefficients:
```

```
## (Intercept)          x1          x2          x3
```



```
##      4.40942      -0.08975      -0.01474      0.98109
```

```
attach(MallData)
```

```
summary(model2)
```

```
##
```

```
## Call:
```

```
## lm(formula = y ~ x1 + x2 + x3, data = MallData)
```

```
##
```

```
## Residuals:
```

```
##      Min       1Q   Median       3Q      Max
## -22.0324  -6.4708  -0.5225   5.9657  31.4818
```

```
##
```

```
## Coefficients:
```

```
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  4.40942     3.02939   1.456  0.1471
## x1          -0.08975     0.04761  -1.885  0.0609 .
## x2          -0.01474     0.02393  -0.616  0.5385
## x3           0.98109     0.02575  38.099 <2e-16 ***
```

```
## ---
```

```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
##
```

```
## Residual standard error: 8.864 on 196 degrees of freedom
```

```
## Multiple R-squared:  0.8956, Adjusted R-squared:  0.894
```

```
## F-statistic: 560.7 on 3 and 196 DF,  p-value: < 2.2e-16
```

*#Intercepts: There are 3 x-values that are independent: these are x1 is the #intercept for Age. x2 is the intercept for Annual Income K and x3 is the #intercept for Spending Score. The estimated Beta for x1 is -0.08975, x2 is #-0.01474 and x3 is 0.98109. There are 3 different variables with different #outcomes. Age, Annual Income K and Spending Score are all continuous variables. #For every unit in x1 (Age), there is a decrease of -0.0875 for the outcome #(LoyaltyScore). Applies for x2, there is a decrease of -0.01474 for the outcome #and for x3, there is an increase of 0.98108 for the outcome.*

*#Based on the summary, x1 and x2 are not significant. They do not contain the #asterick. In this multiple Linear Regression ,x3 significantly predicts the #outcome. For every unit in x3, the expected outcome increases by 0.98109 #which is adjusting for x1 and x2.*

*#Multiple R-Squared: This is used for Multiple Linear Regression to avoid bias #when trying to find associations of x variables with outcome. When adding more #variables R, usually increases, therefore we avoid using Multiple R-squared in #this situation. Thus, they both provide the proportion of variability that #the x's will explain on the variance of the outcome. It is ideal for it to be #high, it should be close to 1. This would mean that all of systematic #variances are being accounted for, in this model, it is high of 0.894.*

*#F-Statistics: The F-statistics is more than 1, it is 560.7, this means that #p-value is significant as well. When F-statistics is more than 1, it means that #the systematic variances are being accounted by these x variables and they #outweigh the unsystematic variances.*

#p-Value: This would be considered a good model because it is less than 0.05,  
#shows that it is statistically significant.

#Income does not predict more Loyalty stronger than other clusters. The cluster  
#that show to have a stronger Loyalty is x3 (Spending Score) with a positive  
#beta of 0.98109. The other betas for x1 and x2 variable show decrease with  
#association of outcome.

#Assessing fit of model  
`plot(model2)`

