

COLLEGE CODE:5127

APPLIED DATA SCIENCE

**CUSTOMER SEGMENTATION USING DATA
SCIENCE**

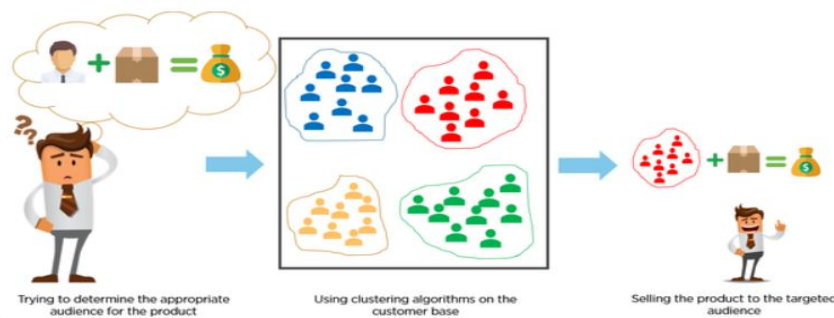
GROUP-5

OBJECTIVES:

- INTRODUCTION
- PROJECT OBJECTIVE
- DESIGN THINKING
- DATASET DETAILS
- PROGRAM CODE
- CORRELATION CODE
- CONCLUSION

INTRODUCTION:

Mall Customers Segmentation Project is based on customers of mall where we predict which customers has earn more money and spending high by using customer behavior and purchasing data.



Customer segmentation involves implementing data science methods to divide the customer base into smaller groups based on certain characteristics. It assists marketing managers in better understanding their customers' preferences and presenting them with better-targeted advertisements. Malls or shopping complexes are often indulged in the race to increase their customers and hence making huge profits. To achieve this task machine learning is being applied by many stores already. This not only increases sales but also makes the complexes efficient.



The use of machine learning can be seen almost everywhere around us, be it Facebook recognizing you or your friends, or YouTube recommending you a video or two based on your history — Machine Learning is everywhere! However, the ‘magic’ of machine learning is not just limited to only these areas. Machine Learning is broadly categorized as **Supervised** and **Unsupervised Learning**. Supervised Learning is one in which we teach the machine by providing both independent and *dependent variables*, for example, Classifying or predicting values. Unsupervised Learning mainly deals with identifying the structure or pattern of the data. In this type of algorithms, we do not have labeled data(or the dependent variable is absent), for example, clustering data, recommendation systems, etc. Unsupervised Learning provides amazing results as one can deduce many hidden relations between different attributes or features.

In this article, I will be discussing a specific problem based on clustering techniques(Unsupervised Learning). However, my main aim in this article is to discuss the opulent use of machine learning in business and profit enhancement.

The Problem Definition:

Malls or shopping complexes are often indulged in the race to increase their customers and hence making huge profits. To achieve this task machine learning is being applied by many stores already. It is amazing to realize the fact that how machine learning can aid in such ambitions. The shopping complexes make use of their customers' data and develop ML models to target the right ones. This not only increases sales but also makes the complexes efficient.

Here we have the following features :

1. CustomerID:

It is the unique ID given to a customer

2. Gender:

Gender of the customer

3. Age:

The age of the customer

4. Annual Income(k\$):

It is the annual income of the customer

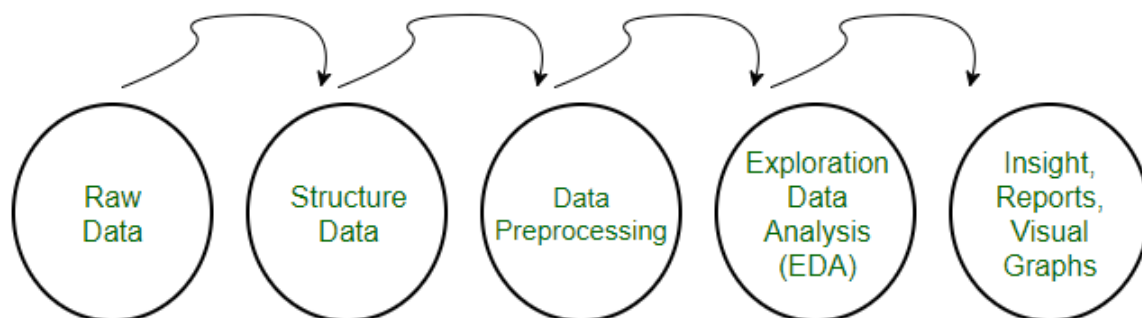
5. Spending Score:

It is the score(out of 100) given to a customer by the mall authorities, based on the money spent and the behavior of the customer.

DESIGN THINKING:

Data Preprocessing

Pre-processing refers to the transformations applied to our data before feeding it to the algorithm. Data preprocessing is a technique that is used to convert the raw data into a clean data set. In other words, whenever the data is gathered from different sources it is collected in raw format which is not feasible for the analysis.



Data Preprocessing

Need of Data Preprocessing:

- For achieving better results from the applied model in Machine Learning projects the format of the data has to be in a proper manner. Some specified Machine Learning model needs information in a specified format, for example, Random Forest algorithm does not support null values, therefore to execute random forest algorithm null values have to be managed from the original raw data set.

- Another aspect is that the data set should be formatted in such a way that more than one Machine Learning and Deep Learning algorithm are executed in one data set, and best out of them is chosen.

DATASET DETAILS:

Dataset **link:**<https://www.kaggle.com/datasets/akram24/mall-customers>

Importing the Dependencies:

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

# loading the data from csv file to a Pandas DataFrame
customer_data = pd.read_csv('/content/Mall_Customers.csv')

# first 5 rows in the dataframe
customer_data.head()
```

	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
```

```
customer_data.shape
(200, 5)
```

```
# getting some informations about the dataset
```

```
customer_data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 200 entries, 0 to 199
```

```
<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   SpendingScore(1-100)  200 non-null   int64
```

```
dtypes: int64(4), object(1)
```

memory usage: 7.9+ KB

checking for missing values

customer_data.isnull().sum()

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

dtype: int64

WCSS -> Within Clusters Sum of Squares

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)

Customer segmentation involves implementing data science methods to divide the customer base into smaller groups based on certain characteristics. It assists marketing managers in better understanding their customers' preferences and presenting them with better-targeted advertisements.

Index	CustomerID	Age	Annual Income (k\$)	Spending Score (1-100)	Gender_Female	Gender_Male
0	1	19	15	39	0	1
1	2	21	15	81	0	1
2	3	20	16	6	1	0
3	4	23	16	77	1	0
4	5	31	17	40	1	0
5	6	22	17	76	1	0
6	7	35	18	6	1	0
7	8	23	18	94	1	0
8	9	64	19	3	0	1
9	10	30	19	72	1	0
10	11	67	19	14	0	1
11	12	35	19	99	1	0
12	13	58	20	15	1	0
13	14	24	20	77	1	0
14	15	37	20	13	0	1
15	16	22	20	79	0	1
16	17	35	21	35	1	0
17	18	20	21	66	0	1
18	19	52	23	29	0	1
19	20	35	23	98	1	0
20	21	35	24	35	0	1

Now the data preprocessing has been done and now let us move on to making the clustering model.

