
Project Title: Customer Segmentation Using K-Means Clustering

Author: Yogapriya P

Internship: Infotact Solutions – Data Analytics Internship

Duration: 2 Months

OBJECTIVE :

This project focuses on customer segmentation using K-Means clustering. The goal is to categorize customers into distinct groups based on their behaviour or purchasing patterns to help businesses personalize marketing strategies.

- Identifying customer segments from raw data
- Use clustering to find similar behaviour groups
- Providing actionable insights for business decisions

TOOLS AND TECHNOLOGIES :

- Programming Language: Python
- Libraries: Pandas, Matplotlib, Seaborn, Sklearn, Streamlit, Joblib
- Jupyter Notebook and VS code
- Visualization: matplotlib & seaborn

DATA UNDERSTANDING AND DATA PREPROCESSING

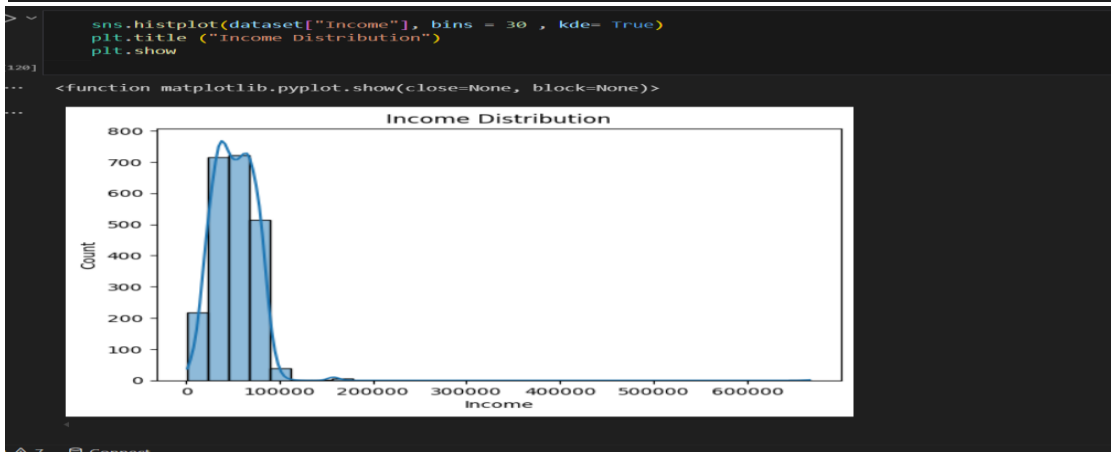
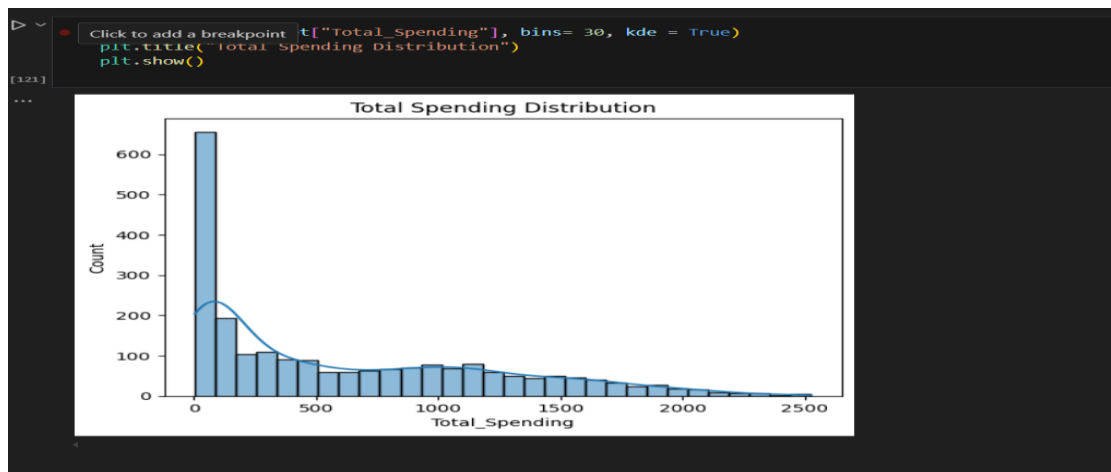
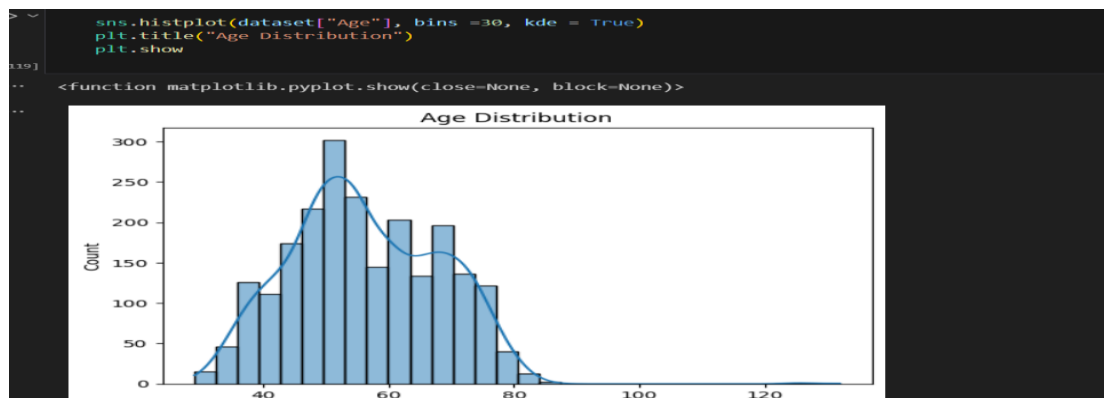
- Dataset : <https://www.kaggle.com/datasets/vishakhapat/customer-segmentation-clustering>
- Initial exploration showed missing values and outliers, which were handled during preprocessing.
- Created new columns for Age, Total Children, Total Spending, Customer Since.

