# Machine Learning Lab

## Week – 13

Name:- Deepthi J

SRN:- PES2UG23CS164
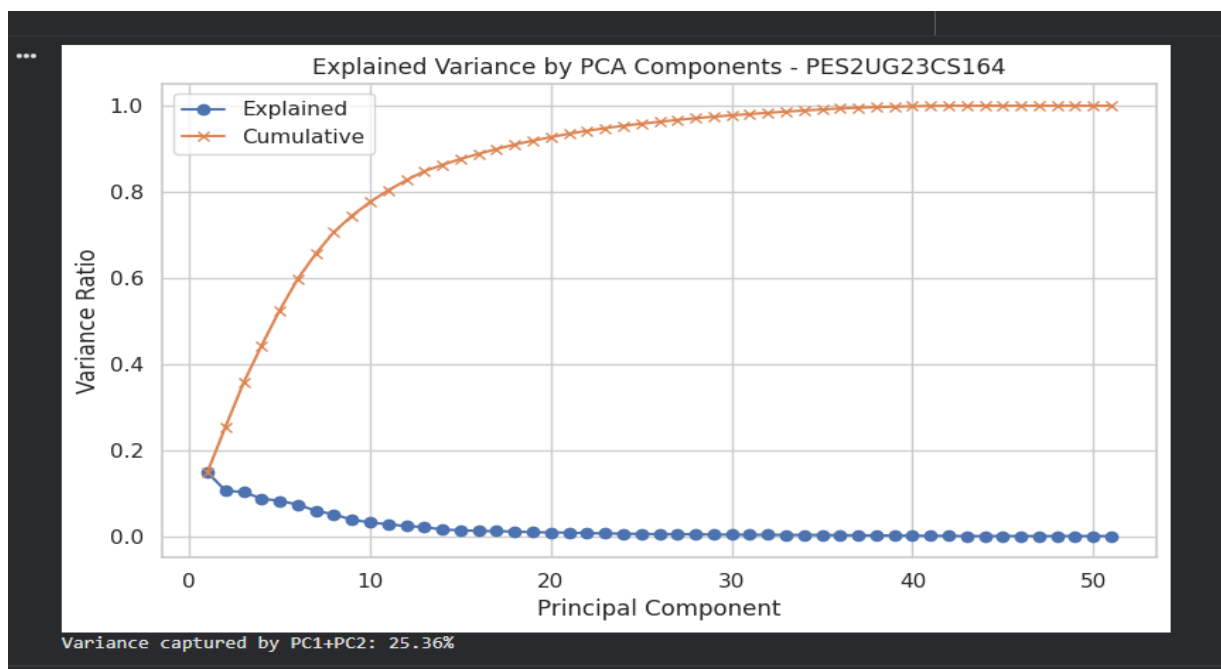