



DATA SCIENCE INTERNSHIP

Project Title:

CUSTOMER SEGMENTATION USING K-MEANS ALGORITHM

Name:

Woody vijaykumar samyuktha

Abstract

In today's fast-moving world of marketing from product-orientation to customer-orientation, the management of customer treatment can be seen as a key to achieve revenue growth and profitability. Knowledge of customer behavior can help marketing managers re-evaluate their strategies with the customers and plan to improve and expand their application of the most effective strategies. B2B or business customers are more complex, their buying process is more complicated and their sales value is greater. The business marketers usually prefer to cooperate with fewer but larger buyers than the final consumer marketer. As a business transaction requires more decision makings and more professional buying effort than the consumer market does, the efficient relationship with business customers is of paramount importance. Most customer segmentation approaches based on customer value fail to account for

the factor of time and the trend of value changes in their analysis. The business done today runs on the basis of innovative ideas as there are large number of potential customers who are confounded to what to buy and what not to buy. The companies doing the business are also not able to diagnose the target potential customers. This is where the machine learning comes into picture, the various algorithms are applied to identify the hidden patterns in the data for better decision making. The concept of which customer segment to target is done using the customer segmentation process using the clustering technique. In this project, the clustering algorithm we are going to use is K-means algorithm which is the partitioning algorithm, to segment the customers according to the similar characteristics.

To determine the optimal clusters, elbow method is used.

Abstract:

- (i) Title of the project.

Customer Segmentation using K-Means Algorithm

- (ii)

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.

The technique of customer segmentation is dependent on several key differentiators that divide customers into groups to be targeted. Data related to demographics, geography, economic status as well as

behavioural patterns play a crucial role in determining the company direction towards addressing the various segment

2. Concept of the project and modules in the project.

- (i) List the Modules in your project
 - ✓ Data Exploration
 - ✓ Data Pre-processing
 - ✓ Data Visualization
 - ✓ Data Analysis
- (ii) List the concepts to be applied in your project.
 - ✓ Un-Supervised Algorithm(K-Means Clustering)
 - ✓ Clustering
 - ✓ Elbow Method
 - ✓ Seaborn,matplotlib,etc.. – Visualizations
 - ✓ All necessary concepts will be used for all the Modules involved in the project
- (iii) Platform & Language to be used for Implementation.

The project will be implemented in Python

- ✓ Python -google collab

Objective:

1. To cluster customers based on common purchasing behaviors for future operations/marketing projects
2. To incorporate best mathematical, visual, programming, and business practices into a thoughtful analysis that is understood across a variety of contexts and disciplines
3. The goal of segmenting customers is to decide how to relate to customers in each segment in order to maximize the value of each customer to the business.

Methodology:

The data set used to implement clustering and Kmeans algorithm was collected from a store of shopping mall. The data set contains 5 attributes and has 200 tuples, representing the data of 200 customers. The attributes in the data set has CustomerId, gender, age, annual income(k\$), spending score on the scale of (1-100).

Literature Review:***Customer Segmentation***

Customer Segmentation Over the years, as there is very strong competition in the business world, the organizations have to enhance their profits and business by satisfying the demands of their customers and attract new customers according to their needs. The identification of customers and satisfying the demands of each customer is a very complex and tedious task. This is because customers may be different according to their demands, tastes, preferences and so on. Instead of “one-size-fits-all” approach, customer segmentation clusters the customers into groups sharing the same properties or behavioral characteristics. Customer segmentation is a strategy of dividing the market into homogenous groups. The data used in customer segmentation technique that divides the customers into groups depends on various factors like, data geographical conditions, economic conditions, demographical conditions as well as behavioral patterns. The customer segmentation technique allows the business to make better use of their marketing budgets, gain a competitive edge over their rival companies, demonstrating the better knowledge of the needs of the customer. It also helps an organization in, increasing their marketing efficiency, determining new market opportunities, making better brand strategy, identifying customers retention.

Clustering and K-Means Algorithm

Clustering algorithms generate clusters such that within the clusters are similar based on some characteristics. Similarity is defined in terms of how close the objects are in space. K-means algorithm is one of the most popular centroid-based algorithms. Suppose data set, D , contains n objects in space. Partitioning methods distribute the objects in D into k clusters, C_1, \dots, C_k , that is, $C_i \subset D$ and $C_i \cap C_j = \emptyset$ for $(1 \leq i, j \leq k)$. A centroid-based partitioning technique uses the centroid of a cluster, C_i , to represent that cluster. Conceptually, the centroid of a cluster is its center point. The difference between an object $p \in C_i$ and c_i , the representative of the cluster, is measured by $\text{dist}(p, c_i)$, where $\text{dist}(x, y)$ is the Euclidean distance between two points x and y .

Algorithm

The k-means algorithm for partitioning, where each cluster's center is represented by the mean value of the objects in the cluster. Input: k : the number of clusters, D : a data set containing n objects. Output: A set of k clusters. Method:

- (1) arbitrarily choose k objects from D as the initial cluster centers;
- (2) repeat
- (3) (re)assign each object to the cluster to which the object is the most similar, based on the mean value of the objects in the cluster;
- (4) update the cluster means, that is, calculate the mean value of the objects for each cluster; (5) until no change.

Implementation of a Customer Segmentation

In the first step of this Fundamental Data Analytics project, we will perform data exploration.

We will import the essential packages required for this role and then read our data.

Finally, we will go through the input data to gain necessary insights about it.

Code

Import dependencies

```
import numpy as np
```

```
import pandas as pd
```

```
import matplotlib.pyplot as plt
```

```
import seaborn as sns
```

```
from sklearn.cluster import KMeans
```

output:

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.cluster import KMeans
```

Data collection and Analysis

- Loading the data from csv file to a pandas dataframe
`customer_Data = pd.read_csv('Mall_Customers.csv')`
- First five rows in the dataframe
`customer_Data.head()`

Output:

	CustomerID	Gender	Age	Annual Income (k\$)	Spending Score (1-100)
0	1	Male	19	15	39
1	2	Male	21	15	81
2	3	Female	20	16	6
3	4	Female	23	16	77
4	5	Female	31	17	40

Finding the number of rows and columns

Code:

```
customer_Data.shape
```

Output:

```
#finding the number of rows and columns
customer_Data.shape
```

(200, 5)

Getting some information about the dataset

Code:

```
customer_Data.info()
```

Output:

```
#getting some information about the dataset
customer_Data.info()
```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 200 entries, 0 to 199
Data columns (total 5 columns):
Column Non-Null Count Dtype
--- ---
0 CustomerID 200 non-null int64
1 Gender 200 non-null object
2 Age 200 non-null int64
3 Annual Income (k\$) 200 non-null int64
4 Spending Score (1-100) 200 non-null int64
dtypes: int64(4), object(1)
memory usage: 7.9+ KB

Checking for missing values

Code:

```
customer_Data.isnull().sum()
```

Output:


```
#checking for missing values
customer_Data.isnull().sum()
```

	0
CustomerID	0
Gender	0
Age	0
Annual Income (k\$)	0
Spending Score (1-100)	0

dtype: int64

Choosing the annual income column and spending score columns

Code:

```
x=customer_Data.iloc[:,[3,4]].values
print(x)
```

Output:

```
[ ] x=customer_Data.iloc[:,[3,4]].values
print(x)
```

[15 39]
[15 81]
[16 6]
[16 77]
[17 40]
[17 76]
[18 6]
[18 94]
[19 3]
[19 72]
[19 14]
[19 99]
[20 15]
[20 77]
[20 13]
[20 79]
[21 35]
[21 66]
[23 29]
[23 98]
[24 35]
[24 73]
[25 5]
[25 73]
[28 14]
[28 82]
[28 32]
[28 61]
[29 31]
[29 87]
[30 4]
[30 73]
[33 4]
[33 92]
[33 14]
[33 81]
[34 17]
[34 73]
[37 26]
[37 75]
[38 35]
[38 92]
[39 36]
[39 61]
[39 28]
[39 65]
[40 55]
[40 47]
[40 42]
[40 42]
[42 52]
[42 60]
[43 54]
[43 60]
[43 45]
[43 41]
[44 50]
[44 46]

```

print(x)
[[ 44 46]
 [ 46 51]
 [ 46 46]
 [ 46 56]
 [ 46 55]
 [ 47 52]
 [ 47 59]
 [ 48 51]
 [ 48 59]
 [ 48 50]
 [ 48 48]
 [ 48 59]
 [ 48 47]
 [ 49 55]
 [ 49 42]
 [ 50 49]
 [ 50 56]
 [ 54 47]
 [ 54 54]
 [ 54 53]
 [ 54 48]
 [ 54 52]
 [ 54 42]
 [ 54 51]
 [ 54 55]
 [ 54 41]
 [ 54 44]
 [ 54 57]
 [ 54 46]
 [ 57 58]
 [ 57 55]
 [ 58 60]
 [ 58 46]
 [ 59 55]
 [ 59 41]
 [ 60 49]
 [ 60 40]
 [ 60 42]
 [ 60 52]
 [ 60 47]
 [ 60 50]
 [ 61 42]
 [ 61 49]
 [ 62 41]
 [ 62 48]
 [ 62 59]
 [ 62 55]
 [ 62 56]
 [ 62 42]
 [ 63 50]
 [ 63 46]
 [ 63 43]
 [ 63 48]
 [ 63 52]
 [ 63 54]
 [ 64 42]
 [ 64 46]
 [ 65 48]

```

```

print(x)
[[ 65 43]
 [ 65 59]
 [ 67 43]
 [ 67 57]
 [ 67 56]
 [ 67 40]
 [ 69 58]
 [ 69 91]
 [ 70 29]
 [ 70 77]
 [ 71 35]
 [ 71 95]
 [ 71 11]
 [ 71 75]
 [ 71 9]
 [ 71 75]
 [ 72 34]
 [ 72 71]
 [ 73 5]
 [ 73 88]
 [ 73 7]
 [ 73 73]
 [ 74 10]
 [ 74 72]
 [ 75 5]
 [ 75 93]
 [ 76 40]
 [ 76 87]
 [ 77 12]
 [ 77 97]
 [ 77 36]
 [ 77 74]
 [ 78 22]
 [ 78 90]
 [ 78 17]
 [ 78 88]
 [ 78 20]
 [ 78 76]
 [ 78 16]
 [ 78 89]
 [ 78 1]
 [ 78 78]
 [ 78 1]
 [ 78 73]
 [ 79 35]
 [ 79 83]
 [ 81 5]
 [ 81 93]
 [ 85 26]
 [ 85 75]
 [ 86 20]
 [ 86 95]
 [ 87 27]
 [ 87 63]
 [ 87 13]
 [ 87 75]
 [ 87 10]
 [ 87 92]
]

```