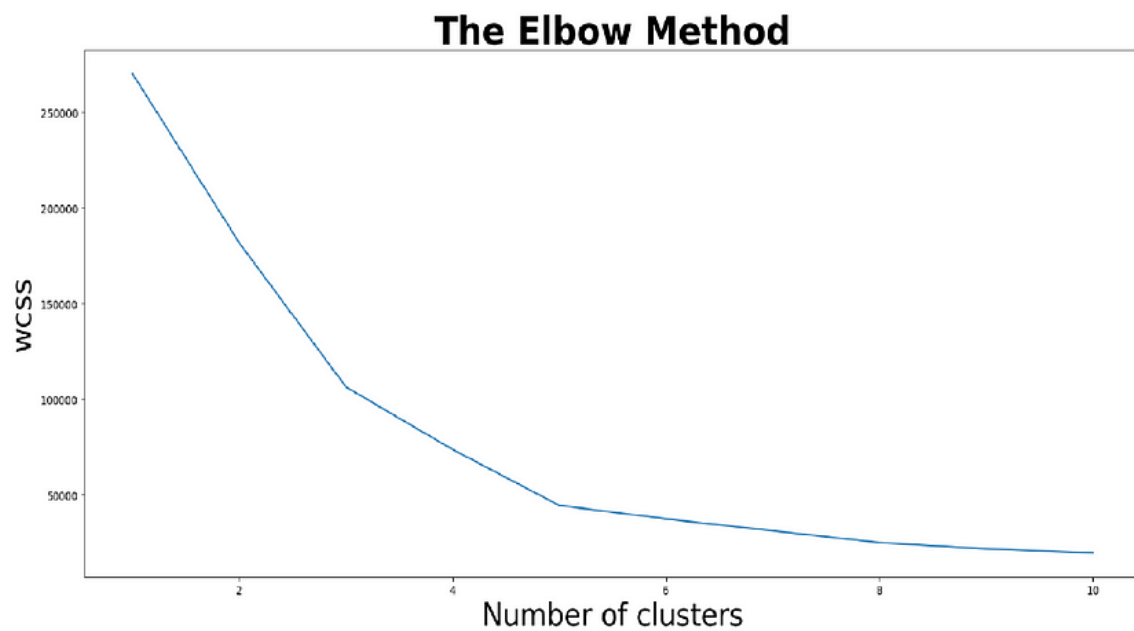
I will use the K-Means Clustering algorithm to cluster the data. To implement K-Means clustering, we need to look at the **Elbow Method**.

The Elbow method is a method of interpretation and validation of consistency within-cluster analysis designed to help to find the appropriate number of clusters in a dataset.

```
# plot an elbow graph

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

The following figure demonstrates the elbow method :



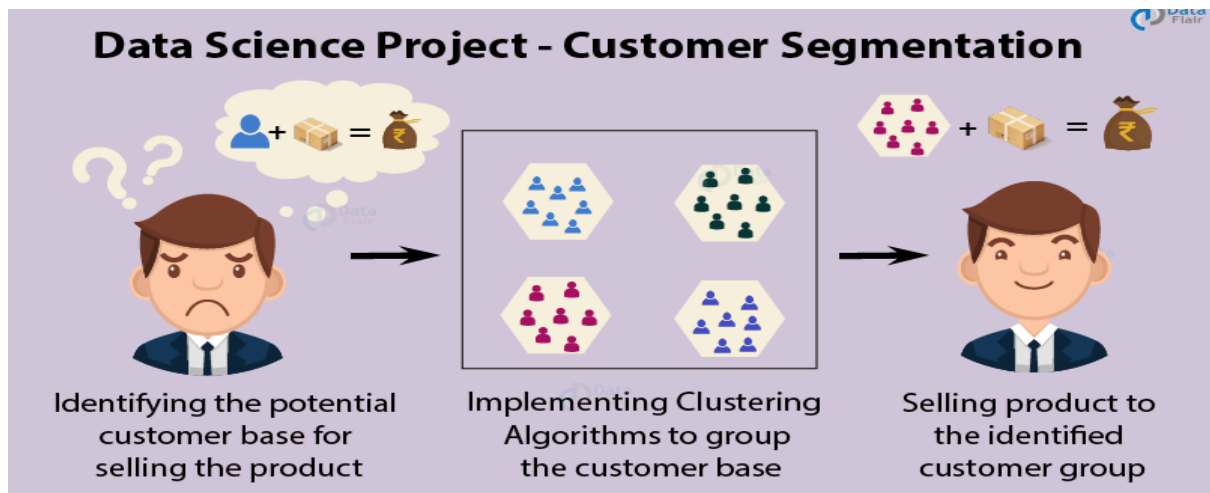
Checking the null values :

Importing the Dependencies:

```
In [2]: df.isnull().sum()
```

```
Out[2]:
```

CustomerID	0
Age	0
Annual Income (k\$)	0
Spending Score (1-100)	0
Gender_Female	0
Gender_Male	0
dtype:	int64



It is clear from the figure that we should take the number of clusters equal to 5, as the slope of the curve is not steep enough after it.

Program:

```
import numpy as np
```

```
import pandas as pd
```

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

```
import seaborn as sns
```

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

```
customer_data = pd.read_csv('/path/to/your/Mall_Customers.csv')
```

```
print(customer_data.head())

print(customer_data.shape)

print(customer_data.info())

print(customer_data.isnull().sum())

X = customer_data.iloc[:, [3, 4]].values

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_)

sns.set()

plt.plot(range(1, 11), wcss)

plt.title('The Elbow Point Graph')

plt.xlabel('Number of Clusters')

plt.ylabel('WCSS')

plt.show()

kmeans = KMeans(n_clusters=5, init='k-means++', random_state=42)
```

```
Y = kmeans.fit_predict(X)

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='red', label='Cluster 2')

plt.scatter(X[Y == 2, 0], X[Y == 2, 1], s=50, c='yellow', label='Cluster 3')

plt.scatter(X[Y == 3, 0], X[Y == 3, 1], s=50, c='violet', label='Cluster 4')

plt.scatter(X[Y == 4, 0], X[Y == 4, 1], s=50, c='blue', label='Cluster 5')

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()
```

OUTPUT:

	CustomerID	Genre	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

(200, 5)

<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	Genre	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

None

CustomerID 0

Genre 0

Age 0

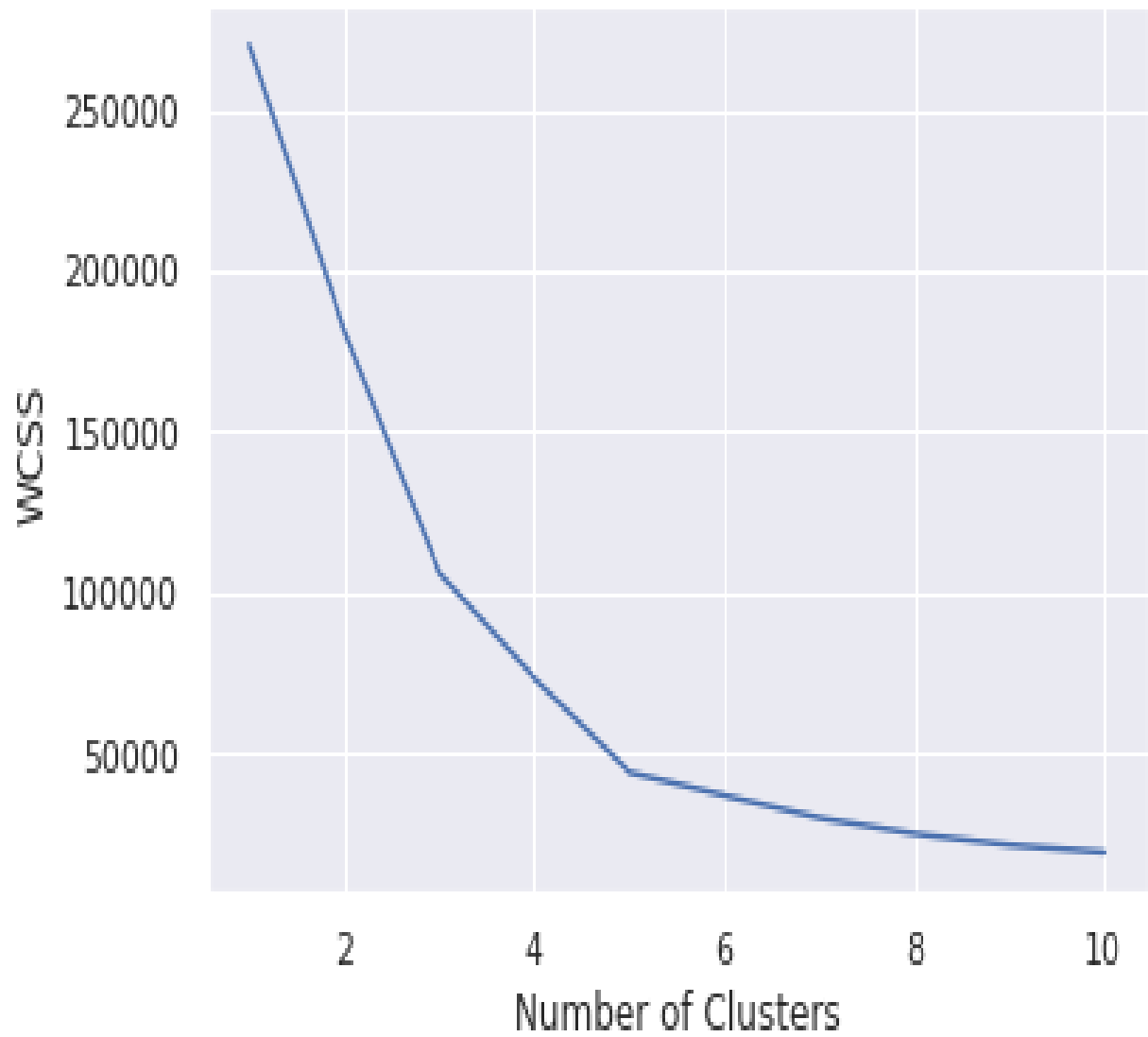
Annual Income (k\$) 0

Spending Score (1-100) 0

dtype: int64

OUTPUT:

The Elbow Point Graph



Optimum Number of Clusters = 5

Training the k-Means Clustering Model

