

REPORT
ON
**CUSTOMER SEGMENTATION USING K-MEANS
CLUSTERING**

BY
NIVETHAA M
B.E COMPUTER SCIENCE ENGINEERING

AT
EXPOSYS DATA LABS

DURATION
1 WEEK INTERNSHIP

TABLE OF CONTENTS

CHAPTER	TITLE	PAGE NO.
1	Abstract	3
2	Introduction	4
3	Existing Method	5
4	Proposed method	5
4.1	System Architecture	6
5	Algorithm	6
5.1	K-Means clustering	6
5.2	Steps of Algorithm	7
6	Methodology	7
7	Implementation	8
7.1	Overview of Dataset	8
7.2	Exploratory data analysis	8
7.2.1	Information of the dataset	8
7.2.2	Description of the data	9
7.3	Elbow method	10
7.4	Visualization of data	11
8	Conclusion	12

1. ABSTRACT

Effective decisions are mandatory for any company to generate good revenue. In these days competition is huge and all companies are moving forward with their own different strategies. We should use data and take a proper decision. Every person is different from one another and we do not know what he/she buys or what their likes are. But, with the help of machine learning technique one can sort out the data and can find the target group by applying several algorithms to the dataset. Without this, it will be very difficult and no better techniques are available to find the group of people with similar character and interests in a large dataset. Here, the customer segmentation using K-Means clustering helps to group the data with same attributes which exactly helps to business the best. We are going to use elbow method to find the number of clusters and at last we visualize the data.

2. INTRODUCTION

Nowadays the competition is vast and lot of technologies came into account for effective growth and revenue generation. For every business the most important component is data. With the help of grouped or ungrouped data, we can perform some operations to find customer interests.

Data mining helpful to extract data from the database in a human readable format. But, we may not know the actual beneficiaries in the whole dataset. Customer Segmentation is useful to divide the large data from dataset into several groups based on their age, demographics, spent, income, gender, etc. These groups are also known as clusters. By this, we can get to know that, which product got huge number of sales and which age group are purchasing etc. And, we can supply that product much for better revenue generation.

Initially we are going to take the old data. As we know that old is gold so, by using the old data we are going to apply K-means clustering algorithm and we have to find the number of clusters first. So, at lastly, we have to visualize the data. One can easily find the potential group of data while observing that visualization.

The goal of this project is to identify customer segments using the data mining approach,using,the partitioning algorithm called as K-means clustering algorithm.The elbow method determines the optimal clusters.

3. EXISTING METHOD

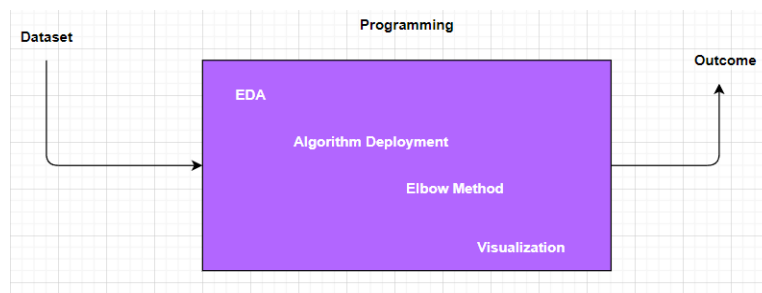
The existing method is storing customer data through paperwork and computer software (digital data) is increasing day by day. At end of the day, they will analyse their data as how many things are sold or actual customer count etc. By analysing the collected data, they got to know who is beneficial to their business and increase their sales. It requires more time and more paperwork. Also, it is not much effective solution to find the desired customers data.

4. PROPOSED METHOD

To overcome the traditional method that is paper work and computerized digital data this new method will play vital role. As we collect a vast data day by day which requires more paperwork and time to do. As new technologies were emerging in today's world. Machine Learning which is powerful innovation which is used to predict the final outcome which has many algorithms. So for our problem statement we will use K-Means Clustering which groups the data into different clusters based on their similar characteristics. And then we will visualize the data.

4.1 SYSTEM ARCHITECTURE

Initially we will see the dataset and then we will perform exploratory data analysis which deals with the missing data, duplicates values and null values. And then we will deploy our algorithm k-means clustering which is unsupervised learning in machine learning.



As in order to find the no of clusters we use elbow method where distance will be calculated through randomly chosen centres and repeat it until there is no change in cluster centres. Thereafter we will analyse the data through data visualization. Finally, we will get the outcome.

5.1.ALGORITHM

5.1.1 K-Means Clustering

- ⦿ K Means algorithm in an iterative algorithm that tries to partition the dataset into K predefined distinct non overlapping sub groups which are called as cluster.
- ⦿ Here K is the total no of clusters.
- ⦿ Every point belongs to only one cluster.
- ⦿ Clusters cannot overlap.

5.1.2 Steps of Algorithm

- ⦿ Arbitrarily choose k objects from D as the initial cluster centers.
- ⦿ Repeat.
- ⦿ 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.
- ⦿ Update the cluster means, i.e. calculate the mean value of the objects for each cluster.
- ⦿ Until no change.

6. METHODOLOGY

1. First of all, we will import all the necessary libraries or modules (pandas, NumPy, seaborn).
2. Then we will read dataset and analyse whether it contains any null values, missing values, and duplicate values. So, we will fix them by dropping or fixing the value with their means, medians etc which is technically named as Data Pre-processing.
3. We will deploy our model algorithm K-Means Clustering, which divides the data into group of clusters based on similar characteristics. To find no.of clusters we will use elbow method.
4. Finally, we will visualize our data using mat plot, which concludes the customers divided into groups who are similar to each other on their group.

7. IMPLEMENTATION

7.1 Overview of a Dataset

This is a mall customer segmentation data which contains 5 columns and 200 rows.