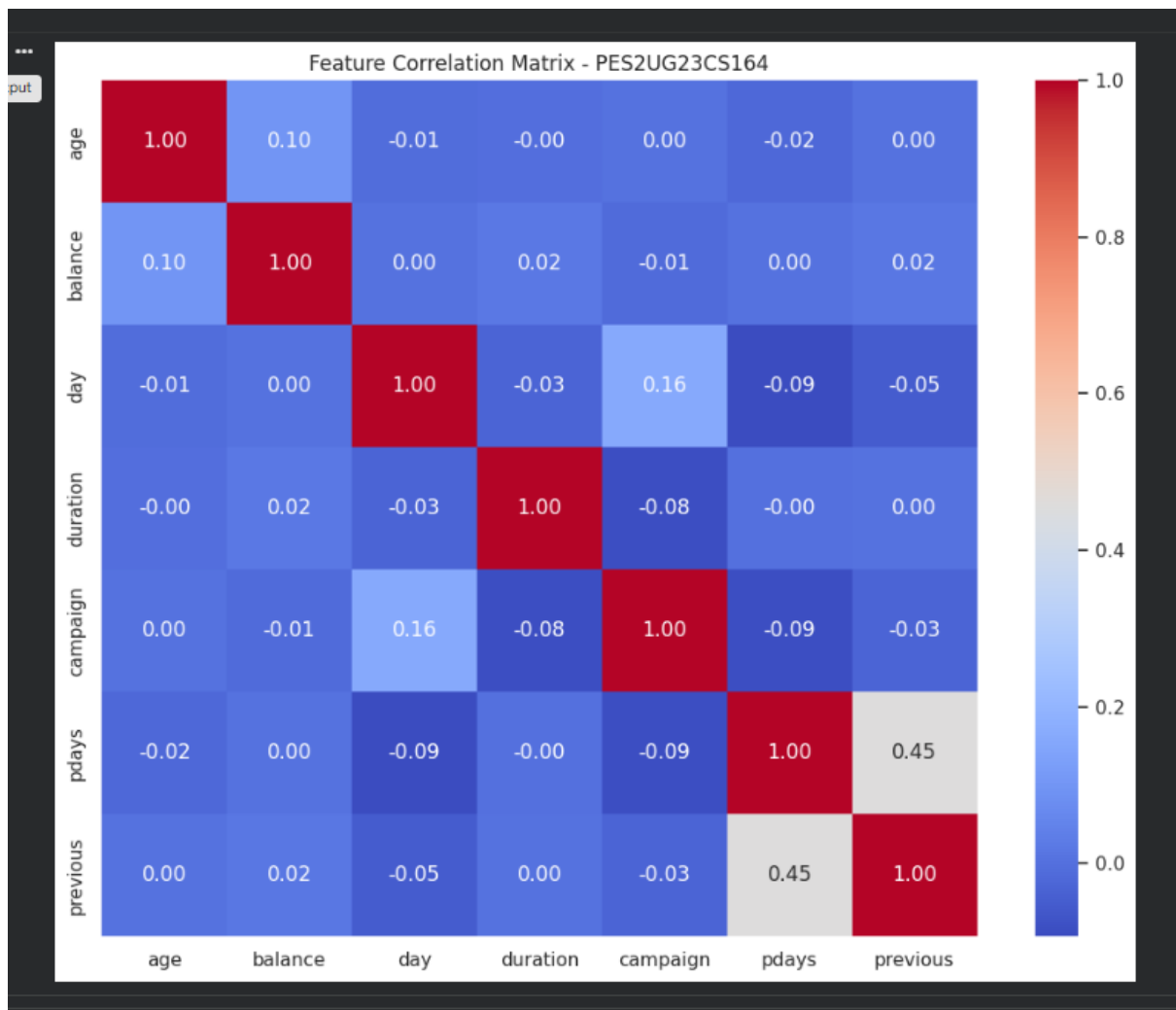
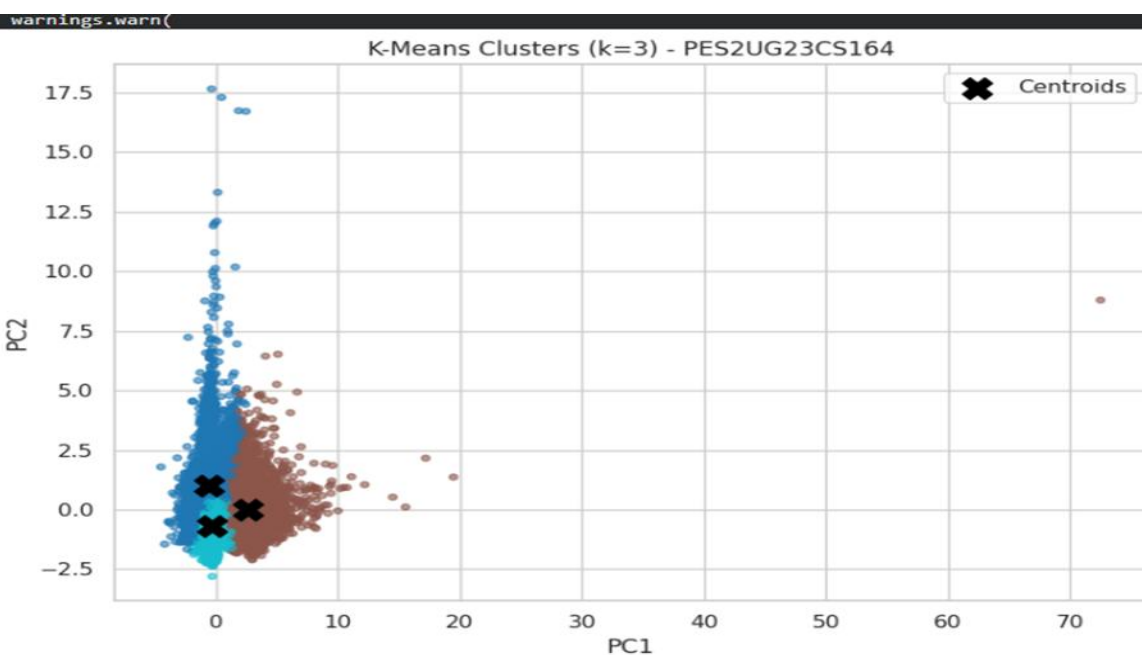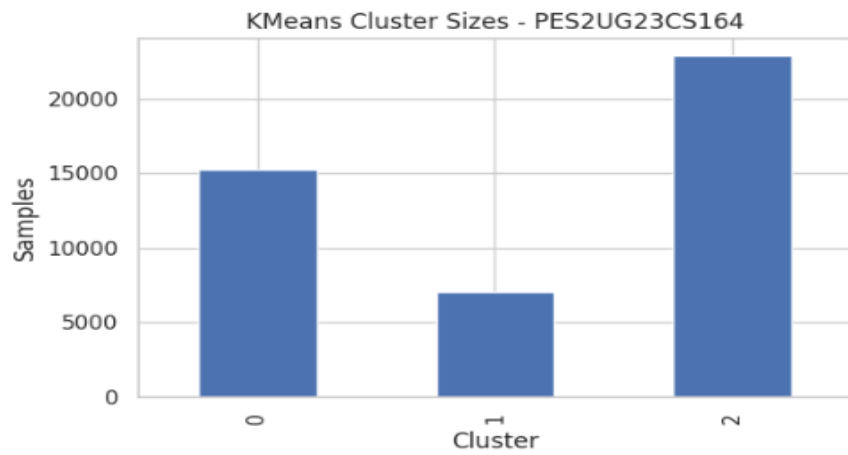Section :- 5C

