

Exploratory Data Analysis (EDA) Summary for Customer Segmentation

1. Data Overview:

- The analysis involves two datasets: `Customers.csv` containing customer profile information (such as age, income, spending score, etc.) and `Transactions.csv` with transaction data (including transaction amount, quantity, and total spent).
- The first step is to load and explore these datasets to understand their structure, identify missing values, and check the data types.

2. Data Cleaning:

- **Missing Values:** Handle any missing values in the datasets by filling them with appropriate statistical measures (mean for numerical features) or dropping unnecessary columns.
- **Data Types:** Ensure that columns are in the correct data type (e.g., date columns as `datetime`, numerical features as integers or floats).

3. Feature Engineering:

- Create new features like `TotalSpent` by multiplying transaction amount and quantity, and `Tenure` to calculate how long each customer has been active since signup.
- Merge the customer and transaction datasets on `CustomerID` to form a comprehensive view of customer profiles along with their transaction history.

4. Univariate Analysis:

- **Age Distribution:** Plot a histogram to understand the age distribution of customers. It could reveal whether the dataset has more younger or older customers.
- **Income Distribution:** Similarly, the income distribution plot could indicate if the customer base is skewed towards higher or lower income groups.
- **Spending Score:** A boxplot or histogram helps identify any outliers or clusters of customers based on their spending behavior.

5. Bivariate Analysis:

- **Income vs. Spending Score:** Scatter plots can help visualize if there is any relationship between a customer's income and their spending score. This could provide insights into whether high-income customers tend to spend more or not.
- **Age vs. Spending Score:** A scatter plot or correlation matrix can show the relationship between customer age and spending, potentially revealing patterns like whether younger customers tend to spend more or less.
- **Correlation Heatmap:** This is crucial for understanding relationships between multiple numerical features. For instance, how spending score correlates with income, total spent, and customer tenure.

6. Multivariate Analysis:

- **PCA for Dimensionality Reduction:** A pair plot or PCA (Principal Component Analysis) can be used to visualize clusters of customers based on several features like age, income, and spending score in a 2D space.
- The results of PCA would allow better visualization of how different customer segments group based on combined features.

7. Outlier Detection:

- Identifying outliers in features like age, income, and spending score helps in understanding if there are any extreme values in the dataset that may need special handling.
- Using methods like the Z-score, extreme values are detected and can either be handled by removal or adjustment depending on their impact.

8. Segment Insights:

- **Cluster Characteristics:** After applying clustering techniques like K-Means, you can gain insights into different customer segments. For instance, clusters may emerge with common characteristics such as:
 - **Cluster 1:** Younger, low-income customers who do not spend much.
 - **Cluster 2:** Older, high-income customers who spend more regularly.
- These insights allow businesses to tailor marketing strategies, such as targeted promotions for high-spending customers or budget-friendly offers for younger, lower-income segments.

9. Conclusion:

- EDA helps in understanding the dataset by visualizing distributions, correlations, and potential outliers. It serves as a critical foundation before applying clustering algorithms for segmentation.
- The resulting customer segments can be leveraged to inform business decisions, such as personalized marketing, customer retention strategies, or product recommendations.