

MINI PROJECT ON

Project Entitled:
‘CUSTOMER SEGMENTATION ‘

Submitted by:
Kaustubh Patil (A-36)
Yash Surana (A-45)

CUSTOMER SEGMENTATION USING R

CUSTOMER SEGEMENTATION -

Customer Segmentation is one the most important applications of unsupervised learning. Using clustering techniques, companies can identify the several segments of customers allowing them to target the potential user base. In this machine learning project, we will make use of K-means Clustering which is the essential algorithm for clustering unlabeled dataset.

Customer Segmentation is the process of division of customer base into several groups of individuals that share a similarity in different ways that are relevant to marketing such as gender, age, interests, and miscellaneous spending habits.

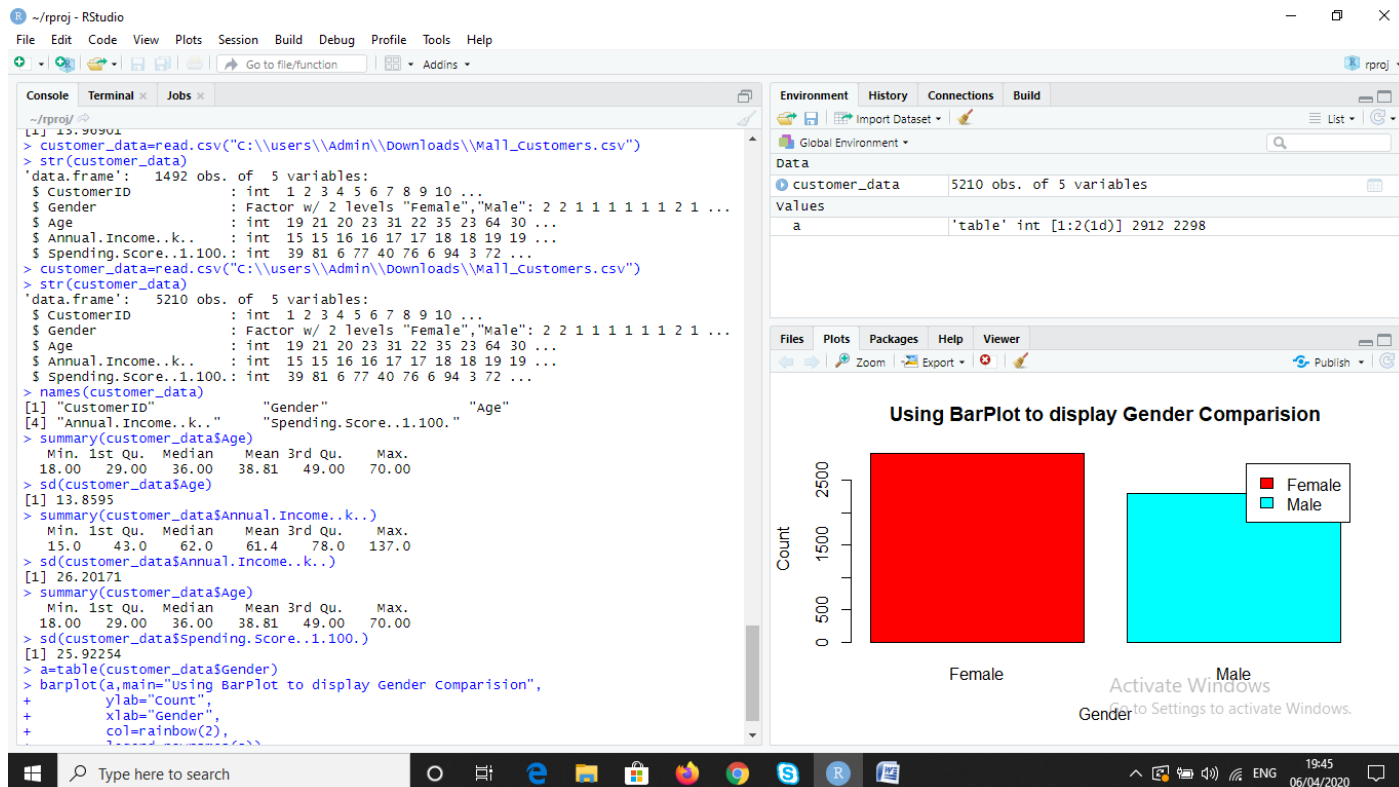
Companies that deploy customer segmentation are under the notion that every customer has different requirements and require a specific marketing effort to address them appropriately. Companies aim to gain a deeper approach of the customer they are targeting. Therefore, their aim has to be specific and should be tailored to address the requirements of each and every individual customer. Furthermore, through the data collected, companies can gain a deeper understanding of customer preferences as well as the requirements for discovering valuable segments that would reap them maximum profit. This way, they can strategize their marketing techniques more efficiently and minimize the possibility of risk to their investment.

