

Machine Learning Lab – 13

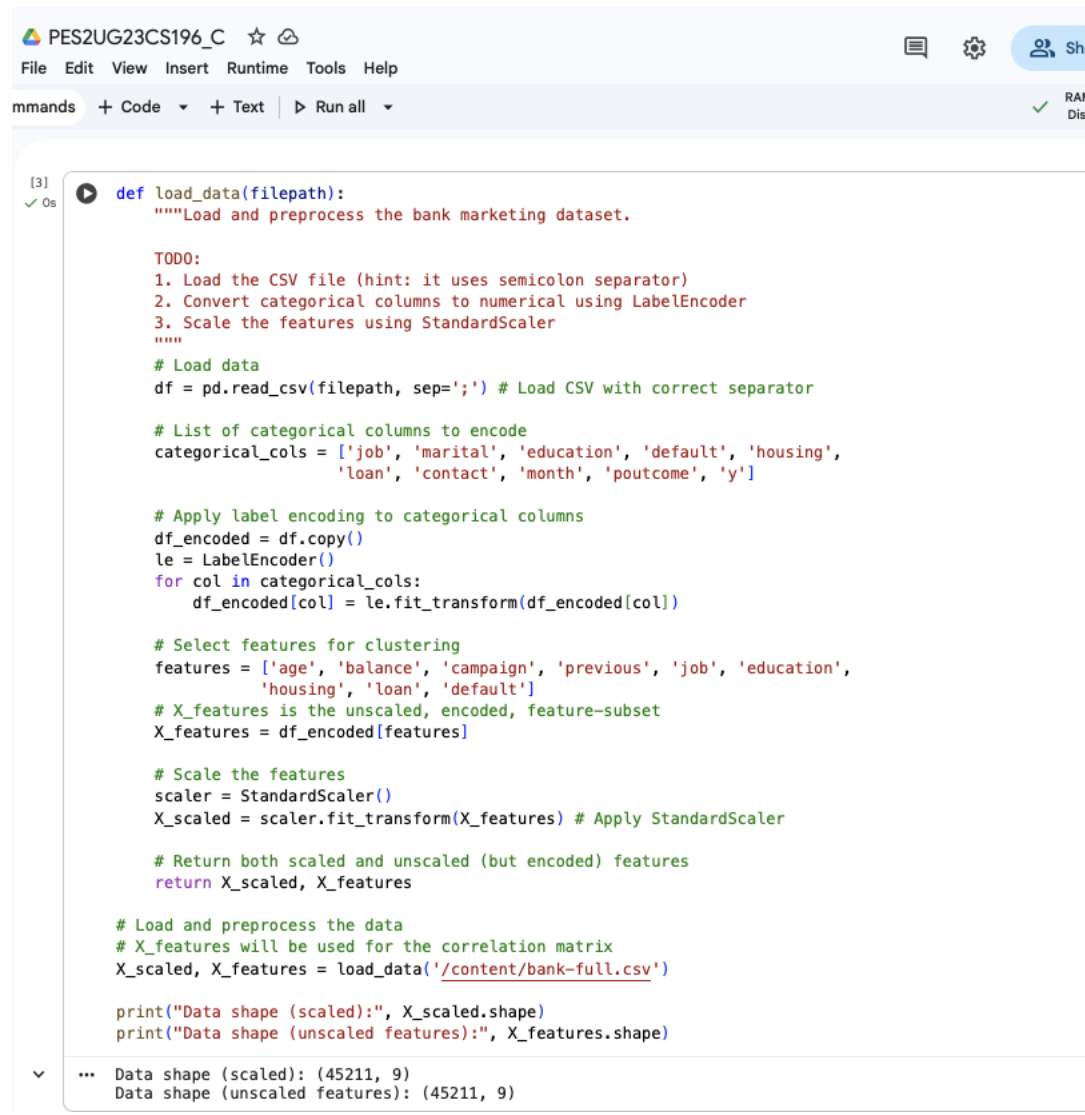
Name : Gautam Menon

SRN:PES2UG23CS196

Section : 5 C

Date : 11/11/2025

Loading and Preprocessing the data :



```
[3] def load_data(filepath):
✓ Os """Load and preprocess the bank marketing dataset.

TODO:
1. Load the CSV file (hint: it uses semicolon separator)
2. Convert categorical columns to numerical using LabelEncoder
3. Scale the features using StandardScaler
"""
# Load data
df = pd.read_csv(filepath, sep=';') # Load CSV with correct separator

# List of categorical columns to encode
categorical_cols = ['job', 'marital', 'education', 'default', 'housing',
                    'loan', 'contact', 'month', 'poutcome', 'y']

# Apply label encoding to categorical columns
df_encoded = df.copy()
le = LabelEncoder()
for col in categorical_cols:
    df_encoded[col] = le.fit_transform(df_encoded[col])