```

[[ 87 92]
 [ 88 13]
 [ 88 86]
 [ 88 15]
 [ 88 69]
 [ 93 14]
 [ 93 90]
 [ 97 32]
 [ 97 86]
 [ 98 15]
 [ 98 88]
 [ 99 39]
 [ 99 97]
 [101 24]
 [101 68]
 [103 17]
 [103 85]
 [103 23]
 [103 69]
 [113 8]
 [113 91]
 [120 16]
 [120 79]
 [126 28]
 [126 74]
 [137 18]
 [137 83]]

```

Choosing the number of clusters parameter

✓ WCSS- WITHIN CLUSTER SUM OF SQUARES

Finding wcss value for different number of clusters

Code:

```

wcss=[]
for i in range(1,11):
    kmeans = KMeans(n_clusters=i, init='k-means++', random_state=42)
    kmeans.fit(x)
    wcss.append(kmeans.inertia_)

```

Output:

```
#finding wcss value for different number of clusters

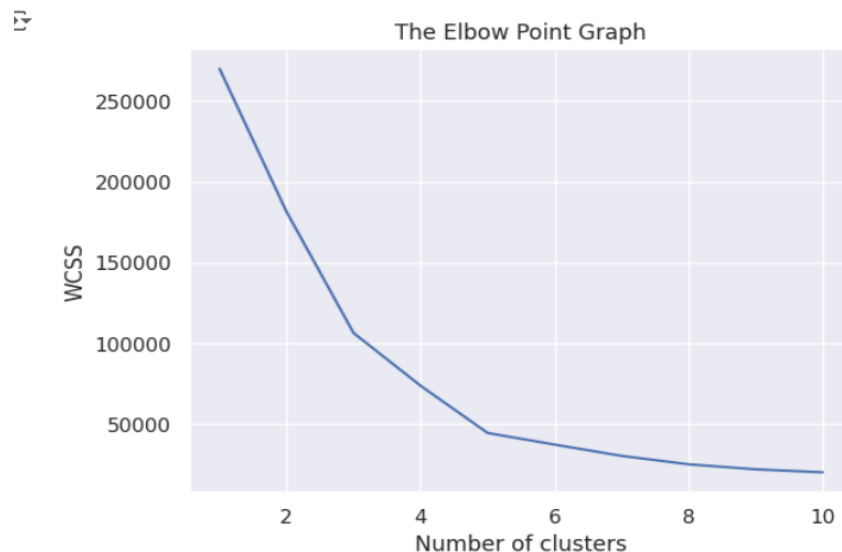
wcss=[]
for i in range(1,11):
    kmeans = KMeans(n_clusters=i, init='k-means++', random_state=42)
    kmeans.fit(x)
    .....
    wcss.append(kmeans.inertia_)
```

Plot an elbow graph to find which cluster has minimum value

Code:

```
sns.set()
plt.plot(range(1,11),wcss)
plt.title('The Elbow Point Graph')
plt.xlabel('Number of clusters')
plt.ylabel('WCSS')
plt.show()
```

Output :



The optimum number of clusters are 5 from elbow point graph

-Training the kmeans clustering model

(Value of k=5)

Code:

```
kmeans=KMeans(n_clusters=5, init='k-means++', random_state=0)
```

```
#return a label for each data point based on their cluster
```

```
y=kmeans.fit_predict(x)
```

```
print(y)
```

Output:

[illegible]

There are 5 clusters= 0,1,2,3,4

Visualize the clusters

Code:

```
#plotting all the clusters and their centriods
```

```
plt.figure(figsize=(8,8))
plt.scatter(x[y==0,0], x[y==0,1], s=50, c='green', label='Cluster1')
plt.scatter(x[y==1,0], x[y==1,1], s=50, c='red', label='Cluster2')
plt.scatter(x[y==2,0], x[y==2,1], s=50, c='pink', label='Cluster3')
plt.scatter(x[y==3,0], x[y==3,1], s=50, c='violet', label='Cluster4')
plt.scatter(x[y==4,0], x[y==4,1], s=50, c='black', label='Cluster5')
```

```
#plot the centroids
plt.scatter(kmeans.cluster_centers_[:,0], kmeans.cluster_centers_[:,1], s=100, c='yellow',
label='Centroids')
```

```
plt.title('Customer Groups')
plt.xlabel('Annual Income')
plt.ylabel('Spending Score')
plt.show()
```

Output:



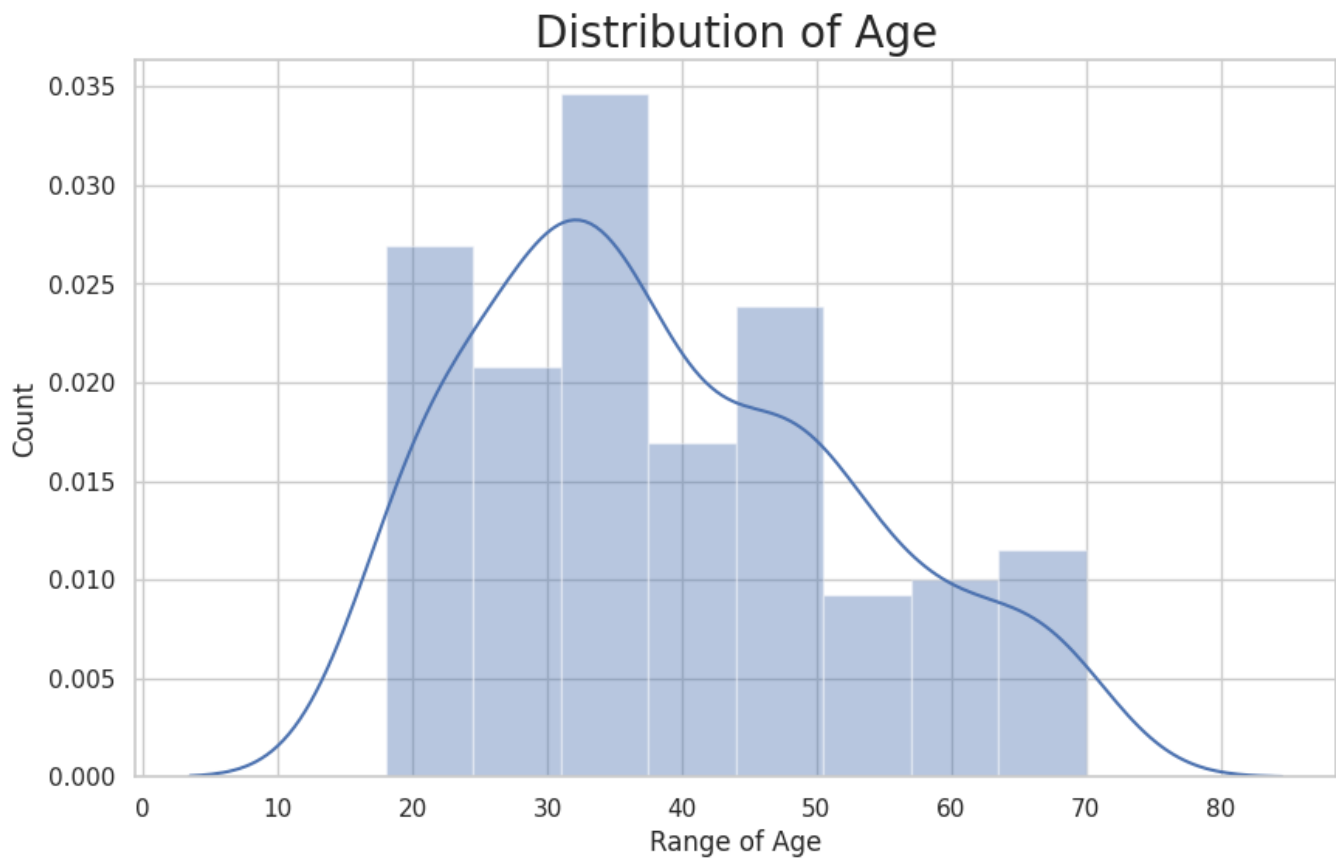
Visualization on Age distribution

Using Distplot

Code:

```
plt.figure(figsize=(10, 6))
sns.set(style = 'whitegrid')
sns.distplot(customer_Data['Age'])
plt.title('Distribution of Age', fontsize = 20)
plt.xlabel('Range of Age')
plt.ylabel('Count')
```

Output:



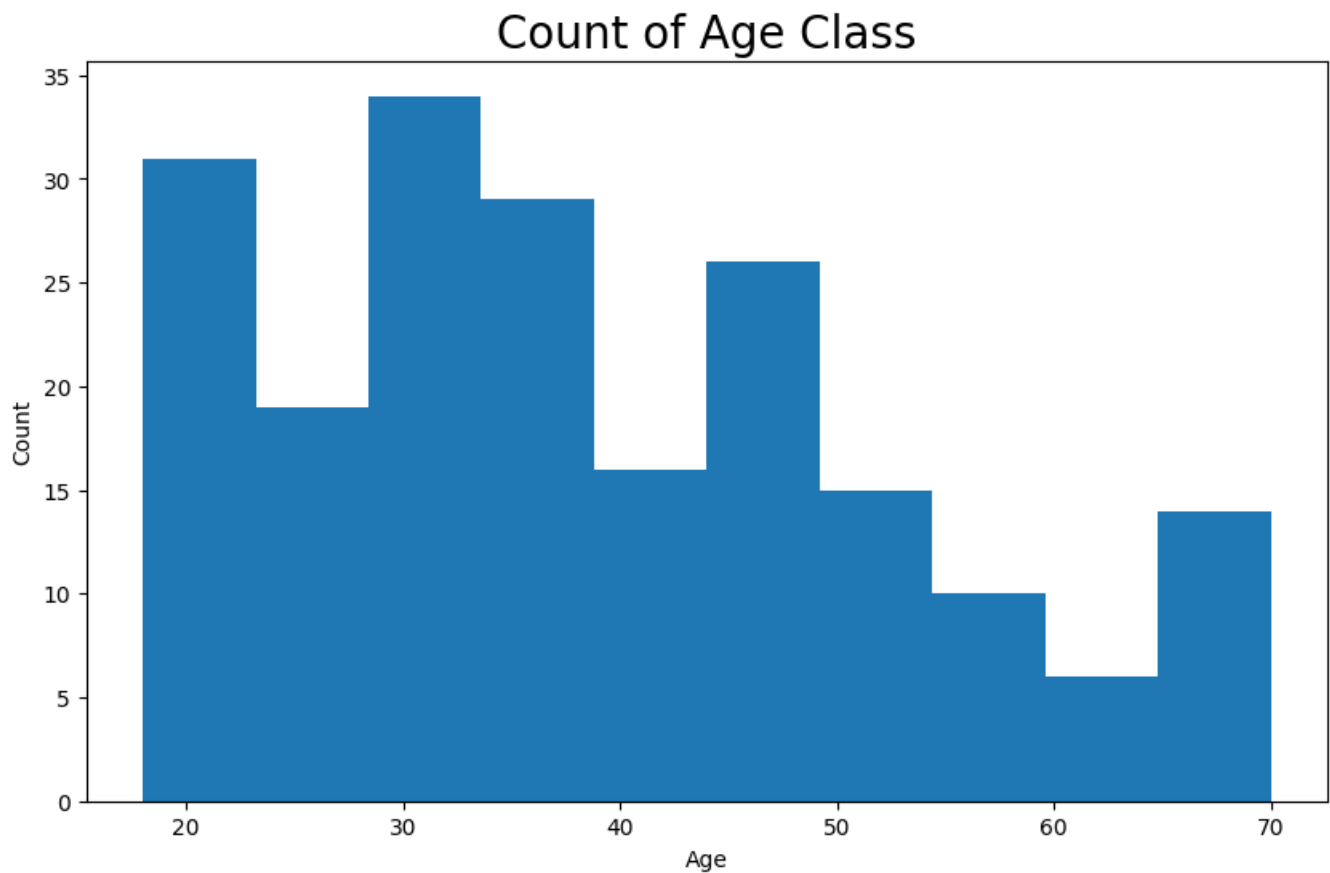
Using Histogram

Code:

```
#histogram to show count of Age class
plt.figure(figsize=(10, 6))
plt.hist(customer_Data['Age'], bins=10)
plt.title('Count of Age Class', fontsize = 20)
plt.xlabel('Age')
plt.ylabel('Count')
```

Output:

Text(0, 0.5, 'Count')



Descriptive Analysis of Age

Using Box Plot

Code:

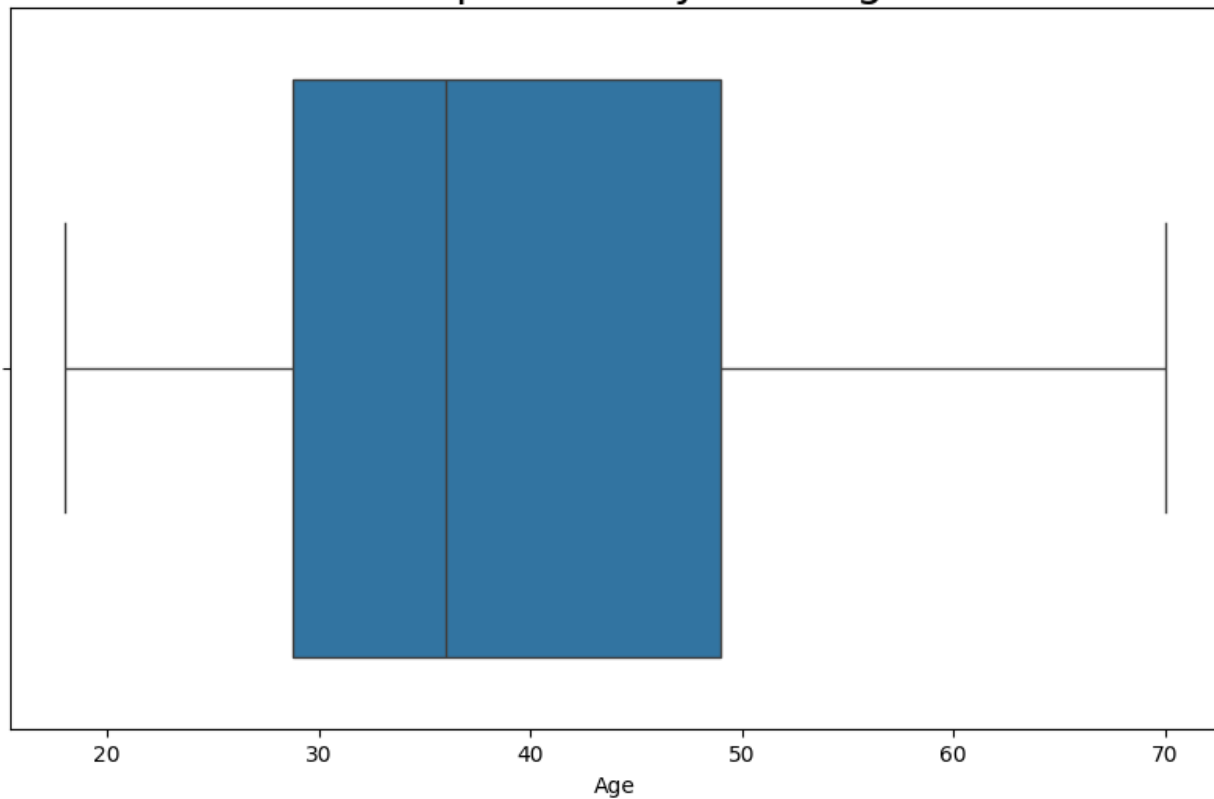
```
#visualize using boxplot for descriptive analysis of age
```

```
import seaborn as sns
plt.figure(figsize=(10, 6))
sns.boxplot(x='Age', data=customer_Data)
plt.title('Descriptive Analysis of Age', fontsize = 20)
plt.xlabel('Age')
```

Output:

```
Text(0.5, 0, 'Age')
```

Descriptive Analysis of Age



Customer Gender Visualization

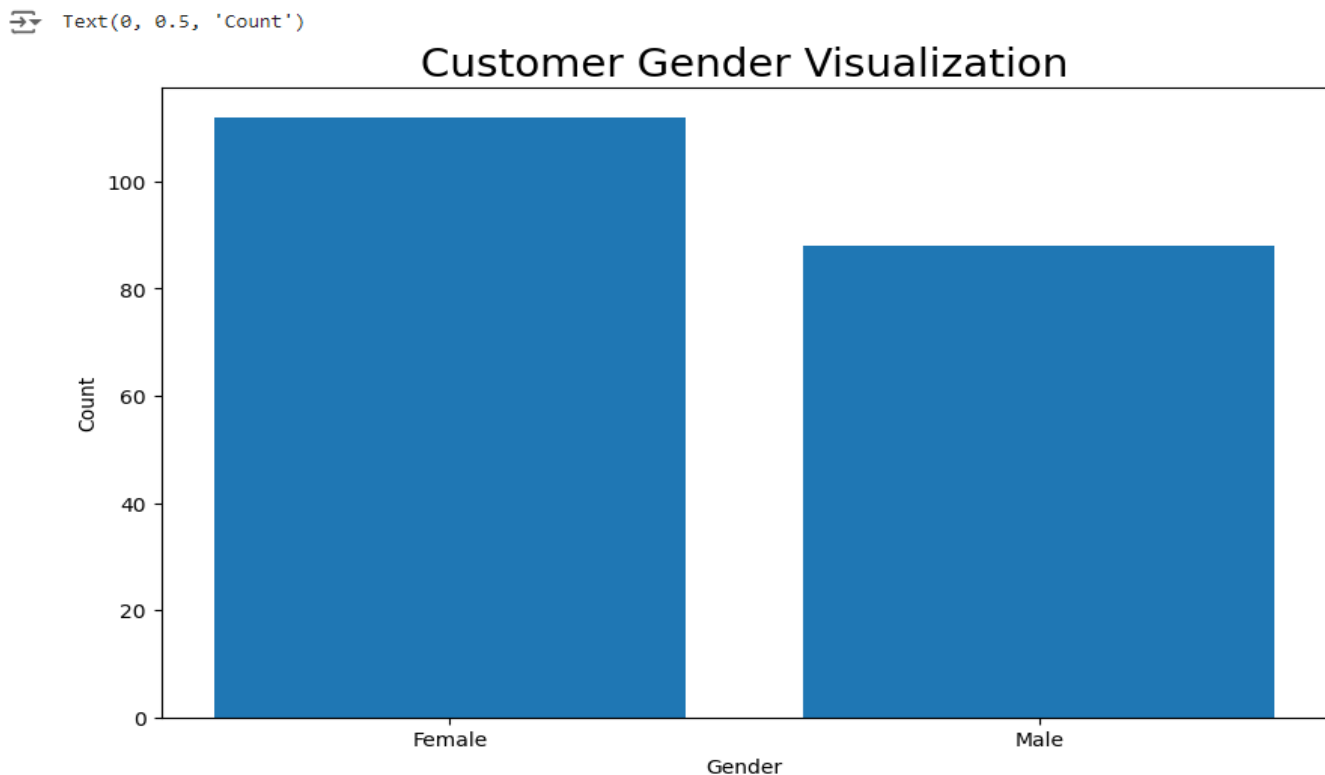
we will create a barplot and a piechart to show the gender distribution across our customer_data dataset.

Using Bar plot

Code:

```
import matplotlib.pyplot as plt
import pandas as pd
customer_Data=pd.read_csv('Mall_Customers.csv')
plt.figure(figsize=(10, 6))
plt.bar(customer_Data['Gender'].value_counts().index, customer_Data['Gender'].value_counts())
plt.title('Customer Gender Visualization', fontsize = 20)
plt.xlabel('Gender')
plt.ylabel('Count')
```

Output:



Inference:- From the above bar plot, we observe that the number of females is higher than the males.

Using Pie Chart

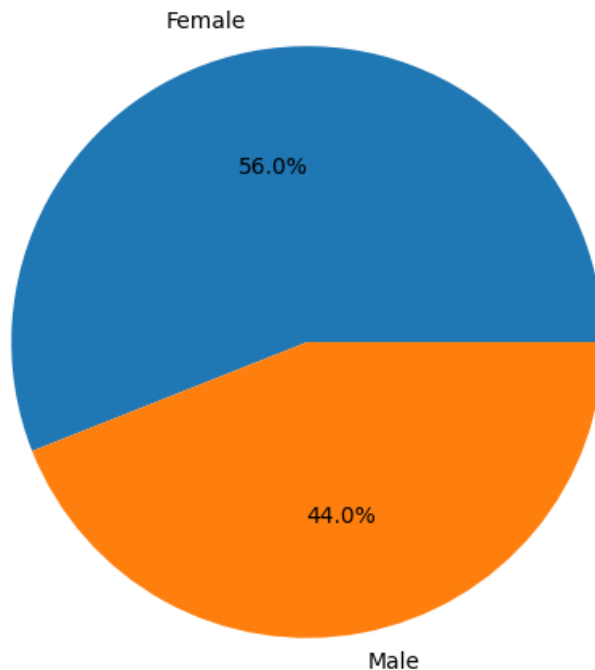
Code:

```
%%') #visualize a pie chart to observe the ratio of male and female distribution
```

```
plt.figure(figsize=(10, 6))  
plt.pie(customer_Data['Gender'].value_counts(), labels=customer_Data['Gender'].value_counts().index,  
autopct='%1.1f%%')
```

Output:

```
[<matplotlib.patches.Wedge at 0x78fa2f49b370>,  
<matplotlib.patches.Wedge at 0x78fa2f49bb80>],  
[Text(-0.20611945413751356, 1.080515974257694, 'Female'),  
Text(0.20611945413751367, -1.080515974257694, 'Male')],  
[Text(-0.11242879316591647, 0.5893723495951058, '56.0%'),  
Text(0.11242879316591654, -0.5893723495951058, '44.0%')]]
```



Inference:- From the above graph, we can conclude that the percentage of females is **56%**, whereas the percentage of male in the customer dataset is **44%**.

Analysis of the Annual Income of the Customers

we will create visualizations to analyse the annual income of the customers. We will plot a histogram and then we will proceed to examine this data using a density plot.

Using Histogram

Code:

```
#histogram for annual income
```

```
plt.figure(figsize=(10, 6))
```

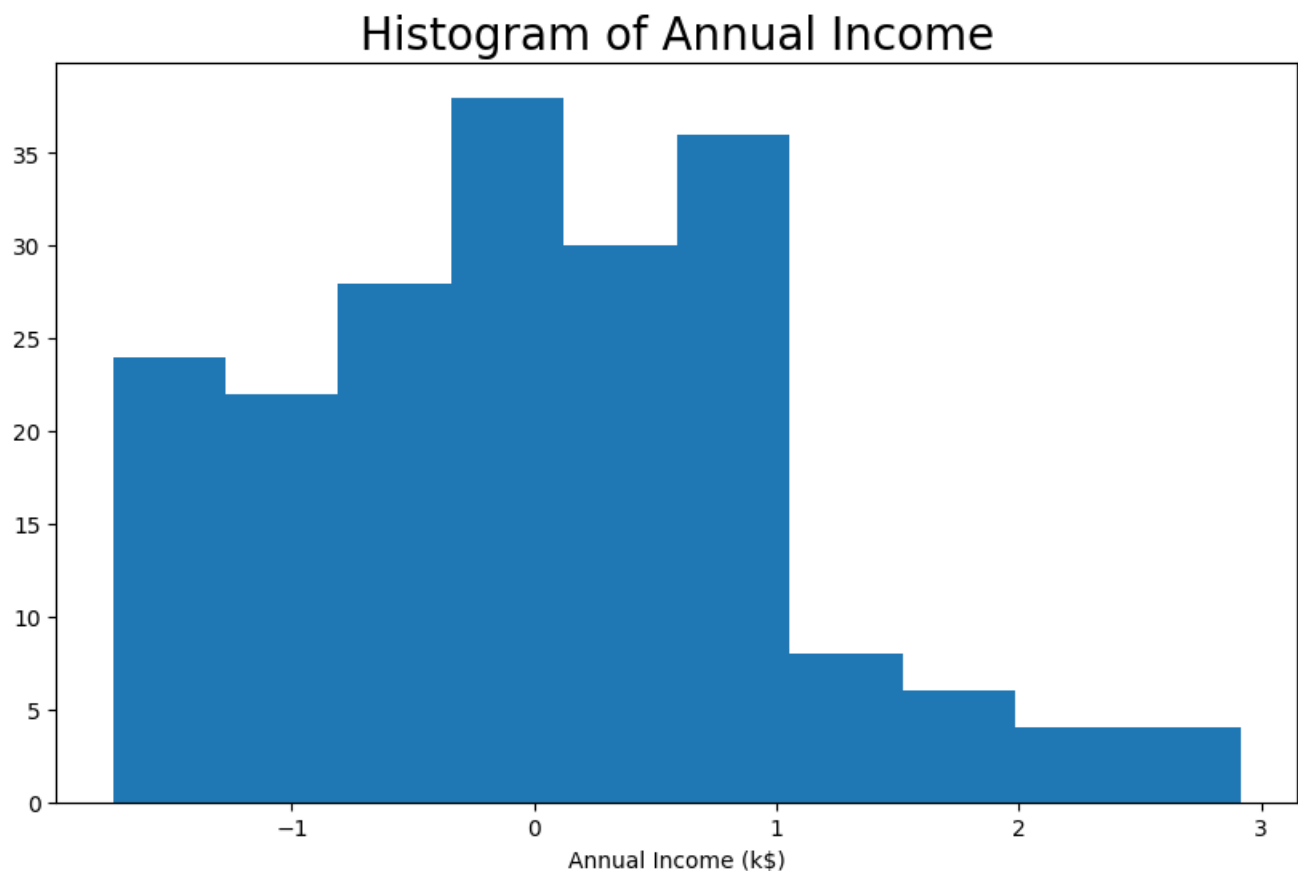
```
plt.hist(customer_Data['Annual Income (k$)'], bins=10)
```

```
plt.title('Histogram of Annual Income', fontsize = 20)
```

```
plt.xlabel('Annual Income (k$)')
```

Output:

```
Text(0.5, 0, 'Annual Income (k$)')
```



Using Box plot

Code:

```
#histogram for each variable for annual income
```

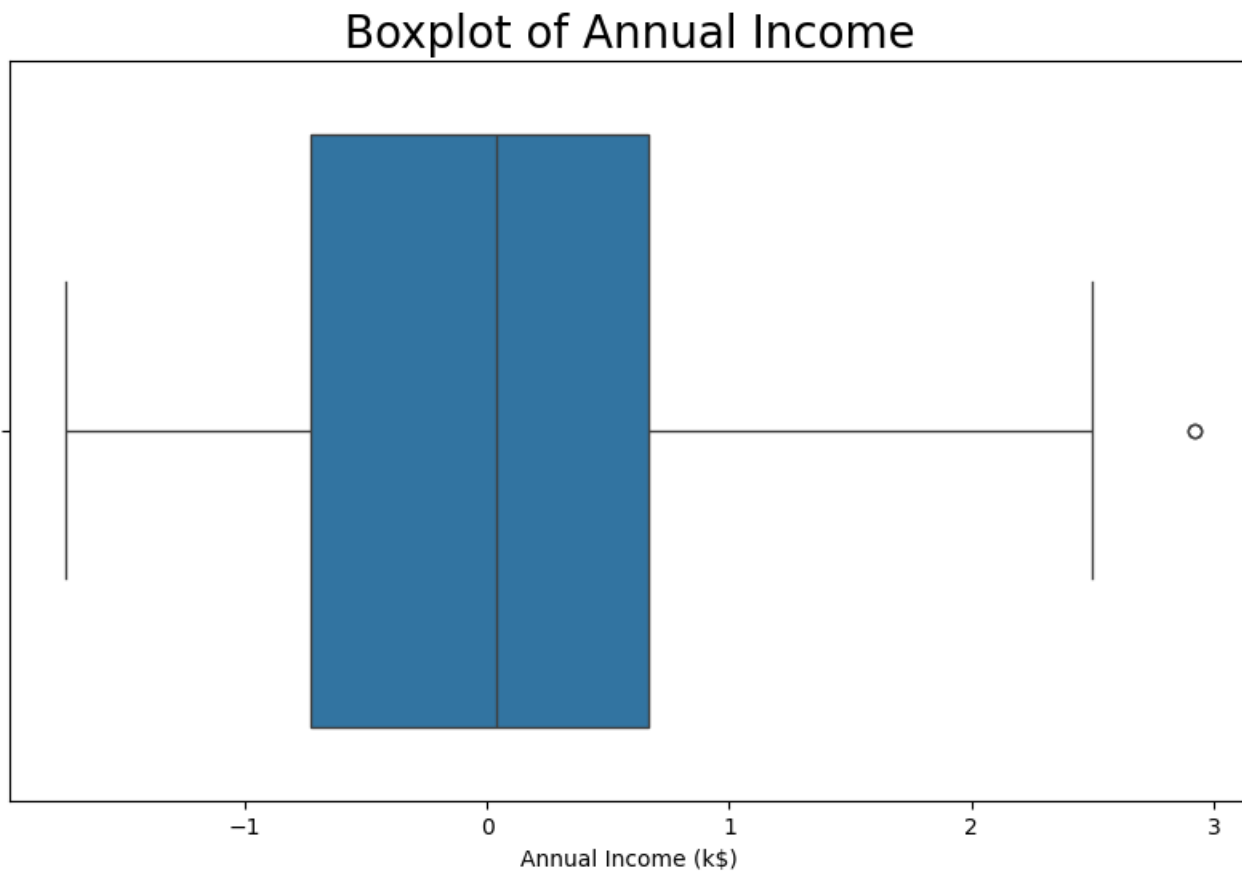
```
plt.figure(figsize=(10, 6))
```

```
sns.boxplot(x='Annual Income (k$)', data=customer_Data)
```

```
plt.title('Boxplot of Annual Income', fontsize = 20)
```

Output:

```
Text(0.5, 1.0, 'Boxplot of Annual Income')
```



Using Density Plot

Code:

