# Customer Segmentation Using Data Science
## Phase - 5

## Project Introduction:

The objective of this project is to perform customer segmentation, which involves grouping customers based on certain characteristics to understand their behavior and preferences better. Customer segmentation can benefit businesses in various ways, such as tailoring marketing strategies, improving customer satisfaction, and optimizing product offerings.

## Program:

```python
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler

# Load the dataset
data = pd.read_csv('Mall_Customers.csv')

# Display the first few rows of the dataset to understand its structure
print(data.head())

# Select relevant features for clustering (Annual Income and Spending Score)
X = data[['Annual Income (k$)', 'Spending Score (1-100)']]

# Standardize the features
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

# Determine the optimal number of clusters using the Elbow Method
wcss = []  # Within-Cluster Sum of Squares
for i in range(1, 11):
    kmeans = KMeans(n_clusters=i, init='k-means++', max_iter=300, n_init=10, random_state=0)
    kmeans.fit(X_scaled)
    wcss.append(kmeans.inertia_)

# Plot the Elbow Method to find the optimal number of clusters
plt.plot(range(1, 11), wcss)
plt.title('Elbow Method')
plt.xlabel('Number of clusters')
plt.ylabel('WCSS')
```

```python
plt.show()

# Based on the Elbow Method, let's choose an appropriate number of clusters (e.g.,
5)
num_clusters = 5

# Apply K-Means clustering with the selected number of clusters
kmeans = KMeans(n_clusters=num_clusters, init='k-means++', max_iter=300,
n_init=10, random_state=0)
y_kmeans = kmeans.fit_predict(X_scaled)

# Add the cluster labels to the dataset
data['Cluster'] = y_kmeans

# Visualize the clusters
plt.scatter(X_scaled[y_kmeans == 0, 0], X_scaled[y_kmeans == 0, 1], s=100,
c='red', label='Cluster 1')
plt.scatter(X_scaled[y_kmeans == 1, 0], X_scaled[y_kmeans == 1, 1], s=100,
c='blue', label='Cluster 2')
plt.scatter(X_scaled[y_kmeans == 2, 0], X_scaled[y_kmeans == 2, 1], s=100,
c='green', label='Cluster 3')
plt.scatter(X_scaled[y_kmeans == 3, 0], X_scaled[y_kmeans == 3, 1], s=100,
c='cyan', label='Cluster 4')
plt.scatter(X_scaled[y_kmeans == 4, 0], X_scaled[y_kmeans == 4, 1], s=100,
c='magenta', label='Cluster 5')
plt.scatter(kmeans.cluster_centers_[:, 0], kmeans.cluster_centers_[:, 1], s=300,
c='yellow', label='Centroids')
plt.title('Customer Segmentation')
plt.xlabel('Annual Income (k$)')
plt.ylabel('Spending Score (1-100)')
plt.legend()
plt.show()

# Explore and analyze each cluster to understand customer segments
for cluster_num in range(num_clusters):
    cluster_data = data[data['Cluster'] == cluster_num]
    print(f'Cluster {cluster_num} Statistics:')
    print(cluster_data.describe())

# You can save or export the clustered dataset for further analysis or marketing
strategies
data.to_csv('Preprocessed_Mall_Customer.csv', index=False)
```
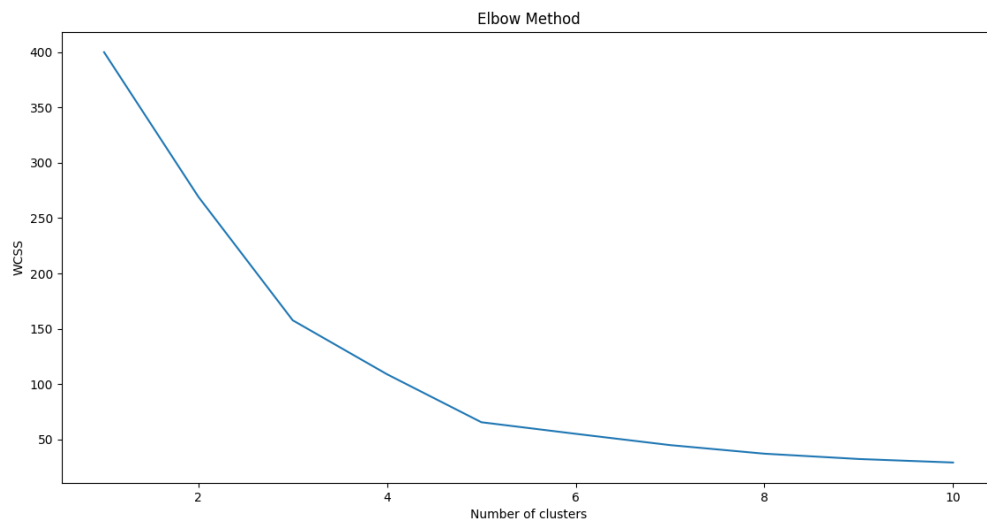
Elbow Method