DATASET SCREENSHOT-

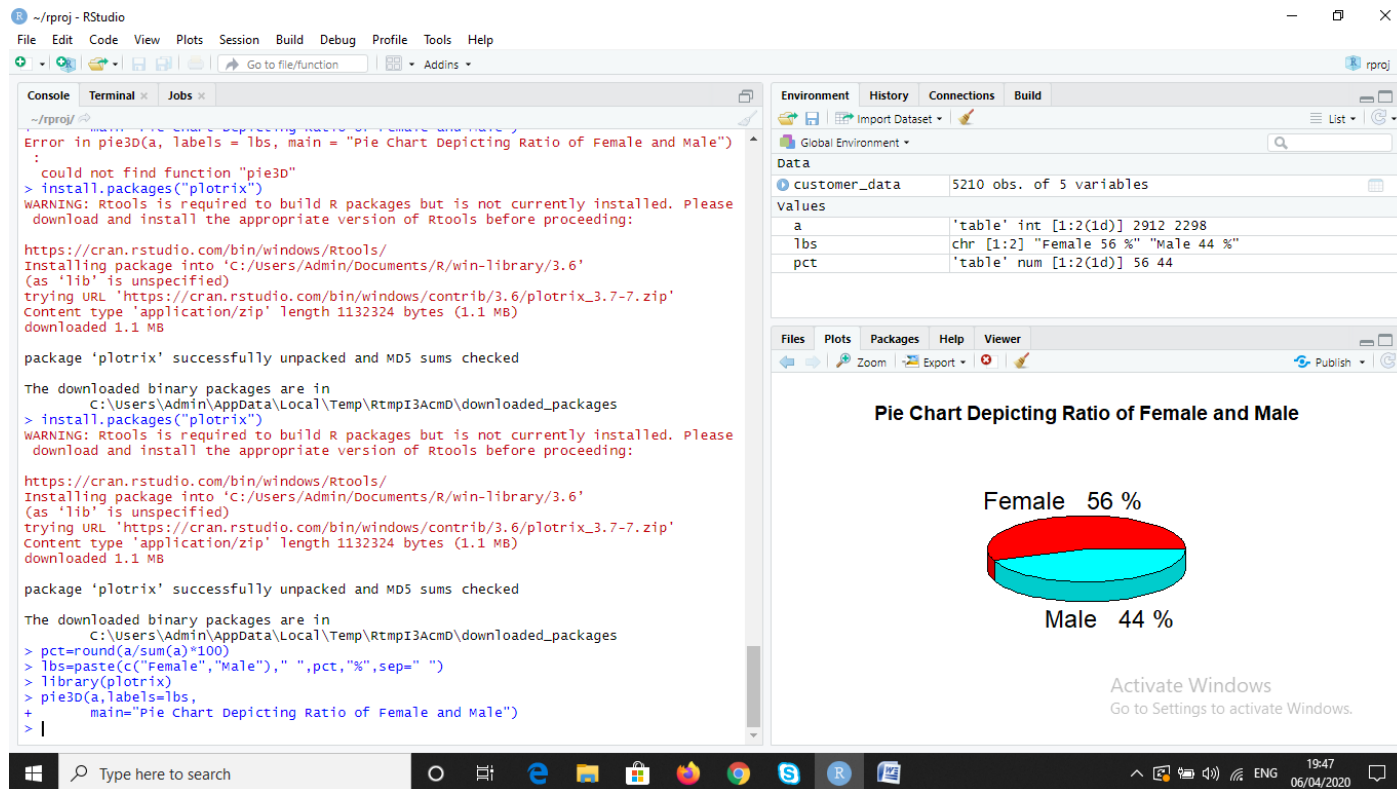
The screenshot shows a Microsoft Excel spreadsheet with the following data:

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q
	Customer	Gender	Age	Annual Income (k\$)	Spending Score (1-100)												
1	1	Male	19	15	39												
2	2	Male	21	15	81												
3	3	Female	20	16	6												
4	4	Female	23	16	77												
5	5	Female	31	17	40												
6	6	Female	22	17	76												
7	7	Female	35	18	6												
8	8	Female	23	18	94												
9	9	Male	64	19	3												
10	10	Female	30	19	72												
11	11	Male	67	19	14												
12	12	Female	35	19	99												
13	13	Female	58	20	15												
14	14	Female	24	20	77												
15	15	Male	37	20	13												
16	16	Male	22	20	79												
17	17	Female	35	21	35												
18	18	Male	20	21	66												
19	19	Male	52	23	29												
20	20	Female	35	23	98												
21	21	Male	35	24	35												
22	22	Male	25	24	73												
23	23	Male	25	24	73												

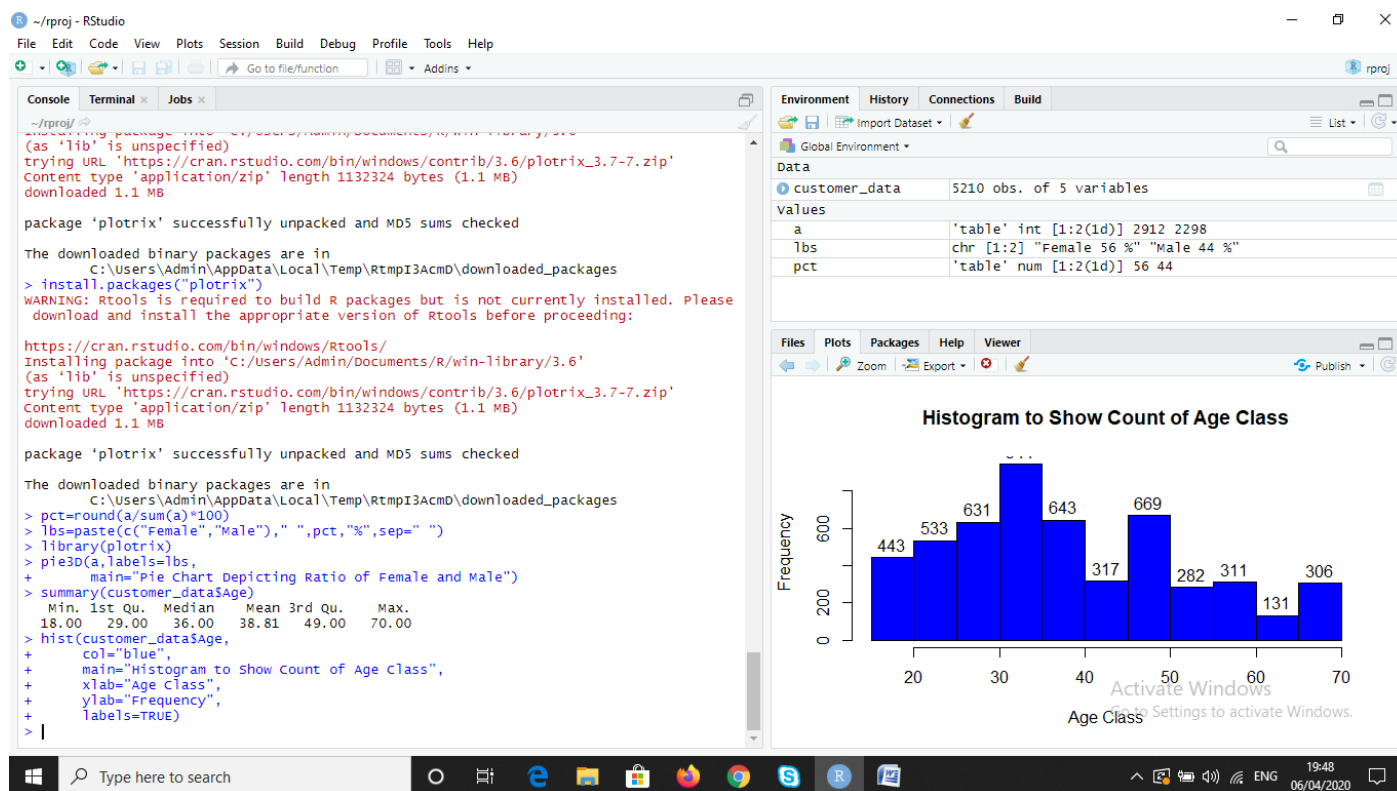
OUTPUT-



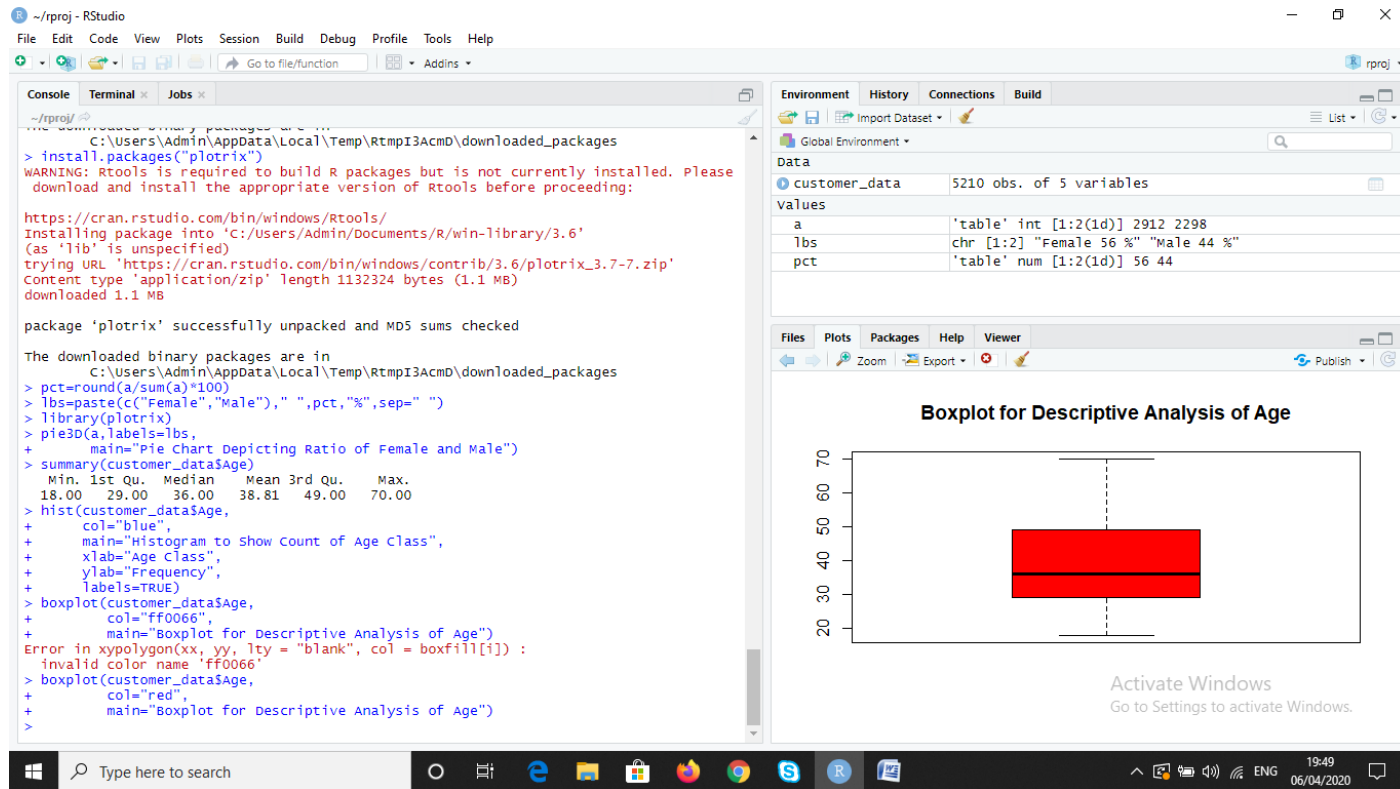
BAR PLOT



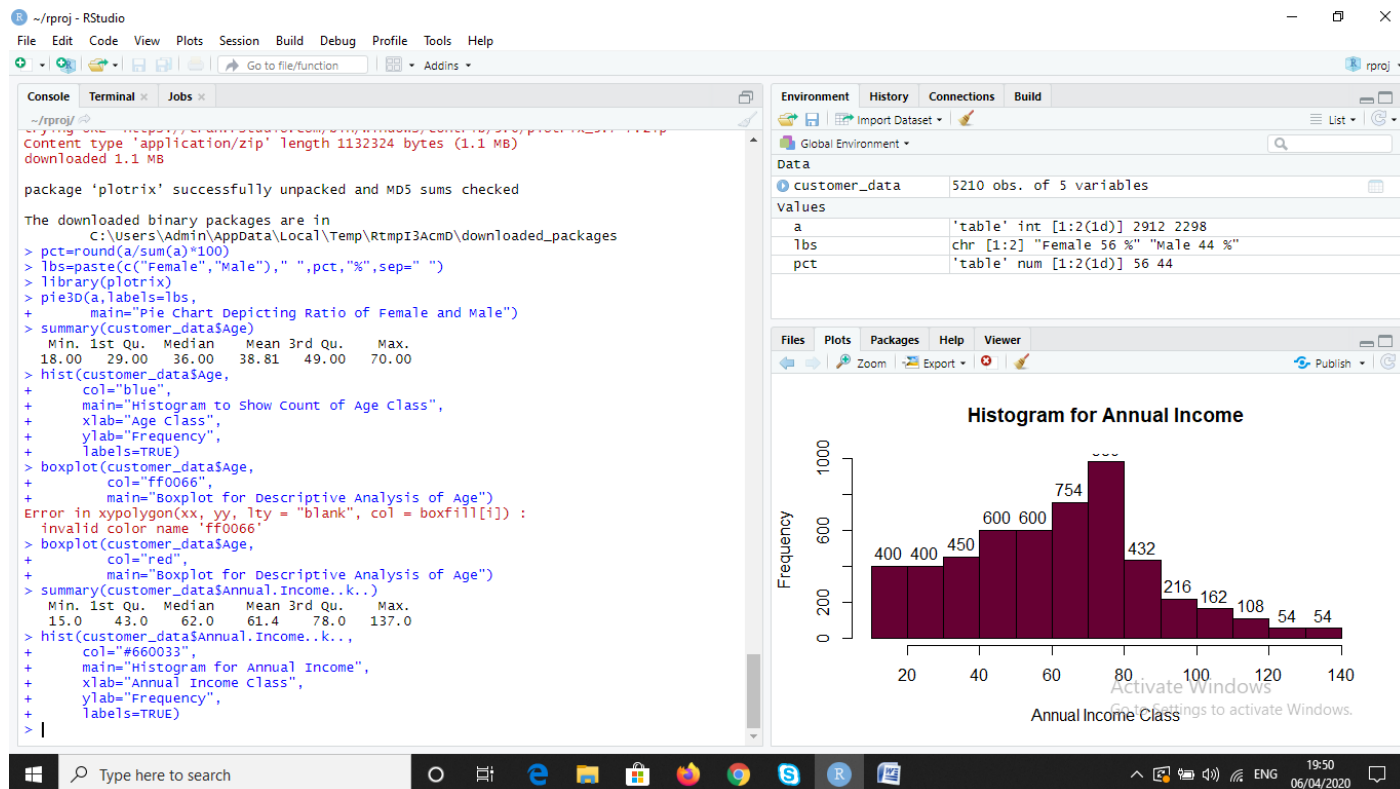
PIE CHART



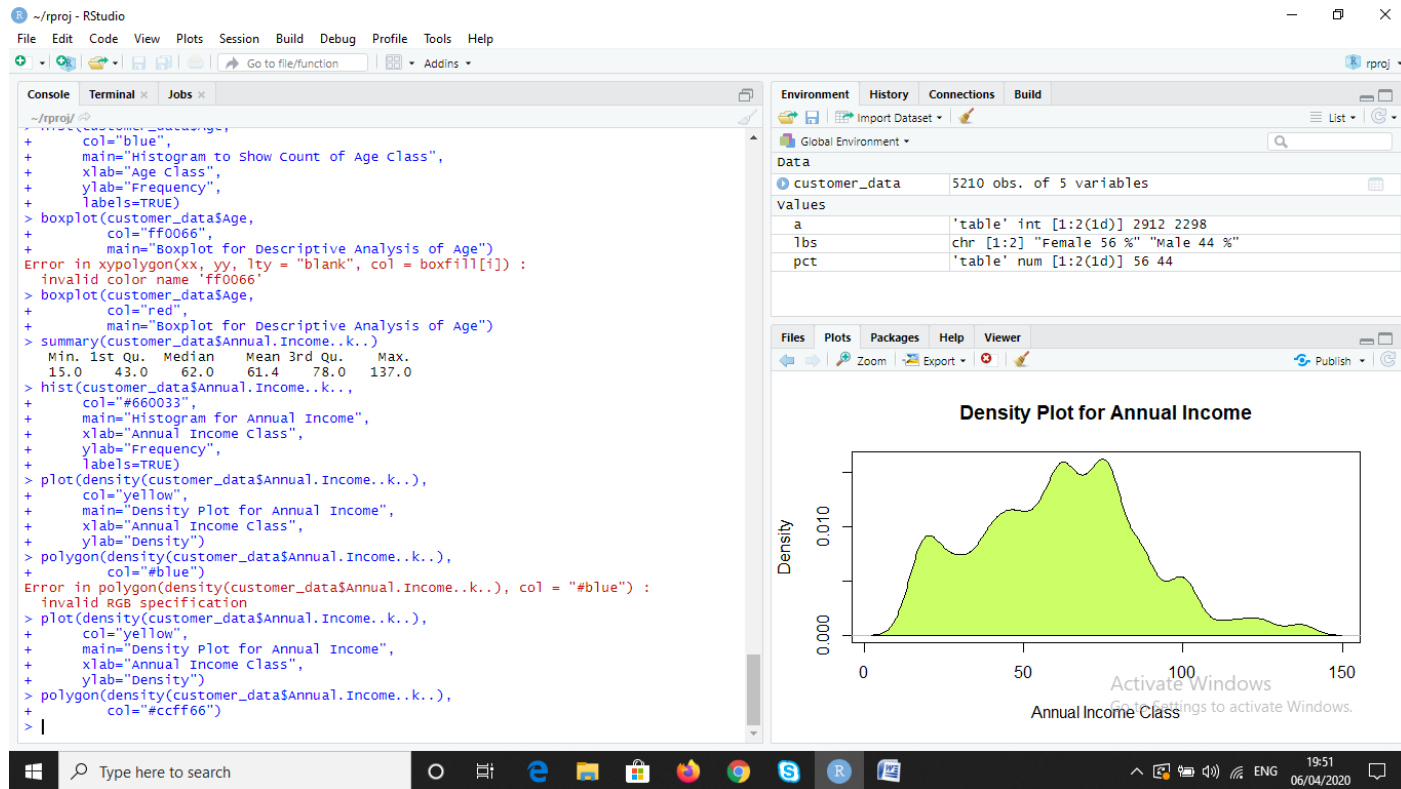
HISTROGRAM



BOX PLOT



HISTOGRAM



DENSITY PLOT

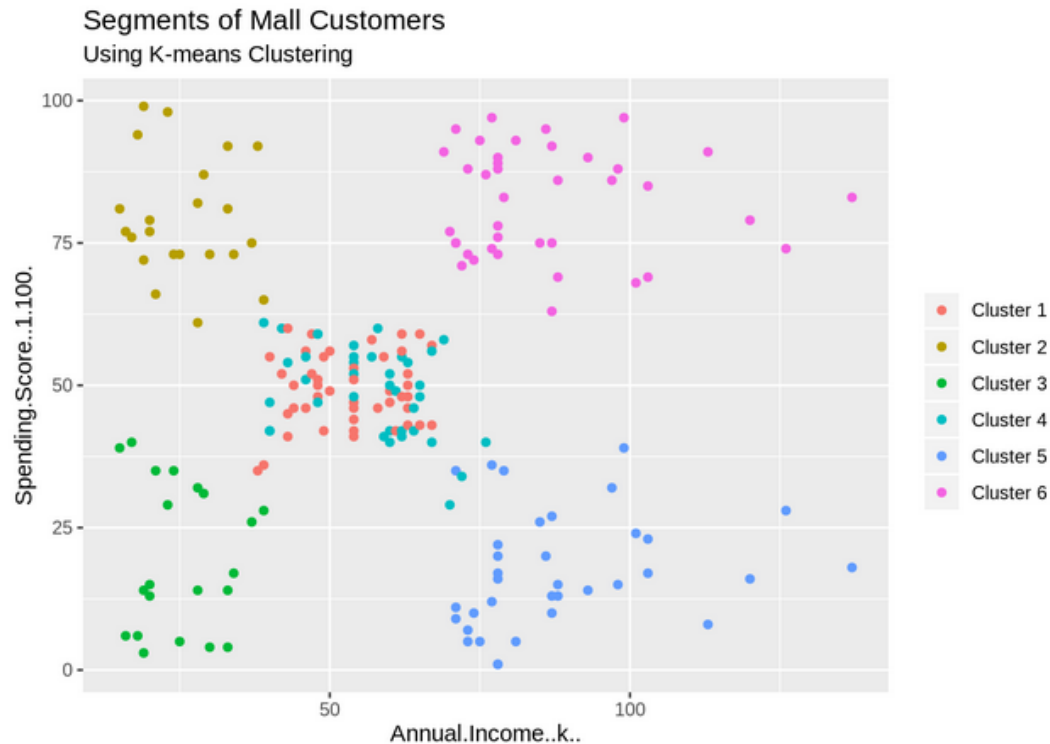
K-means Algorithm-

While using the k-means clustering algorithm, the first step is to indicate the number of clusters (k) that we wish to produce in the final output. The algorithm starts by selecting k objects from dataset randomly that will serve as the initial centers for our clusters. These selected