EXPLORATORY DATA ANALYSIS (EDA)

- Histograms and boxplots were used to understand distributions
- Correlation matrix showed strong relationships between some features. Visualized clusters using 2D scatter plots.

1) Histogram and Boxplots:

A **histogram** shows the **distribution of a single variable** — i.e., how frequently each range of values occurs. Here we used histogram for **Age Distribution, Income Distribution, Total Spending Distribution**.

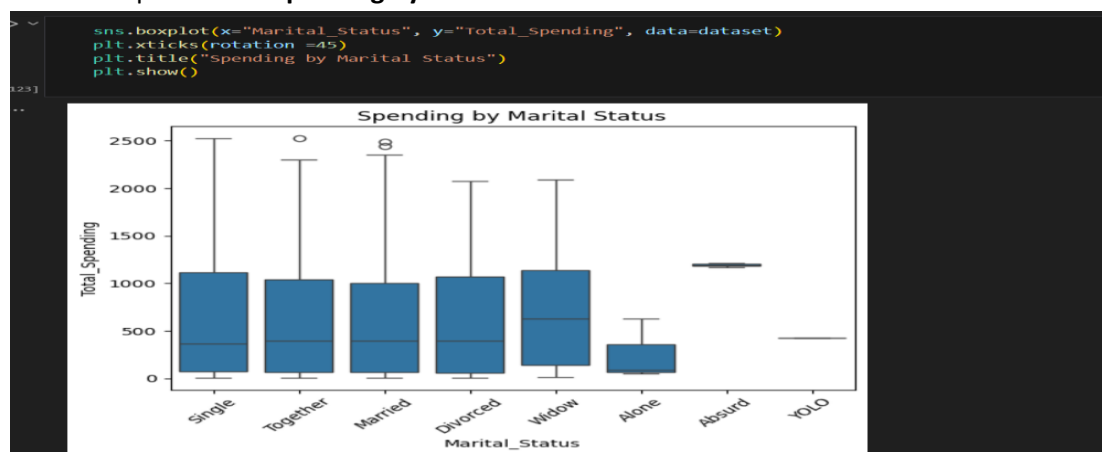


A **box plot** shows the **summary statistics** of a variable. Here it uses to detect **outliers** Compare **distribution between categories** like

- relationship between **Education and Income**.