```
#density plot for annual income
```

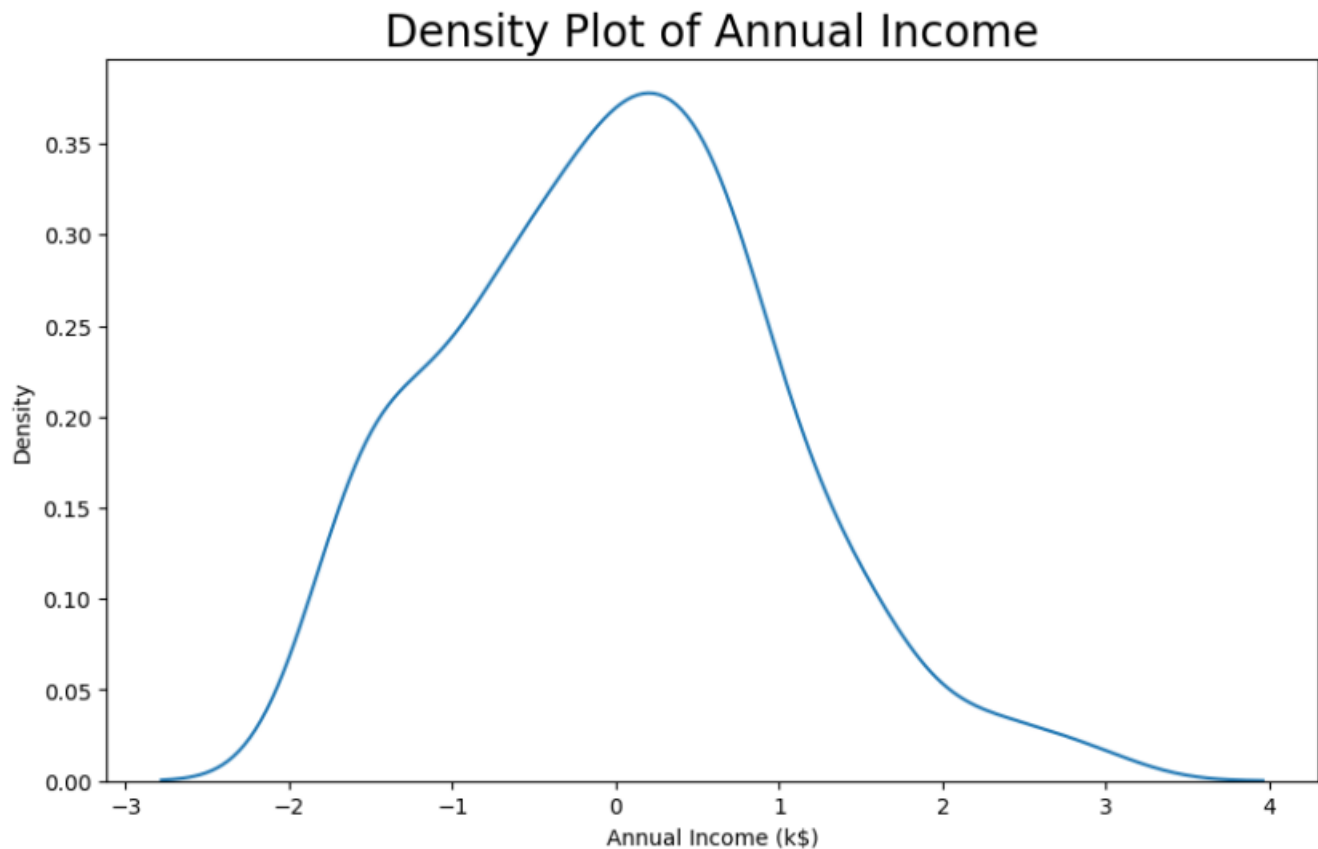
```
plt.figure(figsize=(10, 6))
```

```
sns.kdeplot(customer_Data['Annual Income (k$)'])
```

```
plt.title('Density Plot of Annual Income', fontsize = 20)
```

Output:

```
↔ Text(0.5, 1.0, 'Density Plot of Annual Income')
```



Analyzing Spending Score of the Customers


Using Boxplot

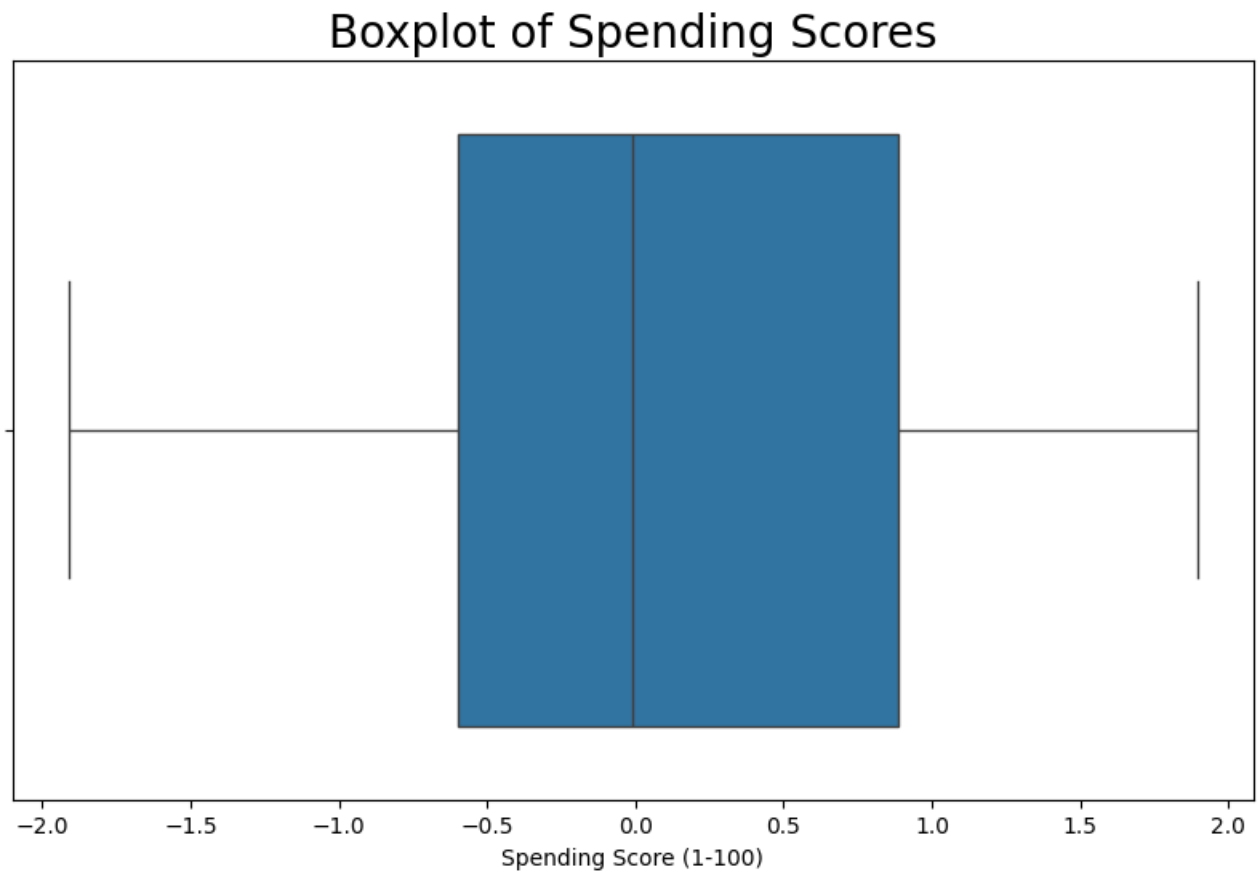
Code:

```
#boxplot for spending scores of customer

plt.figure(figsize=(10, 6))
sns.boxplot(x='Spending Score (1-100)', data=customer_Data)
plt.title('Boxplot of Spending Scores', fontsize = 20)
```

Output:

 Text(0.5, 1.0, 'Boxplot of Spending Scores')



Using Histogram

Code:

```
#histogram for spending scores
```

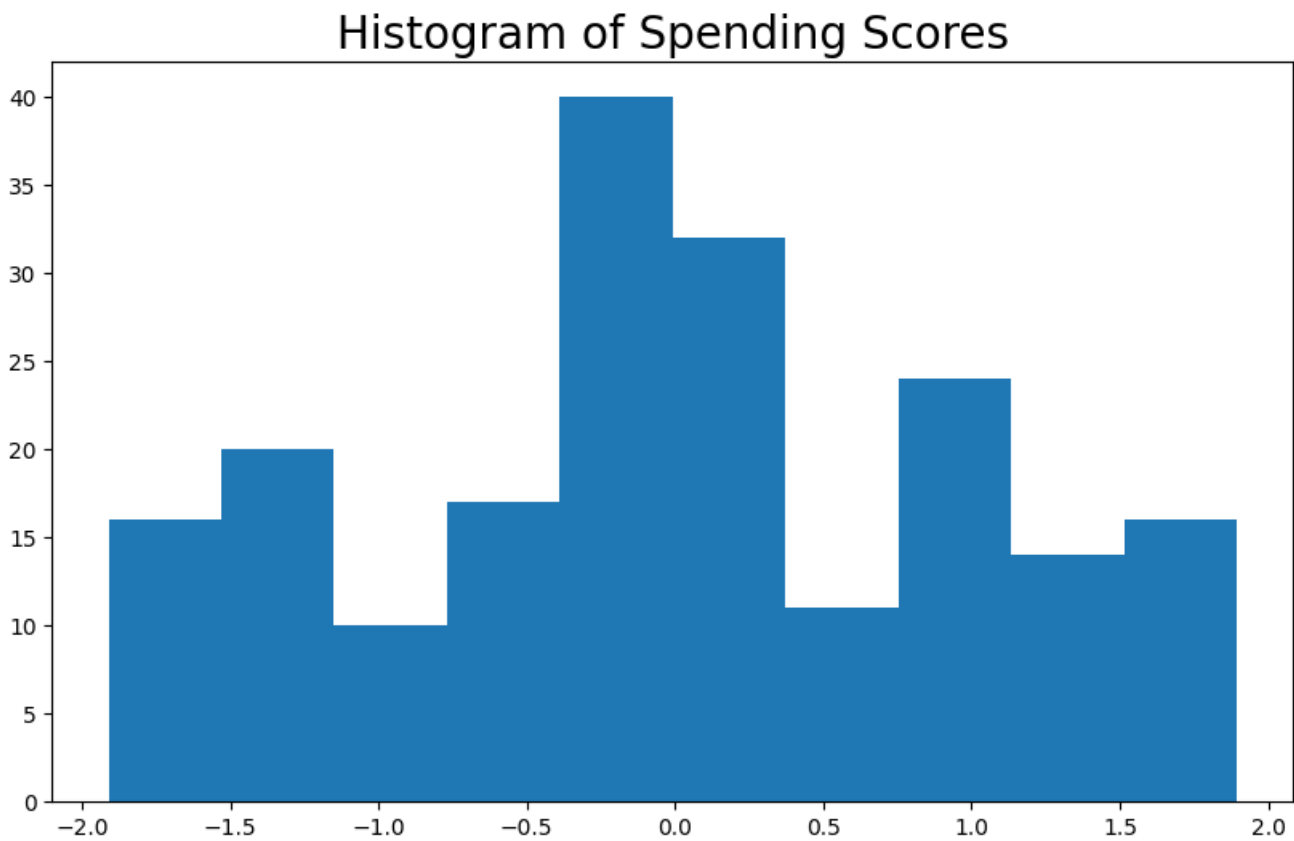
```
plt.figure(figsize=(10, 6))
```

