EXPERIMENT 9

Taniya Ahmed

21BDS0059

1. Loading the libraries used, importing the dataset and renaming the column names.

```
CODE:
library(ggplot2)
df = read.csv("D:\\Downloads\\Mall_Customers.csv")
colnames(df)
colnames(df) = c("CustomerID", "Gender", "Age", "AnnualIncome", "SpendingScore")
head(df)
print("Taniya Ahmed 21BDS0059")
OUTPUT:
> library(ggplot2)
> df = read.csv("D:\\Downloads\\Mall_Customers.csv")
  colnames (df)
[1] "CustomerID" "Gender"
[5] "Spending.Score..1.100."
                                                          "Age"
                                                                                      "Annual.Income..k.."
> spending.score..1.100.
> colnames(df) = c("CustomerID", "Gender", "Age", "AnnualIncome", "SpendingScore")
> head(df)
  CustomerID Gender Age AnnualIncome SpendingScore
         1 Male 19 15
2 Male 21 15
3 Female 20 16
4 Female 23 16
5 Female 31 17
                                                    81
                                                    6
77
                                                    40
76
6 6 Female 22 > print("Taniya Ahmed 21BDS0059")
[1] "Taniya Ahmed 21BDS0059"
```

2. Selecting features based on which the clustering will be done.

```
CODE:

cluster_features = df[, c("AnnualIncome", "SpendingScore")]

head(cluster_features)

print("Taniya Ahmed 21BDS0059")
```

OUTPUT:

CODE:

```
> cluster_features = df[, c("AnnualIncome", "SpendingScore")]
> head(cluster_features)
  AnnualIncome SpendingScore
            15
2
            15
                           81
3
            16
                           6
4
            16
            17
                           40
6
            17
                          76
> print("Taniya Ahmed 21BDS0059")
[1] "Taniya Ahmed 21BDS0059"
```

3. Calculating shortest distance between points using Euclidean distance.

4. Assigning clusters according to shortest distance between the points and the centroids.

```
assign_cluster = function(data, centroids){
  clusters = vector("numeric", nrow(data))

for(i in 1 : nrow(data)){
  distances = apply(centroids, 1, euclid_distance, b= data[i, ])
  clusters[i] = which.min(distances)
}
```

5. Calculating the new centroids as every iteration requires a centroid calculation.

```
CODE:
```

```
new_centroid_calculation = function(data, clusters, k){
  centroids = matrix(NA, nrow = k, ncol = ncol(data))

for(i in 1 : k){
  centroids[i, ] = colMeans(data[clusters == i, , drop = FALSE])
}
  return(centroids)
}
print("Taniya Ahmed 21BDS0059")
```

OUTPUT:

```
> new_centroid_calculation = function(data, clusters, k){
+ centroids = matrix(NA, nrow = k, ncol = ncol(data))
+
+ for(i in 1 : k){
+ centroids[i, ] = colMeans(data[clusters == i, , drop = FALSE])
+ }
+ return(centroids)
+ }
> new_centroid_calculation = function(data, clusters, k){
+ centroids = matrix(NA, nrow = k, ncol = ncol(data))
+ for(i in 1 : k){
+ centroids[i, ] = colMeans(data[clusters == i, , drop = FALSE])
+ }
+ return(centroids)
+ }
> print("Taniya Ahmed 21BDS0059")
[1] "Taniya Ahmed 21BDS0059"
> |
```

6. Designing the custom kmeans algorithm by assigning clusters, , calculating centroids and iterating multiple times until convergence or 150 iterations whichever happens first.

```
CODE:
```

```
kmeans_custom = function(data, k, max_iter = 150){
    set.seed(125)
    ini_centroid = data[sample(1 : nrow(data), k),]
    centroids = ini_centroid
    clusters = assign_cluster(data, centroids)

for(i in 1 : max_iter){
    centroids = new_centroid_calculation(data, clusters, k)
    new_clusters = assign_cluster(data, centroids)

if(all(clusters == new_clusters)){
    break
  }

clusters = new_clusters
}
```

```
return(list(clusters = clusters, centroids = centroids))
}
print("Taniya Ahmed 21BDS0059")
OUTPUT:
> kmeans_custom = function(data, k, max_iter = 150){
   set.seed(125)
    ini_centroid = data[sample(1 : nrow(data), k),]
    centroids = ini_centroid
    clusters = assign_cluster(data, centroids)
    for(i in 1 : max_iter){
     centroids = new_centroid_calculation(data, clusters, k)
      new_clusters = assign_cluster(data, centroids)
      if(all(clusters == new_clusters)){
        break
      clusters = new_clusters
    return(list(clusters = clusters, centroids = centroids))
> print("Taniya Ahmed 21BDS0059")
[1] "Taniya Ahmed 21BDS0059"
```

7. Specifying the number of clusters, checking the clusters and centroids.

```
CODE:
k = 4
kmeans_final = kmeans_custom(cluster_features, k)
print(kmeans_final$clusters)
print(kmeans_final$centroids)
print("Taniya Ahmed 21BDS0059")
```

OUTPUT:

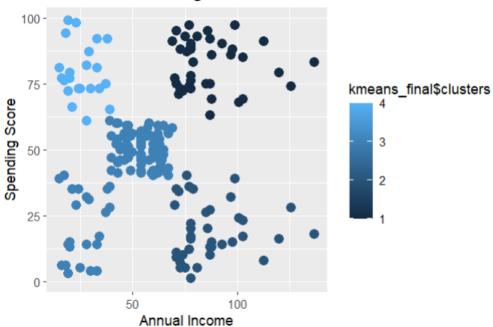
8. Making a scatter plot to evaluate the clusters.

CODE:

```
ggplot(df, aes(x = AnnualIncome, y = SpendingScore, color = kmeans_final$clusters)) +
geom_point(size = 3) +
labs(title = "K - Means clustering with 3 clusters", x = "Annual Income", y = "Spending Score")
print("Taniya Ahmed 21BDS0059")
```

OUTPUT:

K - Means clustering with 3 clusters



```
> ggplot(df, aes(x = AnnualIncome, y = SpendingScore, color = kmeans_final$clusters)) +
+ geom_point(size = 3) +
+ labs(title = "K - Means clustering with 3 clusters", x = "Annual Income", y = "Spending Score")
> print("Taniya Ahmed 21BDS0059")
[1] "Taniya Ahmed 21BDS0059"
> |
```