```
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)
```

```
[3 1 3 1 3 1 3 1 3 1 3 1 3 1 3 1 3 1 3 1 3 1 3 1 3 1 3 1 3
 1 3 1 3 1 3 0 3 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 2 4 2 0 2 4 2 4 2 0 2 4 2 4 2 4 2 4 2
4 2 4 2 4 2 4 2 4 2 4 2 4 2 4 2 4 2 4 2 4 2 4 2 4 2 4 2 4 2
2 4 2 4 2 4 2 4 2 4 2 4 2 4 2]
```

5 Clusters - 0, 1, 2, 3, 4

Visualizing all the Clusters

```
# plotting all the clusters and their Centroids

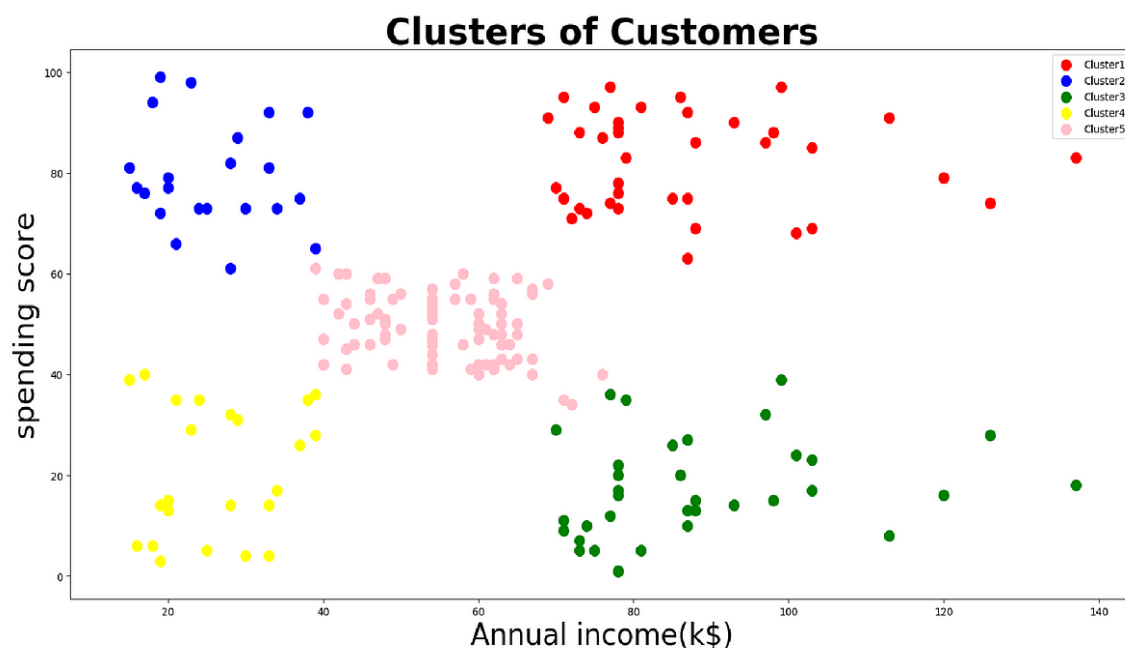
plt.figure(figsize=(8,8))
plt.scatter(X[Y==0,0], X[Y==0,1], s=50, c='green', label='Cluster 1')
```

```
plt.scatter(X[Y==1,0], X[Y==1,1], s=50, c='red', label='Cluster 2')
plt.scatter(X[Y==2,0], X[Y==2,1], s=50, c='yellow', label='Cluster 3')
plt.scatter(X[Y==3,0], X[Y==3,1], s=50, c='violet', 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], kmeans.cluster_centers_[1],
s=100, c='cyan', label='Centroids')

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

Finally, let us plot the clusters :



The data(clusters) are plotted on a spending score Vs annual income curve.

Let us now analyze the results of the model.

Analyzing the Results

We can see that the mall customers can be broadly grouped into 5 groups based on their purchases made in the mall.

In cluster 4(yellow colored) we can see people have low annual income and low spending scores, this is quite reasonable as people having low salaries prefer to buy less, in fact, these are the wise people who know how to spend and save money. The shops/mall will be least interested in people belonging to this cluster.

In cluster 2(blue colored) we can see that people have low income but higher spending scores, these are those people who for some reason love to buy products more often even though they have a low income. Maybe it's because these people are more than satisfied with the mall services. The shops/malls might not target these people that effectively but still will not lose them.

In cluster 5(pink colored) we see that people have average income and an average spending score, these people again will not be the prime targets of the shops or mall, but again they will be considered and other data analysis techniques may be used to increase their spending score.

In cluster 1(red-colored) we see that people have high income and high spending scores, this is the ideal case for the mall or shops as these people are the prime sources of profit. These people might be the regular customers of the mall and are convinced by the mall's facilities.

In cluster 3(green colored) we see that people have high income but low spending scores, this is interesting. Maybe these are the people who are unsatisfied or unhappy by the mall's services. These can be the prime targets of the mall, as they have the potential to spend

money. So, the mall authorities will try to add new facilities so that they can attract these people and can meet their needs.

Finally, based on our machine learning technique we may deduce that to increase the profits of the mall, the mall authorities should target people belonging to cluster 3 and cluster 5 and should also maintain its standards to keep the people belonging to cluster 1 and cluster 2 happy and satisfied.

To conclude, I would like to say that it is amazing to see how machine learning can be used in businesses to enhance profit.

With the help of clustering, we can understand the variables much better, prompting us to take careful decisions. With the identification of customers, companies can release products and services that target customers based on several parameters like income, age, spending patterns, etc. Furthermore, more complex patterns like product reviews are taken into consideration for better segmentation.

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.