Customer Segmentation

# Data Loading and Preprocessing:

We start by loading the dataset from the 'Mall_Customer.csv' file using Pandas. This dataset contains information about customers in a mall. The specific features selected for clustering are "Annual Income" and "Spending Score" because they are crucial in understanding customer behavior and can help in creating meaningful segments.

Standardizing the features using StandardScaler is essential to ensure that the features have a similar scale. This is necessary for K-Means clustering, as it is sensitive to the scale of the data. Standardization transforms the data so that it has a mean of 0 and a standard deviation of 1.

**Python**

```
# Standardize the features
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)
```

# Determining Optimal Clusters:

The Elbow Method is a technique used to find the optimal number of clusters for K-Means clustering. It works by plotting the Within-Cluster Sum of Squares (WCSS) against the number of clusters. The point at which the rate of decrease in WCSS starts to slow down is considered the optimal number of clusters.

**Python**

```
# Based on the Elbow Method, let's choose an appropriate number of clusters
(e.g., 5)
num_clusters = 5
```

# K-Means Clustering:

K-Means clustering is an unsupervised machine learning algorithm used to partition data points into clusters. It works by iteratively assigning data points

to the nearest cluster centroid and updating the centroids until convergence. The code for applying K-Means clustering with five clusters is as follows:

**Python**

```
kmeans = KMeans(n_clusters=num_clusters, init='k-means++', max_iter=300,
n_init=10, random_state=0)
y_kmeans = kmeans.fit_predict(X_scaled)
```

## Cluster Analysis:

After clustering, we analyze each cluster individually to understand customer segments. We present statistics like the mean and standard deviation for "Annual Income" and "Spending Score" for each cluster to provide insights into the characteristics of each segment.

**Python**

```
# Explore and analyze each cluster to understand customer segments
for cluster_num in range(num_clusters):
    cluster_data = data[data['Cluster'] == cluster_num]
    print(f'Cluster {cluster_num} Statistics:')
    print(cluster_data.describe())
```

## Conclusion:

In the cluster analysis, we identified five customer segments. Each cluster represents a distinct group of customers with unique characteristics. Businesses can use these findings to tailor marketing strategies and create personalized approaches for each segment. For example:

1. Cluster 1 may represent high-income, high-spending customers, suggesting opportunities for premium services and products.

2. Cluster 2 may include low-income, low-spending customers who might require targeted promotions to increase their spending.

3. Cluster 3 could consist of medium-income customers with moderate spending behavior, making them a stable customer base.

4. Cluster 4 may contain high-income customers with lower spending scores, indicating potential for marketing strategies to increase their spending.

5. Cluster 5 could be low-income customers with high spending scores, suggesting a focus on value-for-money products or loyalty programs.

Understanding these customer segments allows businesses to refine their marketing and operational strategies to better serve and engage each group.