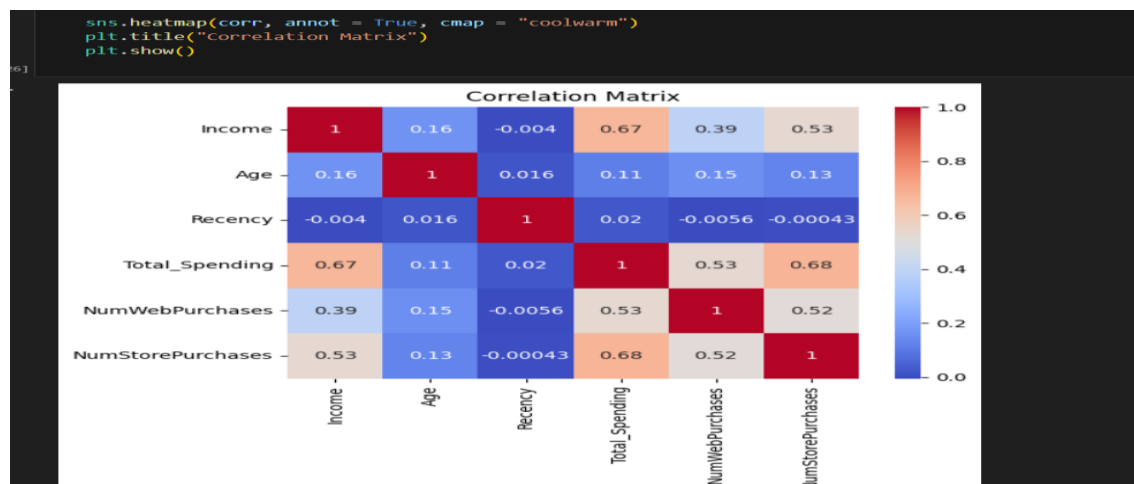


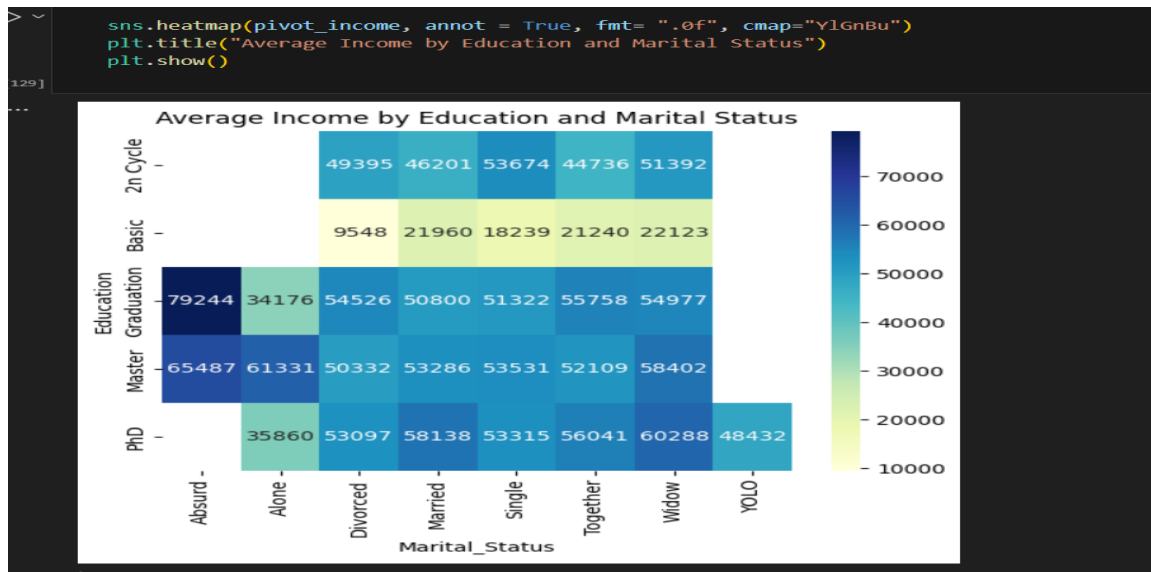
- relationship between Spending by Marital Status and Marital Status.



2) Correlational Matrix :

