

Major Project Synopsis
On
Customer Segmentation using Machine Learning
In partial fulfilment of requirements for the degree
Of
BACHELOR OF TECHNOLOGY
IN
COMPUTER SCIENCE & ENGINEERING

Submitted by :

Prachi Gothwal [20100BTCSDSI07284]

Sakshi Patel [20100BTCSDSI07291]

Under the guidance of

Prof. Om Kant Sharma



DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING
SHRI VAISHNAV INSTITUTE OF INFORMATION TECHNOLOGY
SHRI VAISHNAV VIDYAPEETH VISHWAVIDYALA, INDORE
JUL-DEC – 2022

SHRI VAISHNAV INSTITUTE OF INFORMATION TECHNOLOGY

INTRODUCTION

In this Data Science Project, we will perform one of the most essential applications of machine learning – Customer Segmentation. In this project, we will implement customer segmentation. Whenever we need to find your best customer, customer segmentation is the ideal methodology. Also, in this data science project, we will see the descriptive analysis of our data and then implement several versions of the K-means algorithm. So, follow the complete data science customer segmentation project using machine learning in R and become a pro in Data Science.

PROBLEM STATEMENT

Consider a mall a very famous mall and we are a very experienced data scientist and this mall wants all information about their customers and other expects. As a data scientist you can build a system that can cluster customer in different groups.

One group of customer are those who purchase more from that mall and other group are those who don't purchase too much from that mall. So having these group of customer these is easy to understand and better details to make and understand marketing strategies.



WORK FLOW

First we want to make a small customer data to train or instruct our machine learning project model. So first is getting those customer data and then we have to process these data, we cannot feed these data directly in machine learning model so we need to select key features that particular dataset contain is all come in analysis process and after that we need to choose the current no.

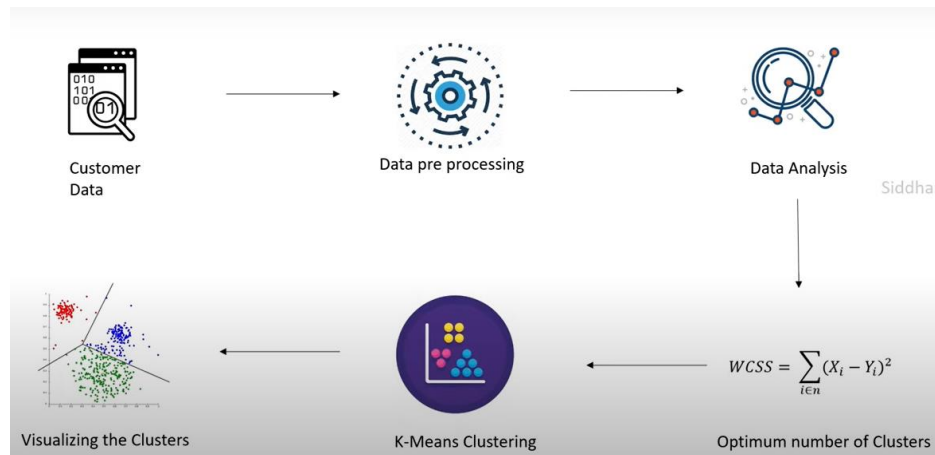
Of clusters and tell the machine learning project how many clusters are exists.

So we find the no. of clusters using method within cluster sum of square, so in this process we find the value of WCSS.

Once we have no. of cluster we can find the data in K-mean cluster Algorithm.

So once we feed this algorithm to this model it can group the data dependent on the similarities or similar expending pattern etc.

After that we can visualizing these cluster data by putting these data on the prediction made by the clustering models to get better insight of data.



K-Means Algorithm

- K-Means Clustering is an Unsupervised Learning algorithm, which groups the unlabeled dataset into different clusters.
- Here K defines the number of pre-defined clusters that need to be created in the process, as if K=2, there will be two clusters, and for K=3, there will be three clusters, and so on.
- It allows us to cluster the data into different groups and a convenient way to discover the categories of groups in the unlabeled dataset on its own without the need for any training.
- It is a centroid-based algorithm, where each cluster is associated with a centroid. The main aim of this algorithm is to minimize the sum of distances between the data point and their corresponding clusters.

Step-1: Data pre-processing Step

Importing the Dependencies

In the above code, the numpy we have imported for the performing mathematics calculation, **matplotlib** is for plotting the graph, and **pandas** are for managing the dataset.

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.cluster import KMeans
```

Importing dataset.

Next, we will import the dataset that we need to use. So here, we are using the Mall_Customer_data.csv dataset. It can be imported using the below code:

```
In [7]: #read dataset
print("Prachi Gothwal")
print("Sakshi Patel")
customer_data = pd.read_csv("C:\\Users\\Abhishek\\Desktop\\Mall_Customers.csv")

Prachi Gothwal
Sakshi Patel
```

Data Analysis

```
In [8]: print("Prachi Gothwal")
print("Sakshi Patel")
customer_data.head()
```

Prachi Gothwal
Sakshi Patel

Out[8]:

	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

```
In [9]: # finding the number of rows and columns
print("Prachi Gothwal")
print("Sakshi Patel")
customer_data.shape
```

Prachi Gothwal
Sakshi Patel

Out[9]: (200, 5)

```
[10]: print("Prachi Gothwal")
      print("Sakshi Patel")
      customer_data.info()
```

```
Prachi Gothwal
Sakshi Patel
<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
```

```
1 [12]: # checking for missing values
      print("Prachi Gothwal")
      print("Sakshi Patel")
      customer_data.isnull().sum()
```

```
Prachi Gothwal
Sakshi Patel
```

```
Out[12]: CustomerID            0
          Gender              0
          Age                 0
          Annual Income (k$)   0
          Spending Score (1-100) 0
          dtype: int64
```

Choosing the Annual Income Column & Spending Score column

