

## ML Lab Week 13

### Clustering Lab

**Name:** Cheruku Manas Ram

**SRN:** PES2UG23CS147

**Section:** 5 'C'

**Date:** 11/11/2025

#### 1. Dimensionality Justification:

Based on the correlation heatmap and explained variance ratio from PCA, why was dimensionality reduction necessary for this dataset? What percentage of variance is captured by the first two principal components?

- A. Dimensionality reduction using PCA was beneficial for this dataset because it helps to reduce the number of features while retaining most of the important information. This can be particularly useful when dealing with potentially correlated features, as PCA transforms them into a set of uncorrelated principal components. Reducing the dimensionality can also help with visualizing the data and can sometimes improve the performance of clustering algorithms by reducing noise and the effects of the "curse of dimensionality." The first two principal components capture approximately 30% of the total variance in the data.

#### 2. Optimal Clusters

Looking at both the elbow curve and silhouette scores, what is the optimal number of clusters for this dataset? Justify your answer using both metrics.

- A. **Elbow Curve:** The elbow method suggests looking for the "elbow point" in the plot of inertia vs. the number of clusters (k). This point is where the rate of decrease in inertia slows down significantly, indicating that adding more clusters beyond this point doesn't provide much benefit in reducing within-cluster variance. Looking at the elbow curve, there is a bend around **3 or 4 clusters**.

**Silhouette Score:** The silhouette score measures how similar an object is to its own cluster compared to other clusters. A higher silhouette score indicates better-defined clusters. The plot shows the silhouette score is approximately **0.39**. While not extremely high, it's a positive score, suggesting that the data points are reasonably well-assigned to their clusters.

Considering both metrics, a choice of **3 or 4 clusters** appears reasonable. The elbow curve shows a potential elbow around these values, and the silhouette score for 3 clusters is positive.

#### 3. Cluster Characteristics: Analyze the size distribution of clusters in both K-means and Bisecting K-means. Why do you think some clusters are larger than others? What might this tell us about the customer segments?

- A. The cluster size distributions for K-means and Bisecting K-means are different. In both cases, some clusters are larger than others because the underlying data has varying densities and structures. These differences in cluster sizes suggest that

certain customer segments are more prevalent than others in the dataset. Analyzing the features within each cluster would reveal the specific characteristics defining these larger and smaller customer groups.

4. Algorithm Comparison: Compare the silhouette scores between K-means and Recursive Bisecting K-means. Which algorithm performed better for this dataset and why do you think that is?