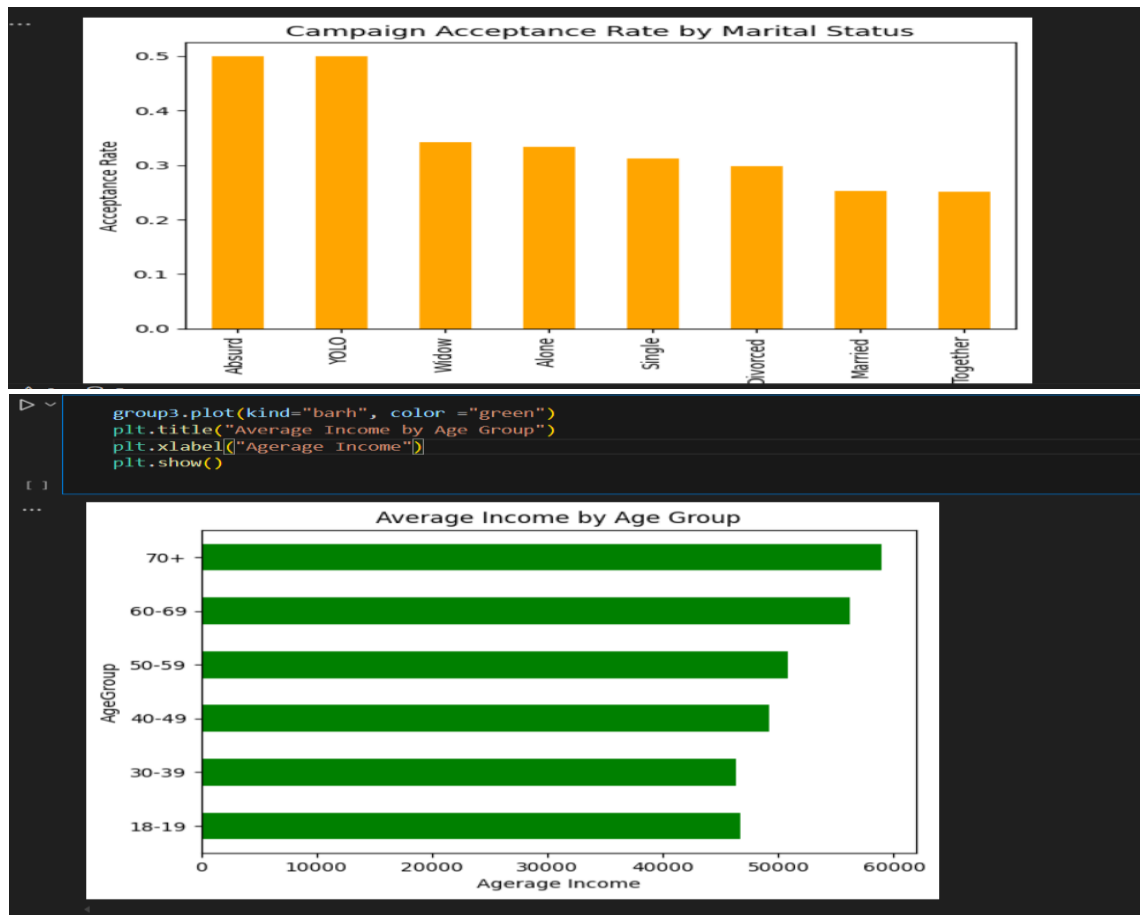
The correlation coefficients between multiple variables in a dataset. Each cell in the matrix shows the relationship strength between two variables. Income, Age, Recency, Total Spending, No. of Purchases in Web and Store.





BUSINESS ANALYSIS :

Forming several groups and their specified plot. Plots for Average Spending by Education and Campaign Acceptance Rate by Marital Status gives customer approach towards their needs. Average Income by age Group provide essential informations.

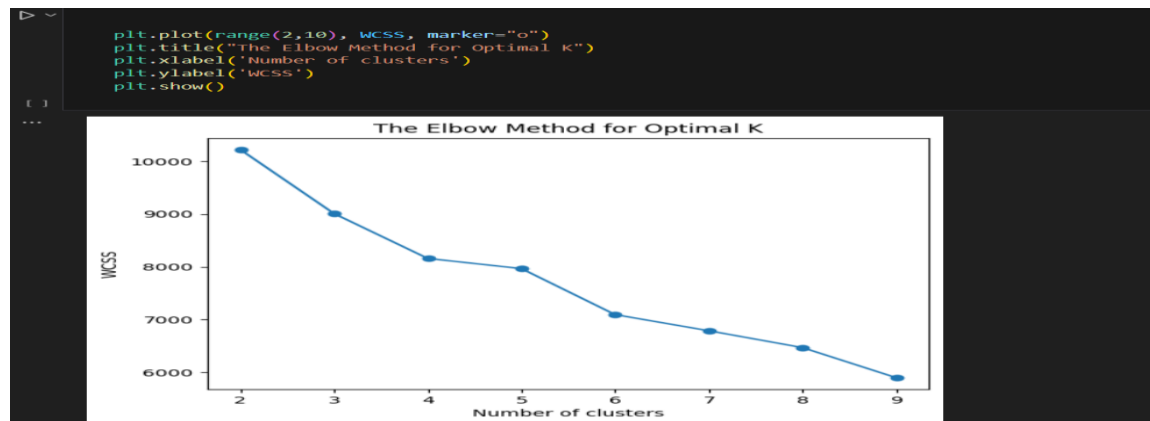


MODELLING K-MEANS CLUSTERING :

Applied StandardScaler for normalization and Selected relevant features for clustering. Finally our features are,

```
features = ["Age", "Income", "Total_Spending", "NumWebPurchases", "NumStorePurchases",  
"NumWebVisitsMonth", "Recency"]
```