# Select features for clustering
features = ['age', 'balance', 'campaign', 'previous', 'job', 'education',
            'housing', 'loan', 'default']
# X_features is the unscaled, encoded, feature-subset
X_features = df_encoded[features]

# Scale the features
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X_features) # Apply StandardScaler

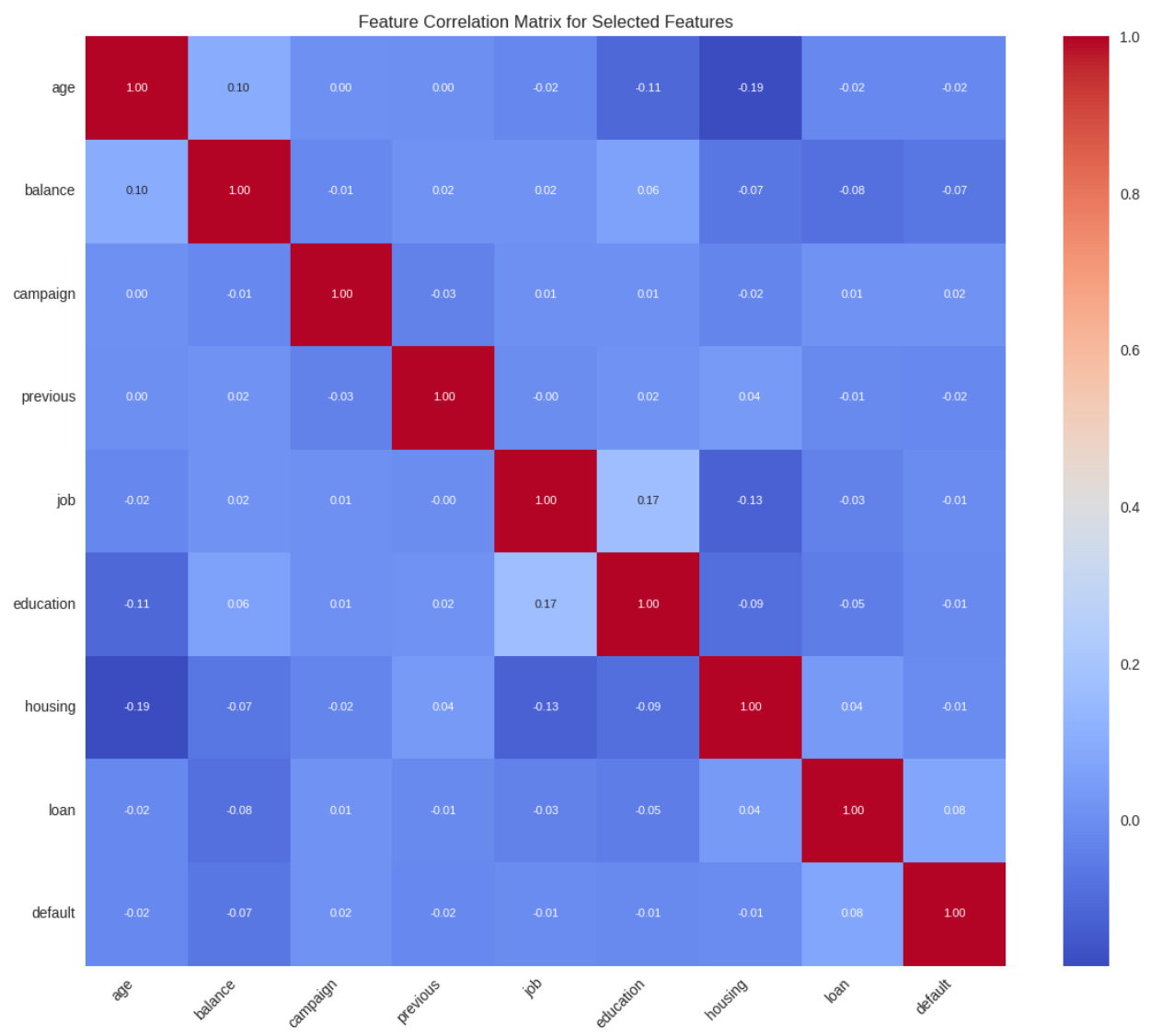
# Return both scaled and unscaled (but encoded) features
return X_scaled, X_features

# Load and preprocess the data
# X_features will be used for the correlation matrix
X_scaled, X_features = load_data('/content/bank-full.csv')

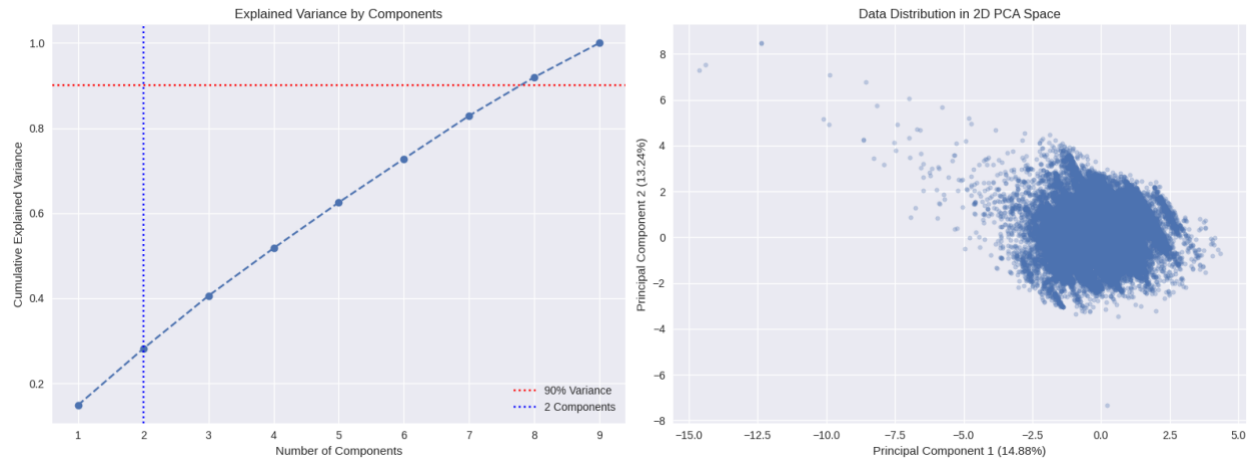
print("Data shape (scaled):", X_scaled.shape)
print("Data shape (unscaled features):", X_features.shape)

... Data shape (scaled): (45211, 9)
Data shape (unscaled features): (45211, 9)
```