```
plt.hist(customer_Data['Spending Score (1-100)'], bins=10)
```

```
plt.title('Histogram of Spending Scores', fontsize = 20)
```

Output:

```
➡ Text(0.5, 1.0, 'Histogram of Spending Scores')
```



Normalising the numerical features

Code:

```
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
customer_Data[['Annual Income (k$)', 'Spending Score (1-100)']] =
scaler.fit_transform(customer_Data[['Annual Income (k$)', 'Spending Score (1-100)']])
```

Encode categorical variables

Code:

```
#encode categorical variables

from sklearn.preprocessing import LabelEncoder
encoder = LabelEncoder()
customer_Data['Gender'] = encoder.fit_transform(customer_Data['Gender'])
```

Provide insights for marketing strategies

Code:

```
customer_Data.groupby('Gender').mean()
```

Output:

```
customer_Data.groupby('Gender').mean()
```

	CustomerID	Age	Annual Income (k\$)	Spending Score (1-100)
Gender				
0	97.562500	38.098214	-0.050002	0.051508
1	104.238636	39.806818	0.063639	-0.065555

Providing a complete report and analysis

Code:

```
customer_Data.describe()
```

Output:

customer_Data.describe()

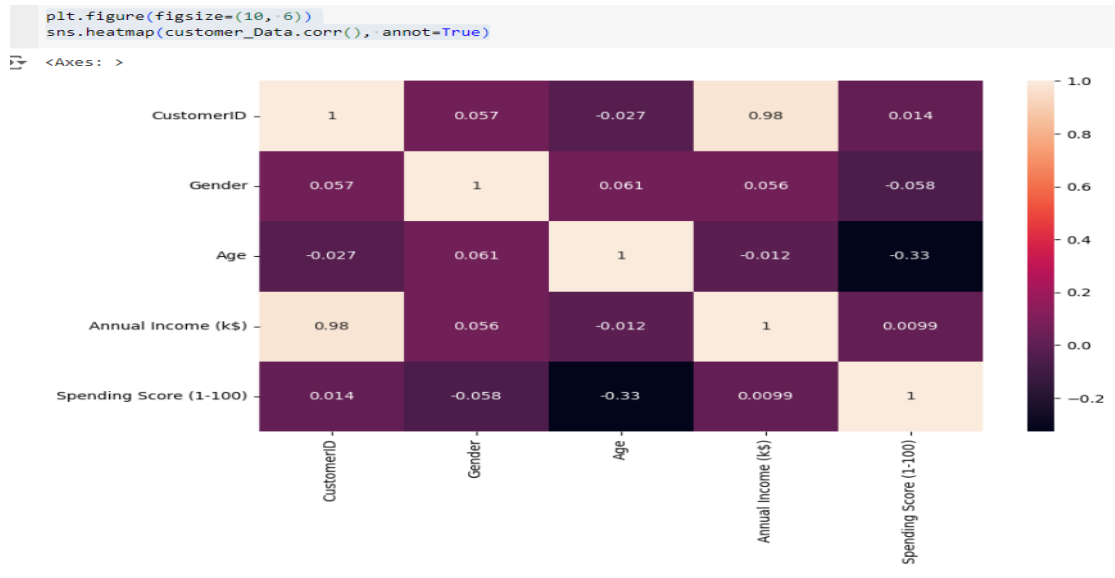
	CustomerID	Gender	Age	Annual Income (k\$)	Spending Score (1-100)
count	200.000000	200.000000	200.000000	2.000000e+02	2.000000e+02
mean	100.500000	0.440000	38.850000	7.105427e-17	7.549517e-17
std	57.879185	0.497633	13.969007	1.002509e+00	1.002509e+00
min	1.000000	0.000000	18.000000	-1.738999e+00	-1.910021e+00
25%	50.750000	0.000000	28.750000	-7.275093e-01	-5.997931e-01
50%	100.500000	0.000000	36.000000	3.587926e-02	-7.764312e-03
75%	150.250000	1.000000	49.000000	6.656748e-01	8.851316e-01
max	200.000000	1.000000	70.000000	2.917671e+00	1.894492e+00

Providing a report using Heatmap

Code:

```
plt.figure(figsize=(10, 6))
sns.heatmap(customer_Data.corr(), annot=True)
```

Output:



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.

Summing up the K-means clustering –

- We specify the number of clusters that we need to create.
- The algorithm selects k objects at random from the dataset. This object is the initial cluster or mean.
- The closest centroid obtains the assignment of a new observation. We base this assignment on the Euclidean Distance between object and the centroid.
- k clusters in the data points update the centroid through calculation of the new mean values present in all the data points of the cluster. The kth cluster's centroid has a length of p that contains means of all variables for observations in the k-th cluster. We denote the number of variables with p.
- Iterative minimization of the total within the sum of squares. Then through the iterative minimization of the total sum of the square, the assignment stop wavering when we achieve maximum iteration. The default value is 10 that the R software uses for the maximum iterations.

Determining Optimal Clusters

While working with clusters, we need to specify the number of clusters to use. We would like to utilize the optimal number of clusters. There are three popular methods to determine the optimal clusters –

1. Elbow method
2. Silhouette method
3. Gap statistic

1. Elbow Method

The main goal behind cluster partitioning methods like k-means is to define the clusters such that the intra-cluster variation stays minimum.

minimize (sum $W(C_k)$), $k=1\dots k$

Where C_k represents the kth cluster and $W(C_k)$ denotes the intra-cluster variation. With the measurement of the total intra-cluster variation, one can evaluate the compactness of the clustering boundary. We can then proceed to define the optimal clusters as follows –

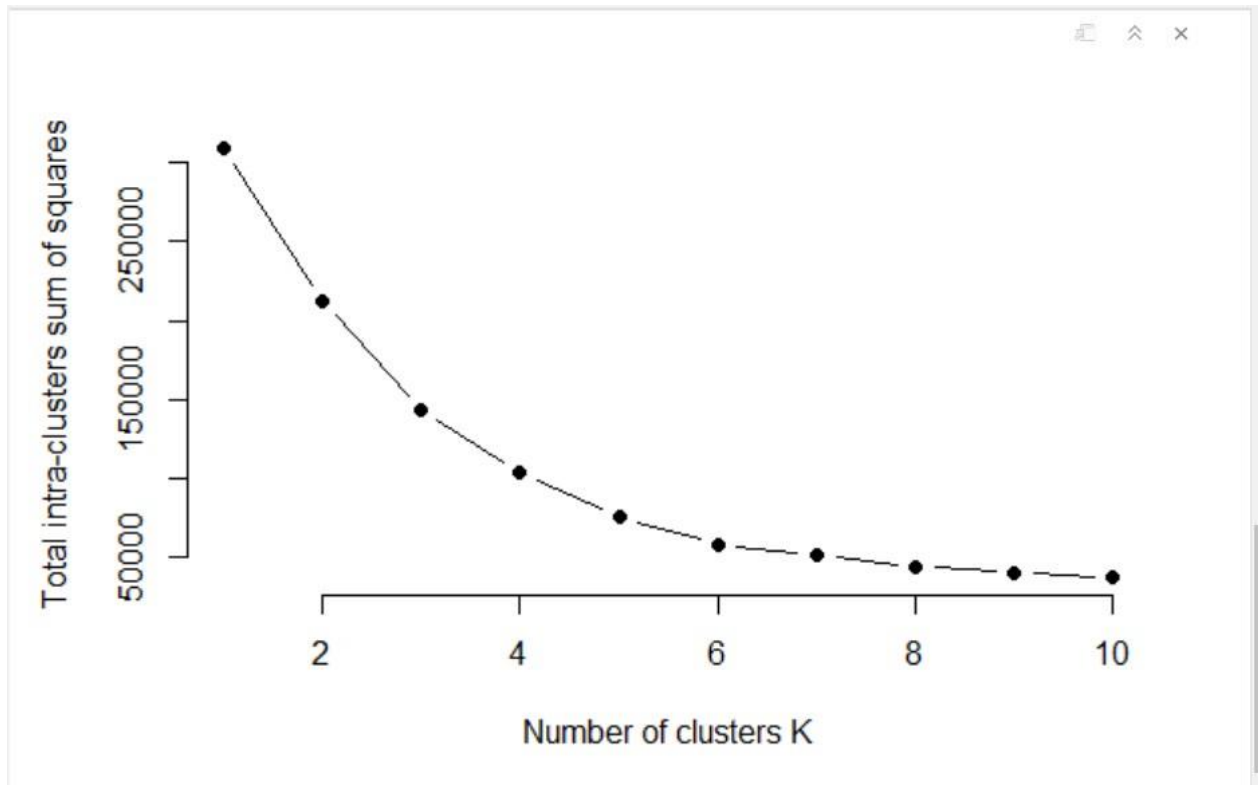
First, we calculate the clustering algorithm for several values of k. This can be done by creating a variation within k from 1 to 10 clusters. We then calculate the total intra-cluster sum of square (iss). Then, we proceed to plot iss based on the number of k clusters. This plot denotes the appropriate number of clusters required in our model. In the plot, the location of a bend or a knee is the indication of the optimum number of clusters.

Code

```
library(purrr)
set.seed(123)
# function to calculate total intra-cluster sum of square
iss <- function(k) {
  kmeans(customer_data[,3:5],k,iter.max=100,nstart=100,algorithm="Lloyd" )$tot.withinss
}
k.values <- 1:10
iss_values <- map_dbl(k.values, iss)
plot(k.values, iss_values,
     type="b", pch = 19, frame = FALSE,
     xlab="Number of clusters K",
     ylab="Total intra-clusters sum of squares")
```

Output

```
library(purrr)
set.seed(123)
# function to calculate total intra-cluster sum of square
iss <- function(k) {
  kmeans(customer_data[,3:5],k,iter.max=100,nstart=100,algorithm="Lloyd"
)$tot.withinss
}
k.values <- 1:10
iss_values <- map_dbl(k.values, iss)
plot(k.values, iss_values,
     type="b", pch = 19, frame = FALSE,
     xlab="Number of clusters K",
     ylab="Total intra-clusters sum of squares")
```



Inference: From the above graph, we can conclude that 4 is the appropriate number of clusters since it seems to be appearing at the bend in the elbow plot.

2. Average Silhouette Method

With the help of the average silhouette method, we can measure the quality of our clustering operation. With this, we can determine how well within the cluster is the data object. If we obtain a high average silhouette width, it means that we have good clustering. The average silhouette method calculates the mean of silhouette observations for different k values. With the optimal number of k clusters, one can maximize the average silhouette over significant values for k clusters.

Using the silhouette function in the cluster package, we can compute the average silhouette width using the kmean function. Here, the optimal cluster will possess highest average.

Code

