

Customer Segmentation using K –Means Clustering

Keywords: Clustering, Customer Segmentation, K –Means Algorithm, Elbow method.

1. Abstract

We live in a world where large and vast amount of data is collected daily. Analyzing such data is an important need. In the modern era of innovation, where there is a large competition to be better than everyone, the business strategy needs to be according to the modern conditions. 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 paper, the clustering algorithm used 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.

2. Introduction

Over the years, the competition amongst businesses is increased and the large historical data that is available has resulted

in the widespread use of data mining techniques in extracting the meaningful and strategic information from the database of the organization. Data mining is the process where methods are applied to extract data patterns in order to present it in the human readable format which can be used for the purpose of decision support. According to, Clustering techniques consider data tuples as objects. They partition the data objects into groups or clusters, so that objects within a cluster are similar to one another and dissimilar to objects in other clusters. Customer Segmentation is the process of division of customer base into several groups called as customer segments such that each customer segment consists of customers who have similar characteristics. The segmentation is based on the similarity in different ways that are relevant to marketing such as gender, age, interests, and miscellaneous spending habits. The customer segmentation has the importance as it includes, the ability to modify the programs of market so that it is suitable to each of the customer segment, support in business decision; identification of products associated with each customer segment and to manage the demand and supply of that product; identifying and targeting the potential customer base, and predicting customer defection, providing directions in finding the solutions. The

thrust of this paper 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 cluster

3. Existing Methods

There are actually variety of methods used for customer segmentation. The most popular methods used for customer segmentation are K - Means clustering, Agglomerative Hierarchical Clustering, Expectation – Maximization Clustering, Density-Based Spatial clustering, and Mean – Shift Clustering.

Agglomerative Hierarchical Clustering (ACM):

Agglomerative Hierarchical Clustering is a clustering (or classification) method which works from the dissimilarities between the objects to be grouped together. A type of dissimilarity can be suited to the subject studied and the nature of the data. One of the results is the dendrogram which shows the progressive grouping of the data. It is then possible to gain an idea of a suitable number of classes into which the data can be grouped.

Expectation – Maximization (EM) clustering:

Expectation-Maximization is an iterative method which alternates between two steps, expectation (E) and maximization (M). For clustering, EM makes use of the finite Gaussian mixtures model and estimates a set of parameters iteratively until a desired convergence value is achieved.

Density – Based spatial clustering:

Density - Based clustering refers to unsupervised machine learning methods that identify distinctive clusters in the data, based on the idea that a cluster/group in a data space is a contiguous region of high point density, separated from other clusters by sparse regions.

Mean - Shift Clustering:

Mean shift clustering algorithm is a centroid-based algorithm that helps in various use cases of unsupervised learning. It is one of the best algorithms to be used in image processing and computer vision. It works by shifting data points towards centroids to be the mean of other points in the region.

4. Proposed Method

Clustering algorithms generates 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 in one of the most popular centroid based algorithm. 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.

5. Methodology

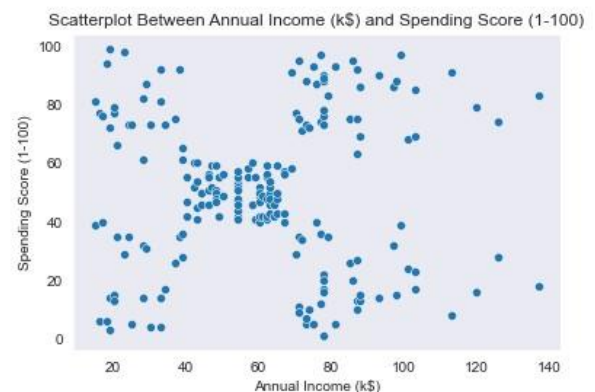
The data set used to implement clustering and K-means algorithm was collected from a store of shopping mall. The data set contains 5 attributes and has 200 tuples, representing the data of 200 customers. The attributes in the data set has CustomerId, gender, age, annual income (k\$), spending score on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