	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
...
195	196	Female	35	120	79
196	197	Female	45	126	28
197	198	Male	32	126	74
198	199	Male	32	137	18
199	200	Male	30	137	83

200 rows × 5 columns

7.2 Exploratory Data Analysis

It deals with the data pre-processing, whether it contains any missing values or null values. There after we will see the information and description of the dataset.

7.2.1 Information of the dataset

`#df.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
```

As here it overview the information of the data. And it gives it doesn't contain any null values.

As we will remove the irrelevant data which is customer id.

```
df.drop(["CustomerID"], axis=1, inplace=True)
```



```
# so here customer data is not required to our analysis. We will drop it.
|
df.drop(["CustomerID"], axis=1, inplace=True)

# printing data frame again (Now, CustomerID column is removed)
df
```

	Gender	Age	Annual Income (k\$)	Spending Score (1-100)
0	Male	19	15	39
1	Male	21	15	81
2	Female	20	16	6
3	Female	23	16	77
4	Female	31	17	40
5	Female	22	17	76
6	Female	35	18	6
7	Female	23	18	94

7.2.2 Description of the data

#df.describe()

	Age	Annual Income (k\$)	Spending Score (1-100)
count	200.000000	200.000000	200.000000
mean	38.850000	60.560000	50.200000
std	13.969007	26.264721	25.823522
min	18.000000	15.000000	1.000000
25%	28.750000	41.500000	34.750000
50%	36.000000	61.500000	50.000000
75%	49.000000	78.000000	73.000000
max	70.000000	137.000000	99.000000

It describes about the count which counts the no of rows in it, mean of the columns, standard deviations, maximum and minimum and percentiles etc.

7.3 Elbow Method

The elbow method is based on the observation that increasing the number of clusters can help to reduce the sum of within-cluster variance of each cluster. This is because having more clusters allows one to capture finer groups of data objects that are more similar to each other. To define the optimal clusters, Firstly, we use the clustering algorithm for various values of k. This is done by ranging k from 1 to 10 clusters. Then we calculate the total intra-cluster sum of square. Then, we proceed to plot intra-cluster sum of square based on the number of clusters. The plot denotes the approximate number of clusters required in our model. The optimum clusters can be found from the graph where there is a bend in the graph. First we will consider the data X which has only two columns they are annual income and spending score.

```
X=df[['Annual Income (k$)','Spending Score (1-100)']]
```

```
X.head()
```

Code:

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

```
wcss = []
```

```
for i in range(1,11):
```

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

```
    kmeans.fit(x)
```

```
    wcss.append(kmeans.inertia_)
```

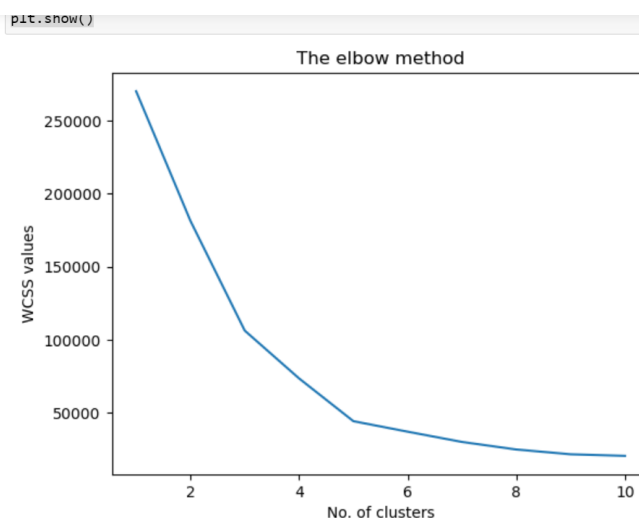
```
plt.plot(range(1,11),wcss)
```

```
plt.title('The elbow method')
```

```
plt.xlabel('No. of clusters')
```

```
plt.ylabel('WCSS values')
```

```
plt.show()
```



7.4 Visualization the clusters

Visualizing the clusters based on Annual Income and Spending Score of the customers. As here we plot a graph named as Clusters of Customers to visualize the data in terms of groups or cluster.

```
plt.scatter(x[y_kmeans == 0,0], x[y_kmeans == 0,1], s=80, c="red", label='Customer 1')
plt.scatter(x[y_kmeans == 1,0], x[y_kmeans == 1,1], s=80, c="blue", label='Customer 2')
plt.scatter(x[y_kmeans == 2,0], x[y_kmeans == 2,1], s=80, c="black", label='Customer 3')
plt.scatter(x[y_kmeans == 3,0], x[y_kmeans == 3,1], s=80, c="green", label='Customer 4')
plt.scatter(x[y_kmeans == 4,0], x[y_kmeans == 4,1], s=80, c="yellow", label='Customer 5')

plt.scatter(kmeans.cluster_centers_[0], kmeans.cluster_centers_[1], s=100, c='magenta',
label='Centroids')

plt.title('Clusters of customers')
plt.xlabel('Annual Income (k$)')
plt.ylabel('Spending Score (1-100)')
plt.legend()
plt.show()
```



8. Conclusion

- ⦿ The Highest income , high spending can be targeting these type of customers as they earn more money and spend as much as they want.
- ⦿ Highest income, low spending can be targeting these types of customers by asking feedback and advertising the product in a better way.
- ⦿ Average income, Average spending may or may not be beneficial to the mall owners of this type of customers.
- ⦿ Low income, High spending can be targeting these types of customers by providing them with low-cost EMI's etc.
- ⦿ Low income, Low spending don't target these types of customers because they earn a bit and spend some amount of money.

So high income, high spending are the most beneficial ones to the mall owners which increases the owner's business. (Cluster 1)