# SEGMENTATION REPORT

## ● Data Preprocessing:

- Merges customers.csv and transactions.csv datasets on CustomerID.
- Selects and preprocesses key features like Quantity, TotalValue, Price, Region, and SignupDate.
- Encodes categorical columns (Region) and converts SignupDate into a numeric format.
- Standardizes the features using StandardScaler.

## ● Optimal Number of Clusters:

- Evaluates the optimal number of clusters based on metrics like Davies-Bouldin index, Silhouette score, and Calinski-Harabasz score.

## ● Clustering Algorithms:

- Applies multiple clustering algorithms: KMeans, DBSCAN, Agglomerative, GaussianMixture (GMM), and SpectralClustering.
- Stores the clustering results for each method.

## ● Cluster Summary:

- Generates summary statistics (mean and standard deviation) for Quantity, TotalValue, and Price per cluster.

## ● Cluster Visualization:

- Creates various visualizations including:
- Distribution of Quantity, TotalValue, and Price for each cluster (box plots).
- Region distribution across clusters (bar plot).
- 3D visualization of clusters using PCA.

## ● Cluster Evaluation:

- Evaluates the performance of clustering methods using multiple metrics like Davies-Bouldin, Silhouette, and Calinski-Harabasz.

## Recommendations for Optimization and Next Steps:

## Handling Outliers:

- Consider adding outlier detection/removal before clustering, as outliers can affect clustering performance, especially in KMeans and GMM.

## DBSCAN Parameter Tuning:

◼ For DBSCAN, the eps (neighborhood size) and min_samples parameters significantly impact the results. It may require fine-tuning based on your dataset's characteristics.

## Clustering with Dimensionality Reduction:

◼ You could use PCA or t-SNE for dimensionality reduction before applying clustering algorithms to improve clustering performance, especially when dealing with high-dimensional data.

## Model Evaluation Metrics:

◼ Consider evaluating clustering results using a few more metrics like Adjusted Rand Index (ARI) or Normalized Mutual Information (NMI) if you have ground truth labels available for comparison.

## Visualizations:

◼ For the 3D visualizations, ensure that the colors correspond to meaningful labels to make the visualizations more informative. You could also include hover information like the customer ID or region for better interaction.

## Code Refactoring:

◼ Avoid redundancy by combining similar code snippets (e.g., plotting functions). You can generalize plotting functions that handle box plots, bar plots, and 3D plots into reusable helper functions.

## Potential Next Steps:

**Explore Clusters:** After performing clustering, you can further analyze the characteristics of each cluster and tailor marketing strategies or customer service approaches based on the cluster profiles.

**Fine-Tune the Models:** Based on the evaluation metrics, you might want to experiment with hyperparameter tuning for KMeans (e.g., using n_init for better results) or try different distance metrics in DBSCAN or AgglomerativeClustering.

**Predictive Modeling:** Once the clusters are formed, you could predict customer behavior or other outcomes based on the cluster labels (e.g., predicting TotalValue for new customers using their cluster information).

```
Model Comparison:
                          KMeans  DBSCAN  Agglomerative     GMM  Spectral
Silhouette Score           0.225   0.049          0.186  -0.042     0.195
Davies-Bouldin Score       1.319   1.322          1.493   2.371     1.284
Calinski-Harabasz Score  263.044  25.184        229.257  76.290   232.748

Best performing model: DBSCAN
```

**Conclusion:**

- **DBSCAN** performs the best overall, with the highest Silhouette score and acceptable Davies-Bouldin and Calinski-Harabasz scores, making it the most suitable clustering model based on these evaluation metrics.