Importing the necessary libraries:

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

Scatterplot between Annual Income (k\$) and Spending Score (1-100):

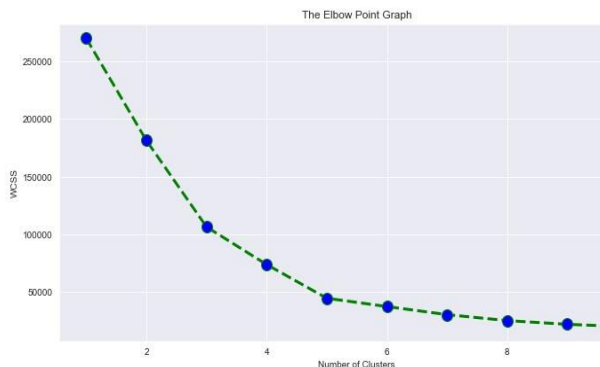
```
sns.set_style('dark')
sns.scatterplot(x = 'Annual Income (k$)', y = 'Spending Score (1-100)', data = df)
plt.xlabel('Annual Income (k$)')
plt.ylabel('Spending Score (1-100)')
plt.title('Scatterplot Between Annual Income (k$) and Spending Score (1-100)')
```



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.

```
plt.figure(figsize = (12,6))
plt.grid()
plt.plot(range(1,11),wcss, color='green', linestyle='dashed', linewidth = 3,
         marker='o', markerfacecolor='blue', markersize=12)
plt.title('The Elbow Point Graph')
plt.xlabel('Number of Clusters')
plt.ylabel('WCSS')
plt.show
```



Visualize the clusters

```
plt.figure(figsize=(8,8))
plt.scatter(X[label == 0,1], X[label == 0,1], s=50, c='green', label='Cluster 1')
plt.scatter(X[label == 1,1], X[label == 1,1], s=50, c='yellow', label='Cluster 2')
plt.scatter(X[label == 2,1], X[label == 2,1], s=50, c='red', label='Cluster 3')
plt.scatter(X[label == 3,1], X[label == 3,1], s=50, c='purple', label='Cluster 4')
plt.scatter(X[label == 4,1], X[label == 4,1], s=50, c='blue', label='Cluster 5')
plt.scatter(kmeans.cluster_centers_[:,0], kmeans.cluster_centers_[:,1], s=100, c='black', marker='*', label='Centroids') #Plotting the centroids
plt.title('Customer groups')
plt.xlabel('Annual Income')
plt.ylabel('Spending Score (1-100)')
plt.legend()
plt.show()
```



Business Insights

The result of the analysis shows that the retail store customers can be group into 5 clusters or segments for targeted marketing.

Cluster 1 (green): These are average income earners with average spending scores. They are cautious with their spending at the store.

Cluster 2 (yellow): The customers in this group are high income earners and with high spending scores. They bring in profit. Discounts and other offers targeted at this group will increase their spending score and maximize profit.

Cluster 3 (red): This group of customers have a higher income but they do not spend more at the store. One of the assumption could be that they are not satisfied with the services rendered at the store. They are another ideal group to be targeted by the marketing team because they have the potential to bring in increased profit for the store.

Cluster 4 (purple): Low income earners with low spending score. I can assume that this is so because people with low income will tend to purchase less item at the store.

Cluster 5 (blue): These are low income earning customers with high spending scores. I can assume that why this group of customers spend more at the retail store despite earning less is because they enjoy and are satisfied with the services rendered at the retail store.

6. Conclusion

This study demonstrates that client segmentation in shopping malls is achievable despite the fact that this form of machine learning application is highly useful in the market, a manager can concentrate all of his or her attention on each cluster that has been discovered and meet all of their requirements. Mall

managers must be able to understand what customers require and, more importantly, how to meet those needs, analyze their purchasing habits, and establish frequent encounters with customers that make them feel comfortable in order to satisfy their demands.