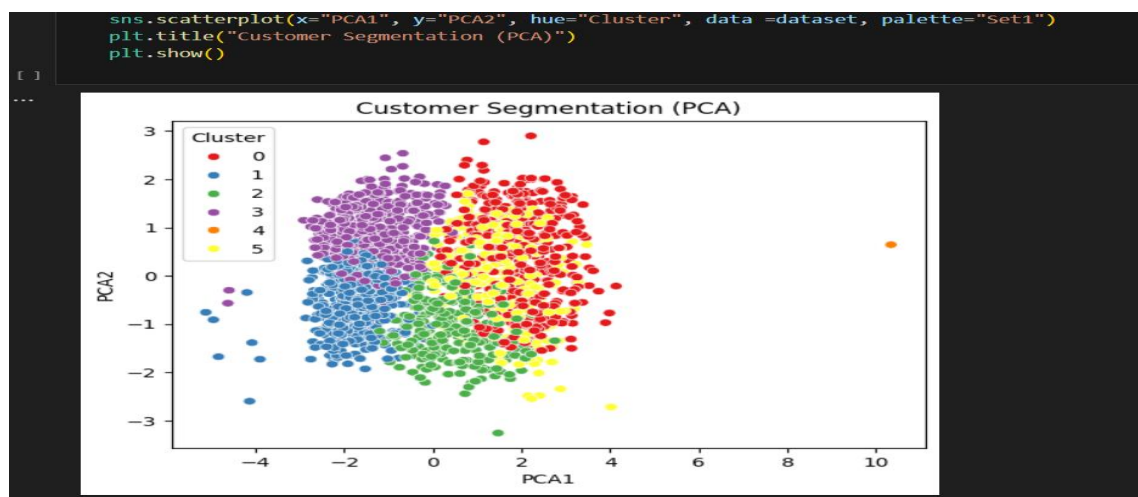
- Used KMeans from sklearn.cluster.
- Chose optimal number of clusters using **WCSS** and **Elbow Method**.
- Applied clustering and visualized the segments.



```
cluster_summary
```

	Age	Income	Total_Spending	NumWebPurchases	NumStorePurchases	NumWebVisitsMonth	Recency
Cluster							
0	55.988395	37082.862669	135.058027	2.446809	3.491296	6.369439	75.835590
1	46.122257	79012.852665	1320.576803	4.564263	8.551724	2.598746	48.952978
2	69.945392	74080.593857	1182.880546	4.436860	8.337884	2.443686	54.180887
3	58.910653	61402.103093	1007.597938	8.333333	8.491409	6.292096	62.384880
4	50.675676	31671.886100	81.113900	1.870656	2.980695	6.646718	27.264479
5	60.964029	55479.517986	620.064748	5.892086	6.705036	6.025180	20.276978

Results and Insights :



Cluster 0: Budget Customers - This cluster has the lowest income & spending. They are also older and have not made a purchase recently.

Cluster 1: Premium Customers - This is the highest-earning cluster, and they are also the highest spenders. They are younger and have made a purchase recently.

Cluster 2: Senior Spenders - This cluster is the oldest and has the second-highest income and spending.

Cluster 3: Digital Buyers - This group has a high number of web purchases and a low number of in-store purchases. They have a high income and are of average age.

Cluster 4: Inactive Customers - This cluster has the lowest recency, meaning they have not made a purchase in a long time. They also have low income and spending.

Cluster 5: Frequent Buyers - This cluster has the highest number of purchases, both online and in-store. They have an average income and are of average age.

FINAL DEPLOYMENT IN STREAMLIT:

Customer Segmentation App

Enter customer details to predict the segments.

Age
25

Income
50000

Total Spending (sum of purchase)
1000

Number of Web Purchases
100

Number of store Purchase
50

Number of visits per Month
5

Recency (days since last purchase)
16

[Predict Segment](#)

Predicted Segment : Cluster 5

Cluster 0: Budget Customers - This cluster has the lowest income & spending. They are also older and have not made a purchase recently.

Cluster 1: Premium Customers - This is the highest-earning cluster, and they are also the highest spenders.

CONCLUSION :

K-Means clustering helped uncover patterns in customer data. The segments can improve decision-making for marketing, product recommendations, and customer engagement.

Links

GitHub: <https://github.com/YogaPriya2000/Customer-Segmentation-KMeans-Clustering->

Linkedin: <https://www.linkedin.com/feed/update/urn:li:activity:7355552242793373696/>

YouTube: <https://lnkd.in/gSih6NUF>