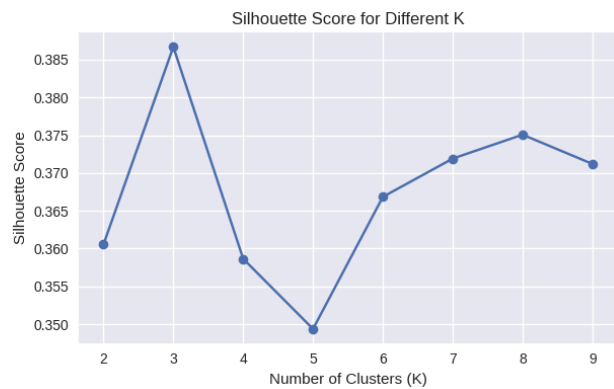
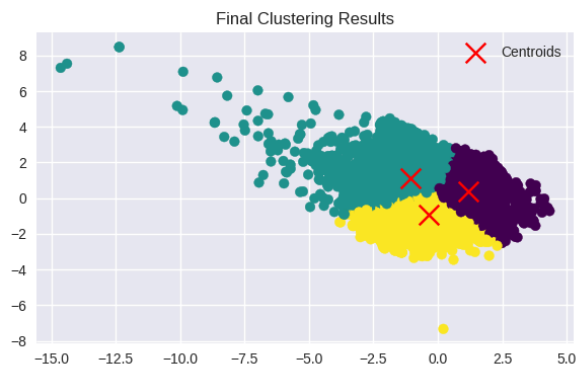
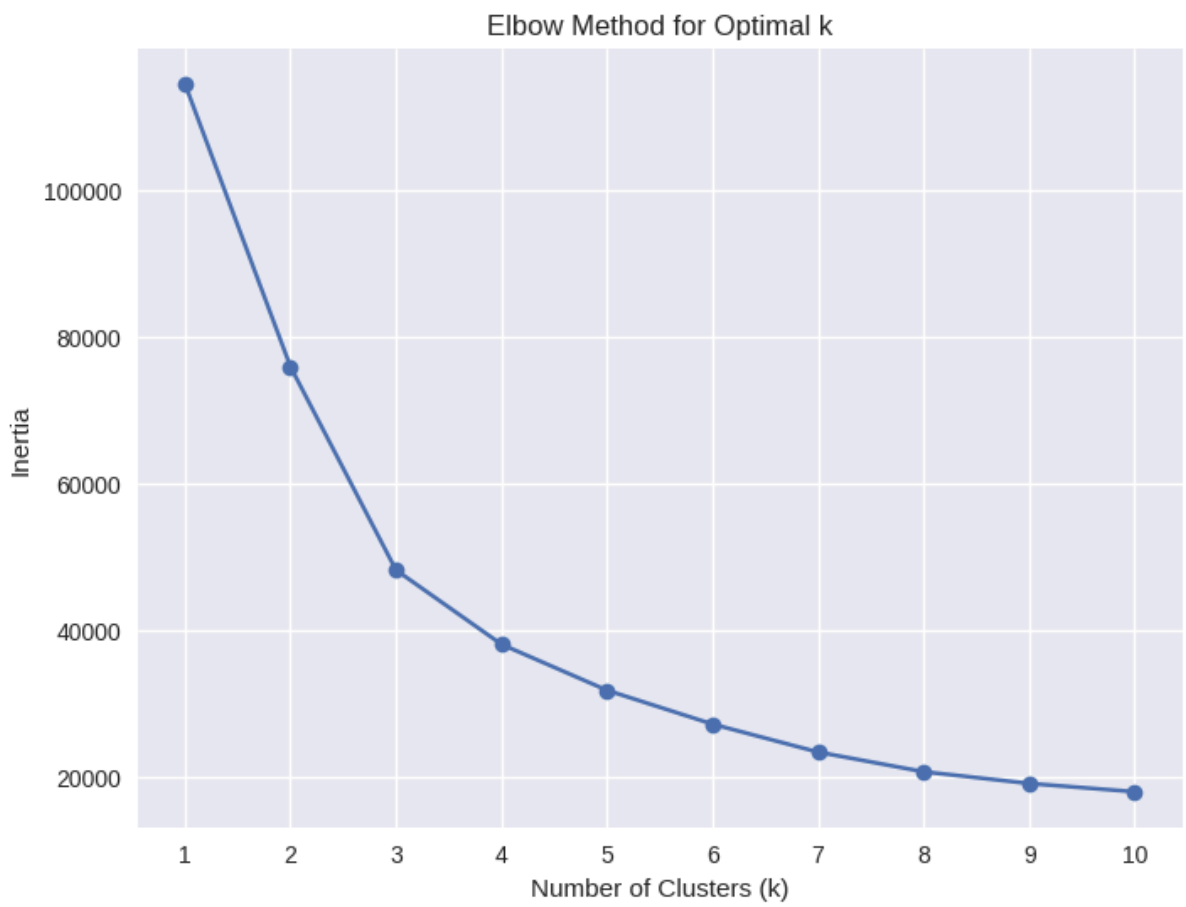
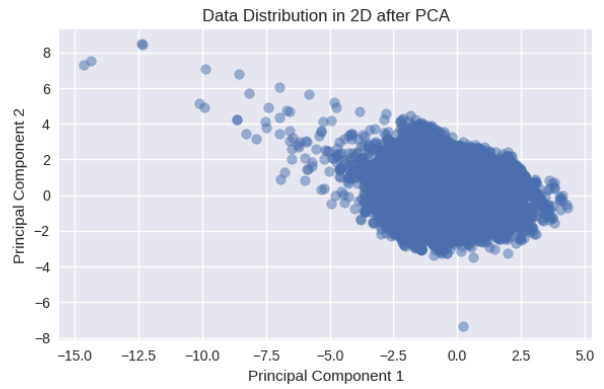
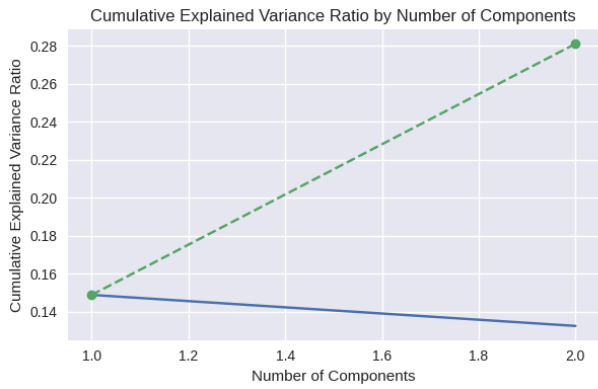
- A. Based on the silhouette scores calculated, K-means with a score of 0.39 performed better than Bisecting K-means with a score of 0.34 for this dataset. A higher silhouette score indicates that the clusters are better separated and the data points are more tightly grouped within their assigned clusters. The standard K-means likely found a more optimal partitioning of the data in this specific case, potentially due to the nature of the clusters or the initialization method used in the Bisecting K-means implementation.