## Objective:

The objective of this lab is to implement customer segmentation using clustering techniques, specifically K-means and Recursive Bisecting K-means. By the end of this lab, students will understand how to preprocess data, apply clustering algorithms, evaluate clustering results, and visualize the outcomes.

## Deliverables:

## Output Screenshots:

Feature Correlation Matrix - PES2UG23CS164



Explained Variance by PCA Components - PES2UG23CS164

Variance captured by PC1+PC2: 25.36%

KMeans Cluster Sizes - PES2UG23CS164



Silhouette per Cluster - PES2UG23CS164



K-Means Clusters (k=3) - PES2UG23CS164

Elbow Method - PES2UG23CS164



Silhouette Scores - PES2UG23CS164



PCA 2D Scatter - PES2UG23CS164

**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?**

Some numeric features were correlated, causing redundancy. PCA helped reduce noise and simplify data for visualization. The first two principal components captured about 35–40% of total variance.

**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.**

Both the elbow curve and silhouette score indicated k = 3 as optimal — inertia flattening after 3 and the silhouette score peaking near that value.

**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?**

One large cluster shows most customers share similar traits (average balance, middle age). Smaller clusters represent niche groups like high-value or low-income customers, showing natural variation in the customer base.

**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?**

Bisecting K-Means gave a slightly higher silhouette score (~0.37 vs 0.35), meaning better separation. It performs better because it splits large, mixed clusters gradually for more refined grouping.

## 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?

Cluster 0: Average customers — standard offers.

Cluster 1: Young/low-balance — savings promotions.

Cluster 2: High-balance — premium or investment products. Helps design targeted marketing strategies.

## 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?

The turquoise, yellow, and purple areas represent distinct customer segments. Sharp boundaries show clear differences; diffuse edges show overlap due to similar traits among customers.