objects are the cluster means, also known as centroids. Then, the remaining objects have an assignment of the closest centroid. This centroid is defined by the Euclidean Distance present between the object and the cluster mean. We refer to this step as “cluster assignment”. When the assignment is complete, the algorithm proceeds to calculate new mean value of each cluster present in the data. After the recalculation of the centers, the observations are checked if they are closer to a different cluster. Using the updated cluster mean, the objects undergo reassignment. This goes on repeatedly through several iterations until the cluster assignments stop altering. The clusters that are present in the current iteration are the same as the ones obtained in the previous iteration.

CODE-

```
set.seed(1)

ggplot(customer_data, aes(x = Annual.Income..k., y = Spending.Score..1.100.)) +
  geom_point(stat = "identity", aes(color = as.factor(k6$cluster))) +
  scale_color_discrete(name=" ",
  breaks=c("1", "2", "3", "4", "5", "6"),
  labels=c("Cluster 1", "Cluster 2", "Cluster 3", "Cluster 4", "Cluster 5", "Cluster 6")) +
  ggtitle("Segments of Mall Customers", subtitle = "Using K-means Clustering")
```



Cluster 1 – This cluster represents the customer_data having a high annual income as well as a high annual spend.

Cluster 2 – This cluster denotes the customer_data with low annual income as well as low yearly spend of income.

Cluster 3 – These clusters represent the customer_data with the medium income salary as well as the medium annual spend of salary.

Cluster 4 – This cluster denotes a high annual income and low yearly spend.

Cluster 5 – This cluster represents a low annual income but its high yearly expenditure

Cluster 6 – These clusters represent the customer_data with the medium income salary as well as the medium annual spend of salary.

Summary-

In this data science project, we went through the customer segmentation model. We developed this using a class of machine learning known as unsupervised learning. Specifically, we made use of a clustering algorithm called K-means clustering. We analyzed and visualized the data and then proceeded to implement our algorithm.