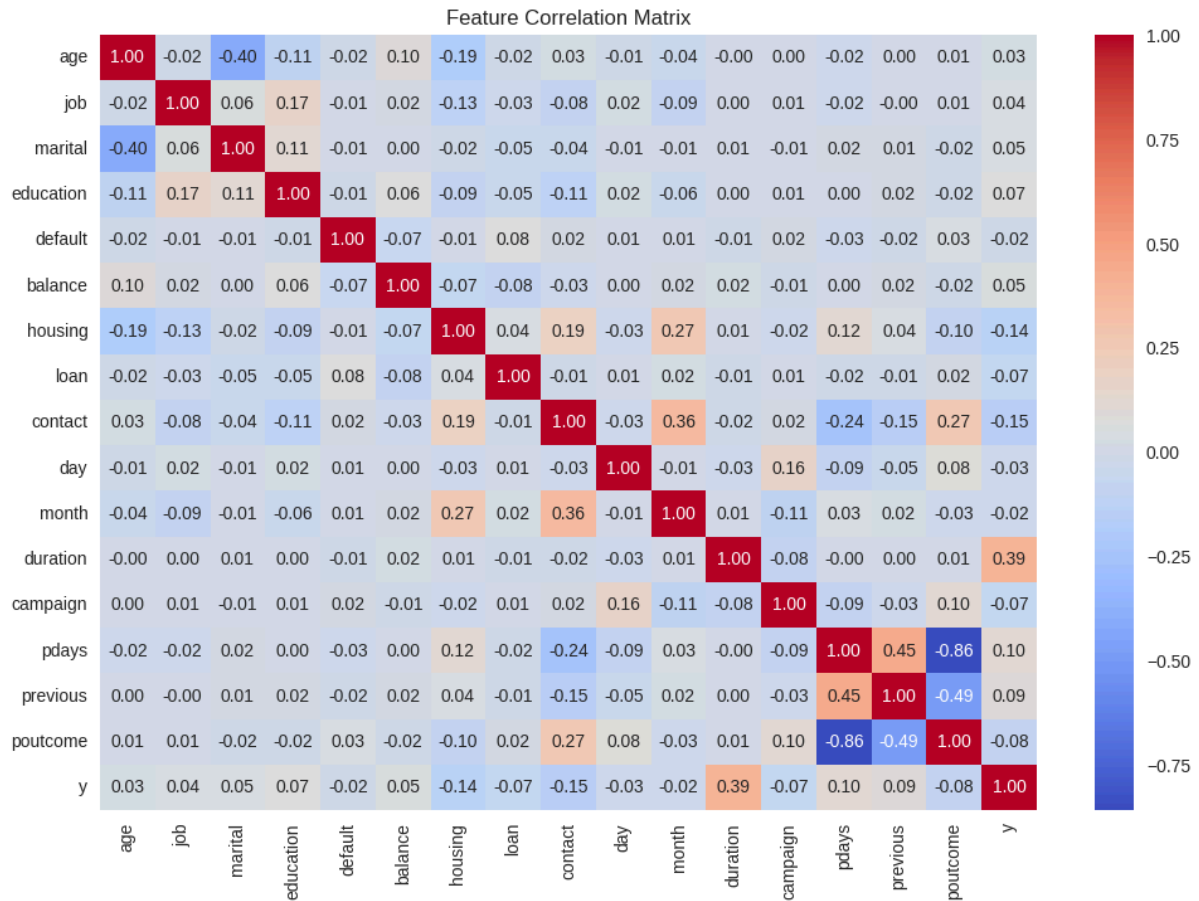
5. Business Insights: Based on the clustering results in the PCA space, what insights can you draw about customer segmentation that might be valuable for the bank's marketing strategy?