```
library(cluster)
library(gridExtra)
library(grid)
```

When n=2

```
k2<-kmeans(customer_data[,3:5],2,iter.max=100,nstart=50,algorithm="Lloyd")
```

```
s2<-plot(silhouette(k2$cluster,dist(customer_data[,3:5],"euclidean")))
```

```
library(cluster)
library(gridExtra)
library(grid)
k2<-kmeans(customer_data[,3:5],2,iter.max=100,nstart=50,algorithm="Lloyd")
s2<-plot(silhouette(k2$cluster,dist(customer_data[,3:5],"euclidean")))
```

Silhouette plot of (x = k2\$cluster, dist = dist(customer_data[, 3:5],

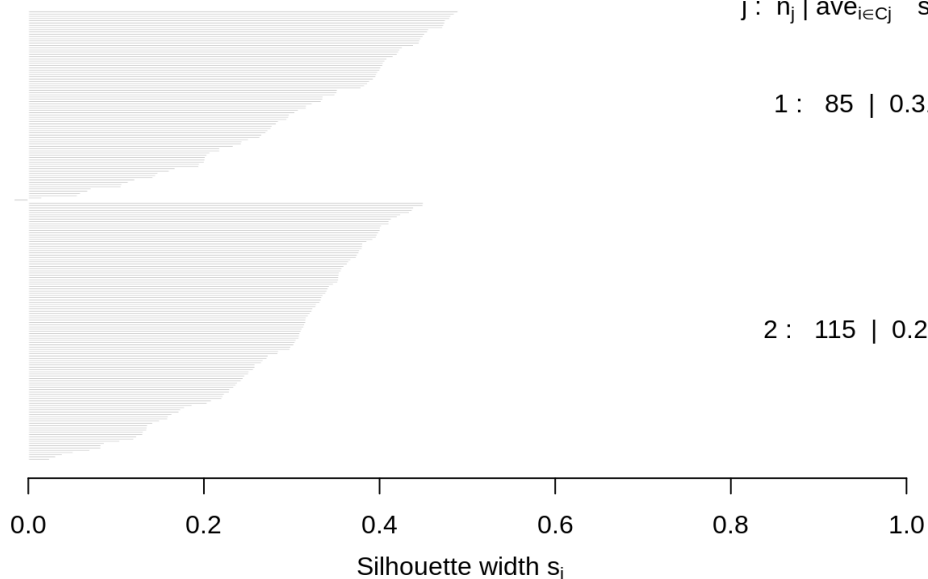
n = 200

2 clusters C_j

$j : n_j \mid \text{ave}_{i \in C_j} s_i$

1 : 85 | 0.31

2 : 115 | 0.28



Average silhouette width : 0.29

When n=3

```
k3<-kmeans(customer_data[,3:5],3,iter.max=100,nstart=50,algorithm="Lloyd")
```

```
s3<-plot(silhouette(k3$cluster,dist(customer_data[,3:5],"euclidean")))
```

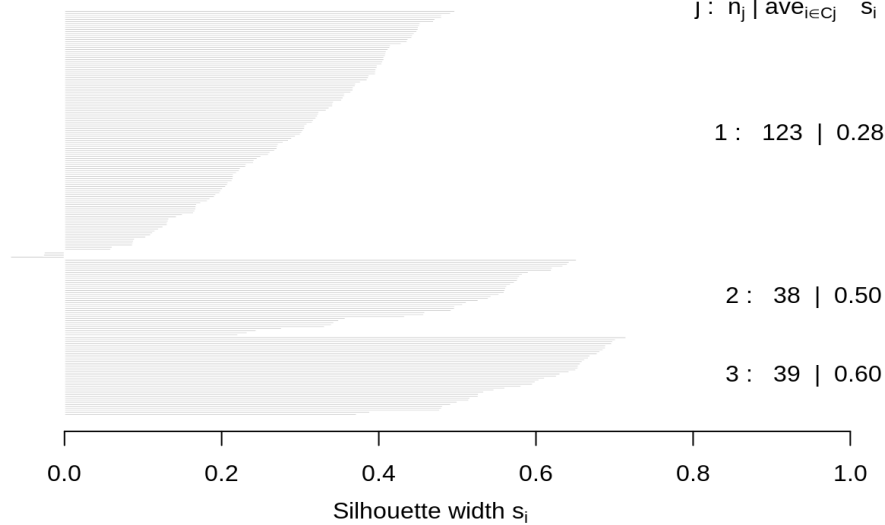
```
library(cluster)
library(gridExtra)
library(grid)
k3<-kmeans(customer_data[,3:5],3,iter.max=100,nstart=50,algorithm="Lloyd")
s3<-plot(silhouette(k3$cluster,dist(customer_data[,3:5],"euclidean")))
```

Silhouette plot of (x = k3\$cluster, dist = dist(customer_data[, 3:5],

n = 200

3 clusters C_j

$j : n_j \mid \text{ave}_{i \in C_j} s_i$



When n=4

```
k4<-kmeans(customer_data[,3:5],4,iter.max=100,nstart=50,algorithm="Lloyd")
```

```
s4<-plot(silhouette(k4$cluster,dist(customer_data[,3:5],"euclidean")))
```

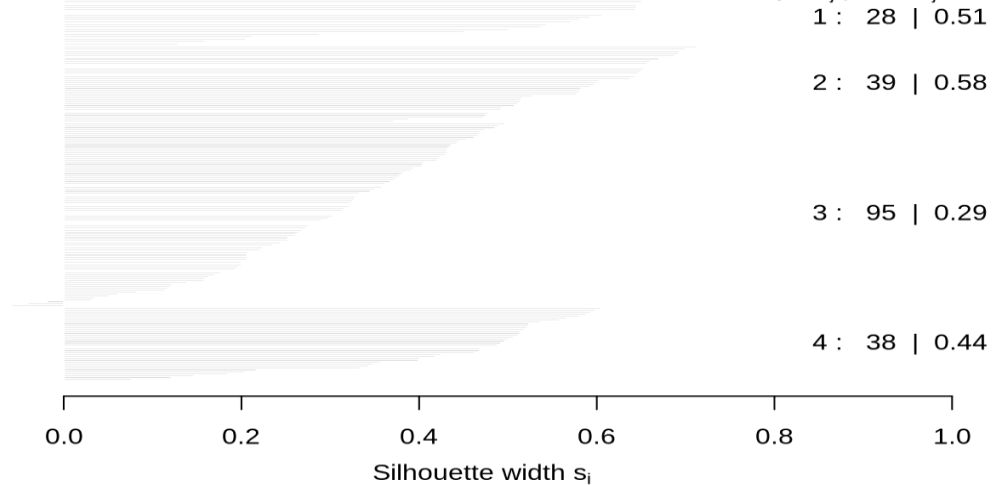
```
k3<-kmeans(customer_data[,3:5],4,iter.max=100,nstart=50,algorithm="Lloyd")
s3<-plot(silhouette(k3$cluster,dist(customer_data[,3:5],"euclidean")))
```

Silhouette plot of (x = k4\$cluster, dist = dist(customer_data[, 3:5],

n = 200

4 clusters C_j

$j : n_j \mid \text{ave}_{i \in C_j} s_i$



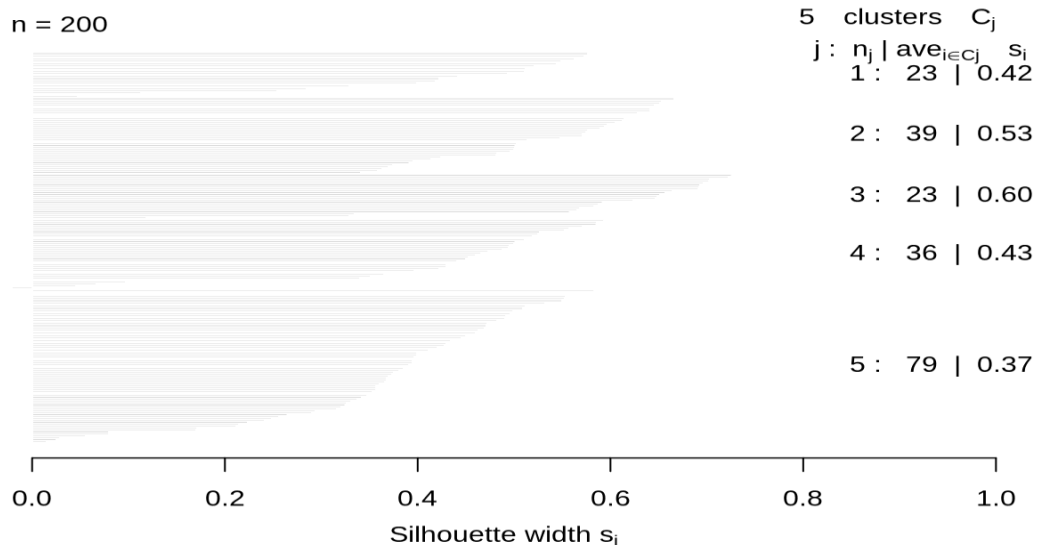
When n=5

```
k5<-kmeans(customer_data[,3:5],5,iter.max=100,nstart=50,algorithm="Lloyd")
s5<-plot(silhouette(k5$cluster,dist(customer_data[,3:5],"euclidean")))
```

```
k3<-kmeans(customer_data[,3:5],5,iter.max=100,nstart=50,algorithm="Lloyd")
s3<-plot(silhouette(k3$cluster,dist(customer_data[,3:5],"euclidean")))
```

Silhouette plot of (x = k5\$cluster, dist = dist(customer_data[, 3:5],

n = 200

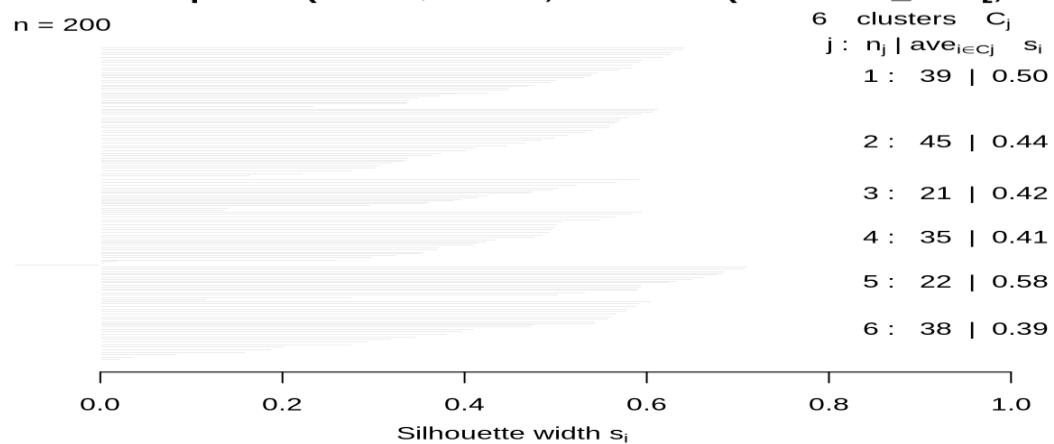


When n=6