```
In [13]: print("Prachi Gothwal")
print("Sakshi Patel")
X = customer_data.iloc[:,[3,4]].values
```

```
Prachi Gothwal
Sakshi Patel
```

```
In [14]: 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  88]]
```

Step-2: Finding the optimal number of clusters using the elbow method

In the second step, we will try to find the optimal number of clusters for our clustering problem.

- Choosing the number of clusters
- WCSS -> Within Clusters Sum of Squares

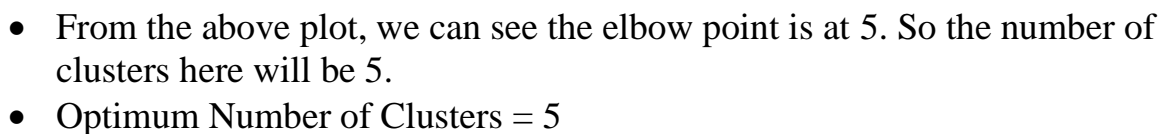
```
In [15]: # finding wcss value for different number of clusters
print("Prachi Gothwal")
print("Sakshi Patel")
wcss = []

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

```
Prachi Gothwal
Sakshi Patel
```

```
C:\Users\Abhishek\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1036: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
  warnings.warn(
```

Prachi Gothwal
Sakshi Patel



As we have got the number of clusters, so we can now train the model on the dataset.

[illegible]

7

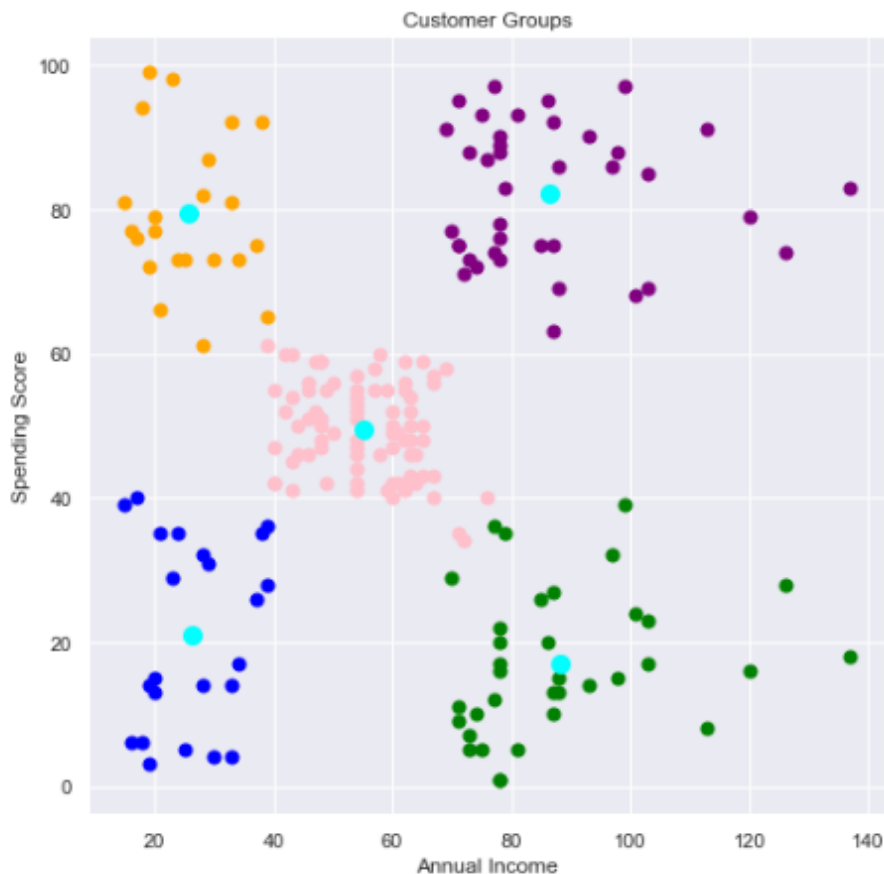
Step-4: Visualizing the Clusters

The last step is to visualize the clusters. As we have 5 clusters for our model, so we will visualize each cluster one by one.

```
# plotting all the clusters and their Centroids
print("Prachi Gothwal")
print("Sakshi Patel")
plt.figure(figsize=(8,8))
plt.scatter(X[Y==0,0], X[Y==0,1], s=50, c='green', label='Cluster 1')
plt.scatter(X[Y==1,0], X[Y==1,1], s=50, c='pink', label='Cluster 2')
plt.scatter(X[Y==2,0], X[Y==2,1], s=50, c='purple', label='Cluster 3')
plt.scatter(X[Y==3,0], X[Y==3,1], s=50, c='orange', label='Cluster 4')
plt.scatter(X[Y==4,0], X[Y==4,1], s=50, c='blue', label='Cluster 5')

# plot the centroids
plt.scatter(kmeans.cluster_centers_[0,0], kmeans.cluster_centers_[0,1], s=100, c='cyan', label='Centroids')
plt.title('Customer Groups')
plt.xlabel('Annual Income')
plt.ylabel('Spending Score')
plt.show()
```

Prachi Gothwal
Sakshi Patel



CONCLUSION

As we see we have multiple cluster here in different different color. All the cluster are portioned differently on the screen and the cluster represent the group of customer we are getting .As we focus on particular cluster it means a group of people we can see not too many have more annual income but they are regular customer of the mall except these others are not regular customers of mall.

So through these survey we can focus on that group of customer who are not buying to much from the mall, we can give them offers so they can be regular customers.

Through these survey we can improve the market value and the profit of mall and improve their status.