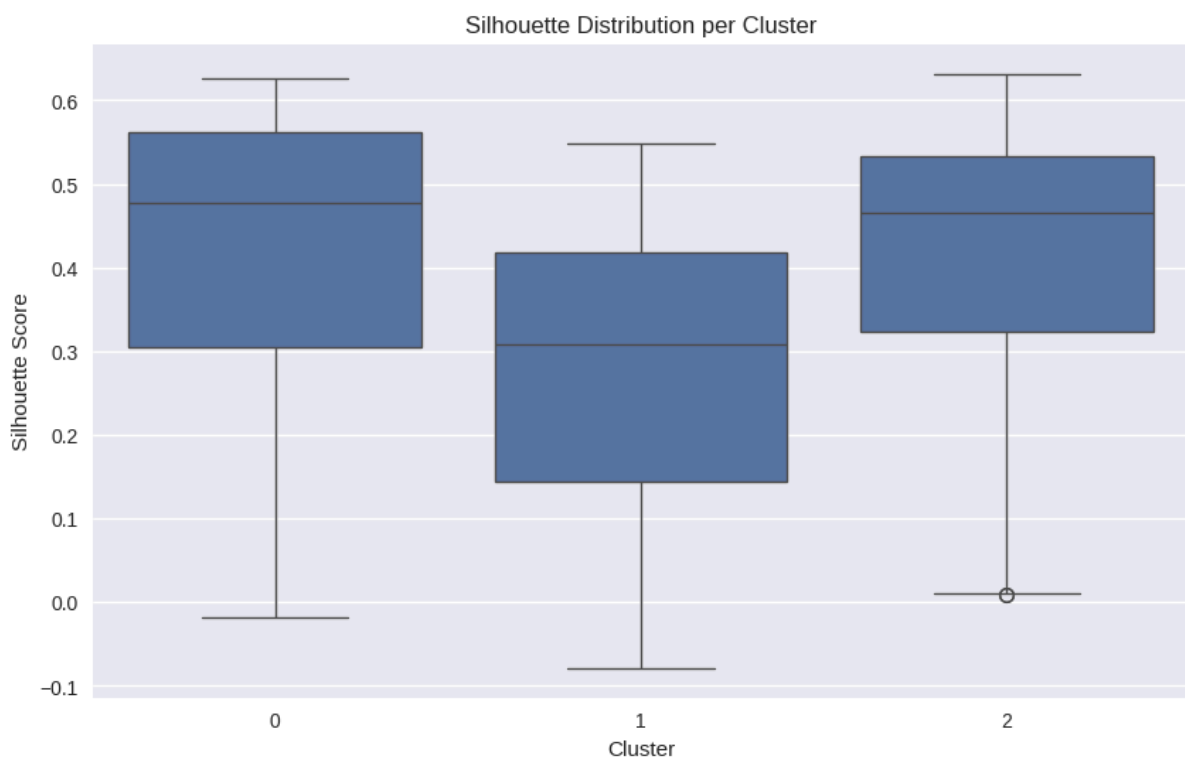
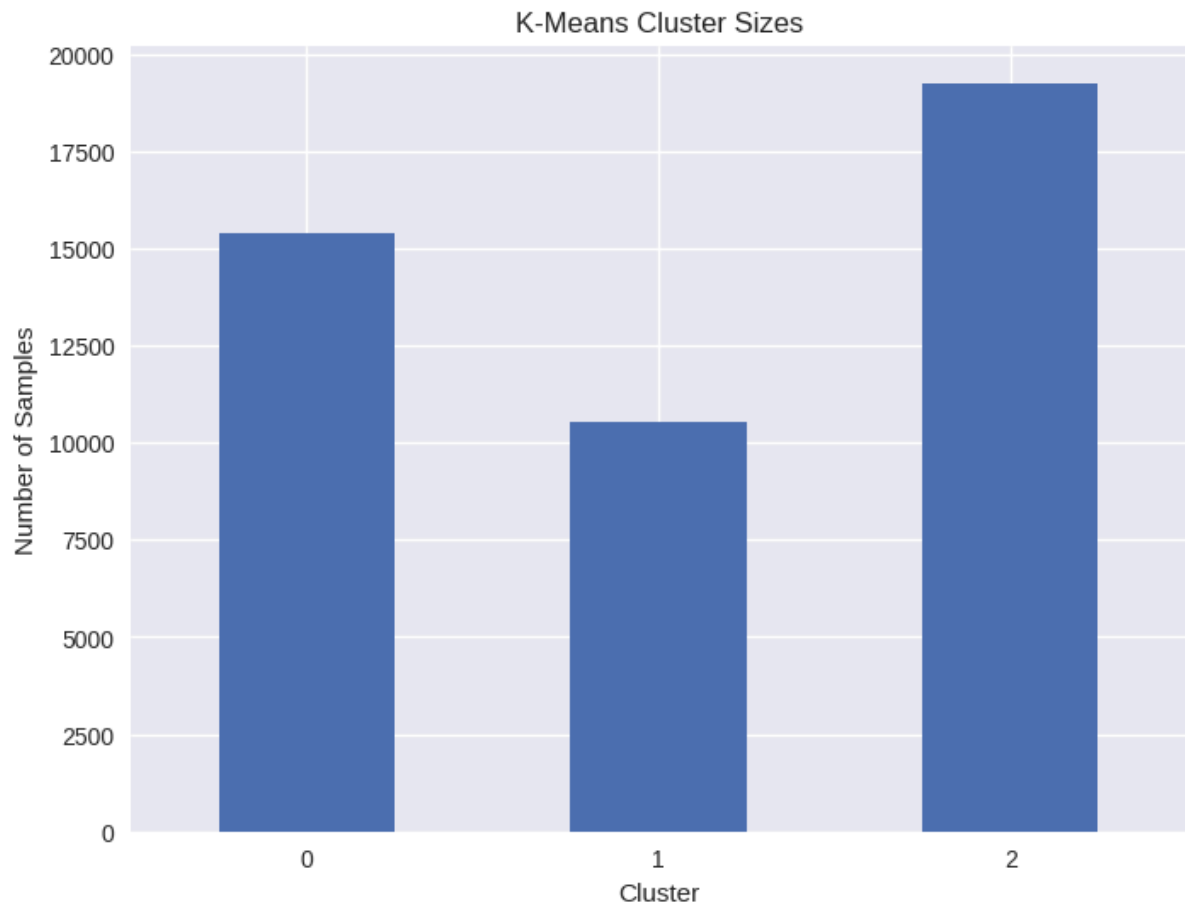
- A. Based on the clustering results visualized in the PCA space, we can observe distinct groupings of customers. These clusters represent different customer segments with unique characteristics. To draw valuable insights for the bank's marketing strategy, we would need to further analyse the original features (like age, balance, job, etc.) of the customers within each cluster. For example, one cluster might represent younger customers with lower balances, while another could consist of older, wealthier individuals. Understanding these differences will allow the bank to tailor marketing campaigns, products, and services to the specific needs and preferences of each segment, potentially leading to increased engagement and profitability.

6. Visual Pattern Recognition: In the PCA scatter plot, we see three distinct colored regions (turquoise, yellow, and purple). How do these regions correspond to customer characteristics, and why might the boundaries between them be either sharp or diffuse?

- A. In the PCA scatter plot, the three distinct colored regions (turquoise, yellow, and purple) represent the customer clusters identified by the K-means algorithm. Each color corresponds to a different cluster. To understand how these regions correspond to specific customer characteristics, you would need to analyze the original features of the data points within each cluster (e.g., what are the typical age, job, and balance of customers in the turquoise cluster?). The boundaries between these clusters might appear sharp if there are clear distinctions in characteristics between groups. The boundaries might be diffuse if there is overlap in the feature space, indicating that some customers share characteristics across different segments, making the separation less distinct.







## Recursive Bisecting K means