```
k6<-kmeans(customer_data[,3:5],6,iter.max=100,nstart=50,algorithm="Lloyd")
s6<-plot(silhouette(k6$cluster,dist(customer_data[,3:5],"euclidean")))
```

```
k3<-kmeans(customer_data[,3:5],6,iter.max=100,nstart=50,algorithm="Lloyd")
s3<-plot(silhouette(k3$cluster,dist(customer_data[,3:5],"euclidean")))
```

Silhouette plot of (x = k6\$cluster, dist = dist(customer_data[, 3:5],
n = 200

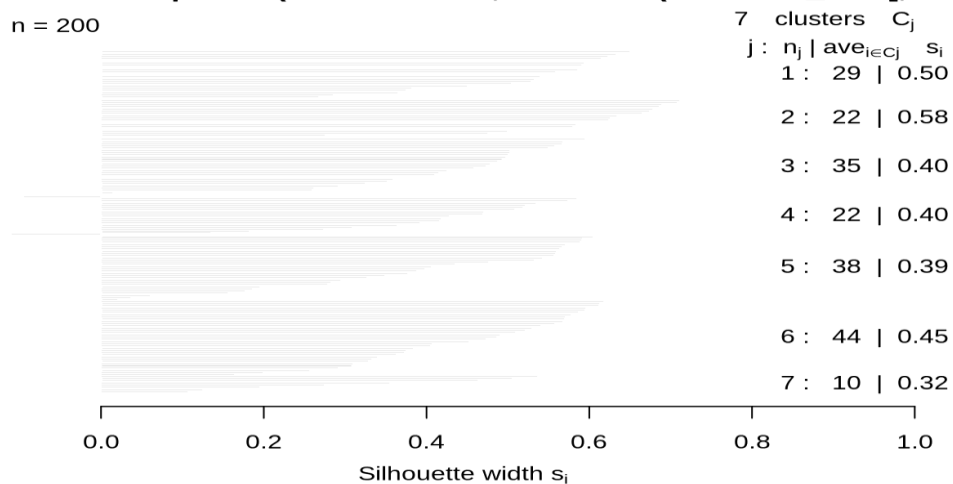


Average silhouette width : 0.45

When n=7

```
k7<-kmeans(customer_data[,3:5],7,iter.max=100,nstart=50,algorithm="Lloyd")
s7<-plot(silhouette(k7$cluster,dist(customer_data[,3:5],"euclidean"))))
k3<-kmeans(customer_data[,3:5],7,iter.max=100,nstart=50,algorithm="Lloyd")
s3<-plot(silhouette(k3$cluster,dist(customer_data[,3:5],"euclidean"))))
```

Silhouette plot of (x = k7\$cluster, dist = dist(customer_data[, 3:5],
n = 200

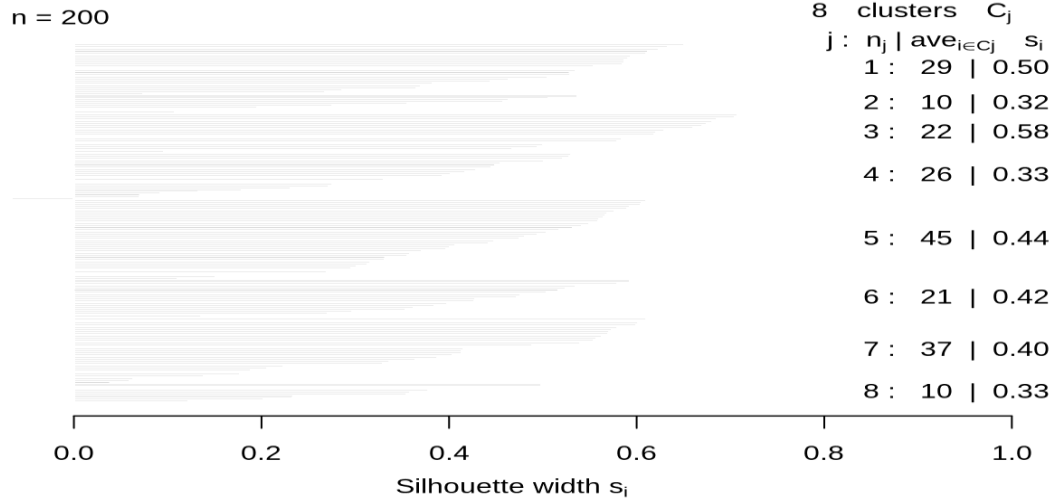


Average silhouette width : 0.44

When n=8

```
k8<-kmeans(customer_data[,3:5],8,iter.max=100,nstart=50,algorithm="Lloyd")
s8<-plot(silhouette(k8$cluster,dist(customer_data[,3:5],"euclidean"))))
k3<-kmeans(customer_data[,3:5],8,iter.max=100,nstart=50,algorithm="Lloyd")
s3<-plot(silhouette(k3$cluster,dist(customer_data[,3:5],"euclidean"))))
```


Silhouette plot of (x = k8\$cluster, dist = dist(customer_data[, 3:5],
n = 200

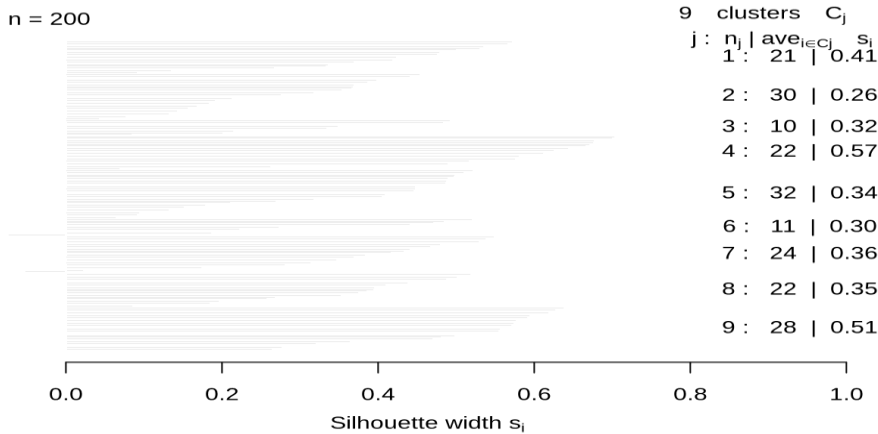


Average silhouette width : 0.43

When n=9

```
k9<-kmeans(customer_data[,3:5],9,iter.max=100,nstart=50,algorithm="Lloyd")
s9<-plot(silhouette(k9$cluster,dist(customer_data[,3:5],"euclidean"))))
k3<-kmeans(customer_data[, 3:5],9,iter.max=100,nstart=50,algorithm="Lloyd")
s3<-plot(silhouette(k3$cluster,dist(customer_data[, 3:5],"euclidean"))))
```

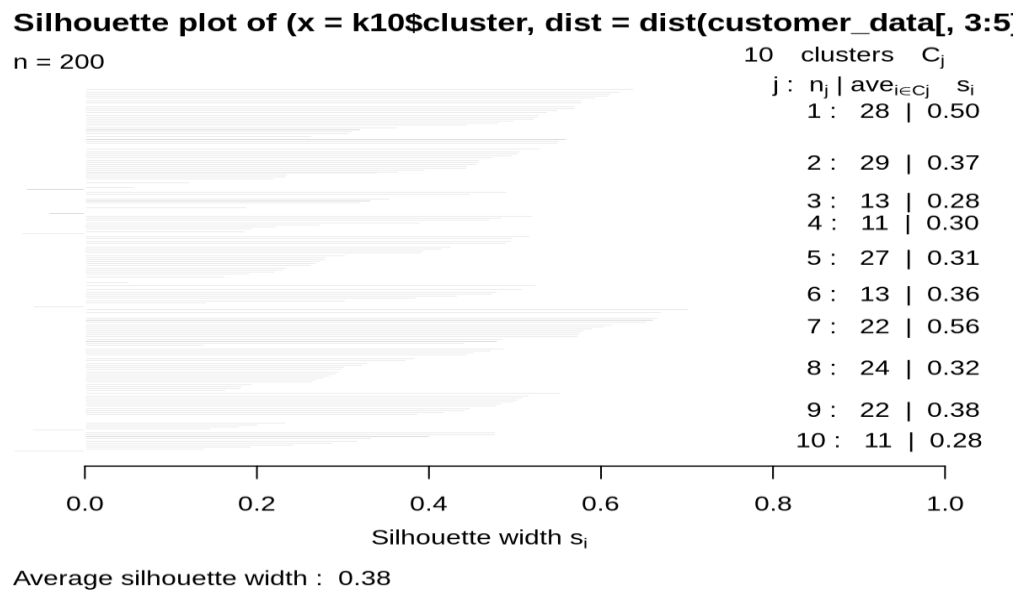
Silhouette plot of (x = k9\$cluster, dist = dist(customer_data[, 3:5],



Average silhouette width : 0.39

When n=10

```
k10<-kmeans(customer_data[,3:5],10,iter.max=100,nstart=50,algorithm="Lloyd")
s10<-plot(silhouette(k10$cluster,dist(customer_data[,3:5],"euclidean"))))
k3<-kmeans(customer_data[, 3:5],10,iter.max=100,nstart=50,algorithm="Lloyd")
s3<-plot(silhouette(k3$cluster,dist(customer_data[, 3:5],"euclidean"))))
```



Now, we make use of the `fviz_nbclust()` function to determine and visualize the optimal number of clusters as follows –

Code

```
library(NbClust)
library(factoextra)

fviz_nbclust(customer_data[,3:5], kmeans, method = "silhouette")
```

Output

```
library(NbClust)
library(factoextra)
fviz_nbclust(customer_data[,3:5], kmeans, method = "silhouette")
```