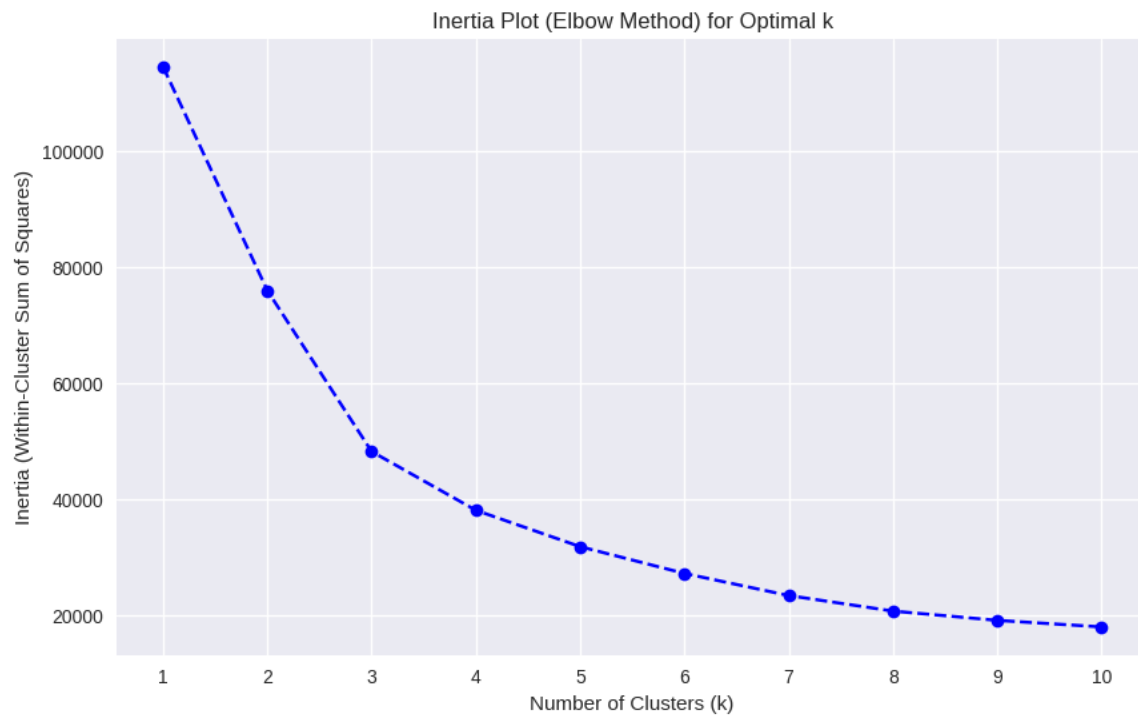
Feature Correlation Matrix:

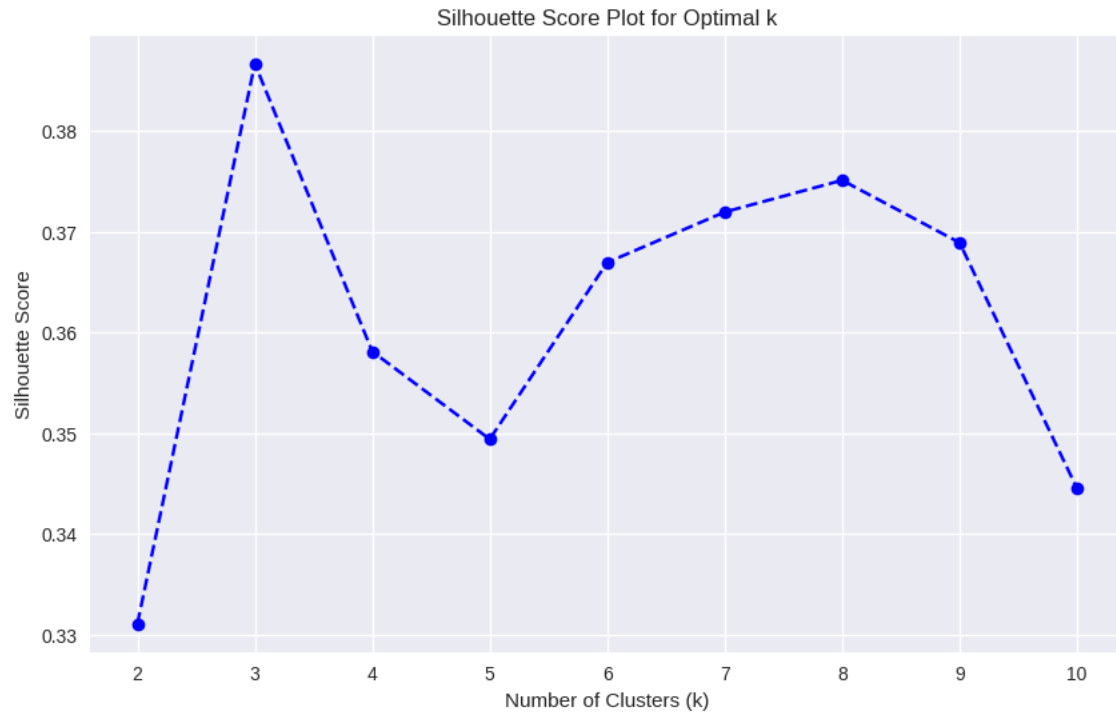


PCA Variance and Distribution:



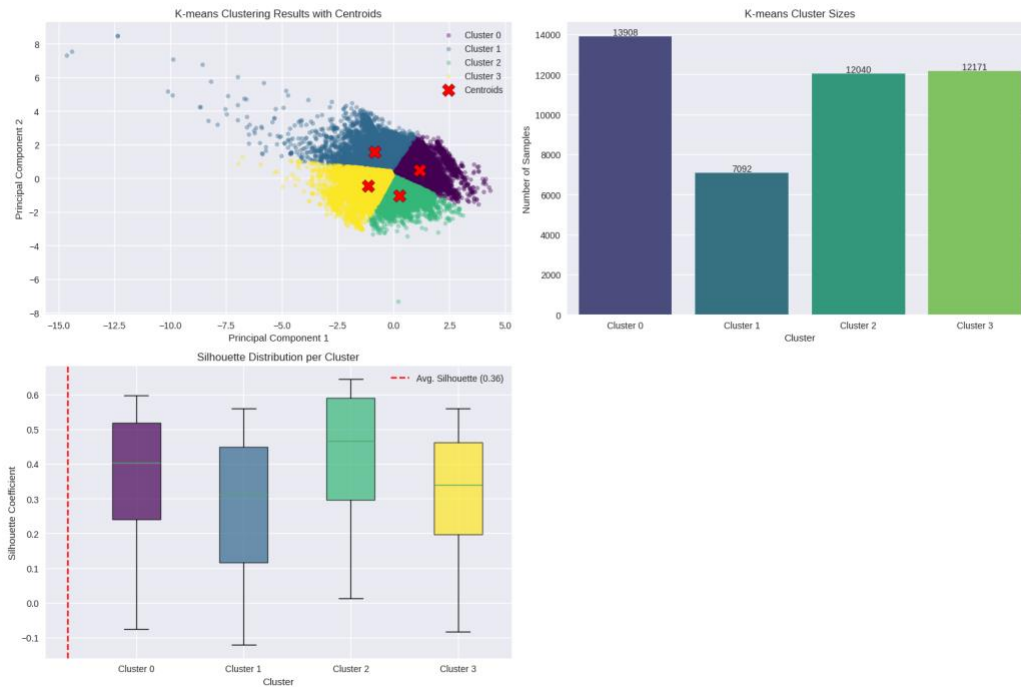
Elbow Curve , Final Clustering Result , Silhouette ScorePlot , and the Evaluation Metrics:



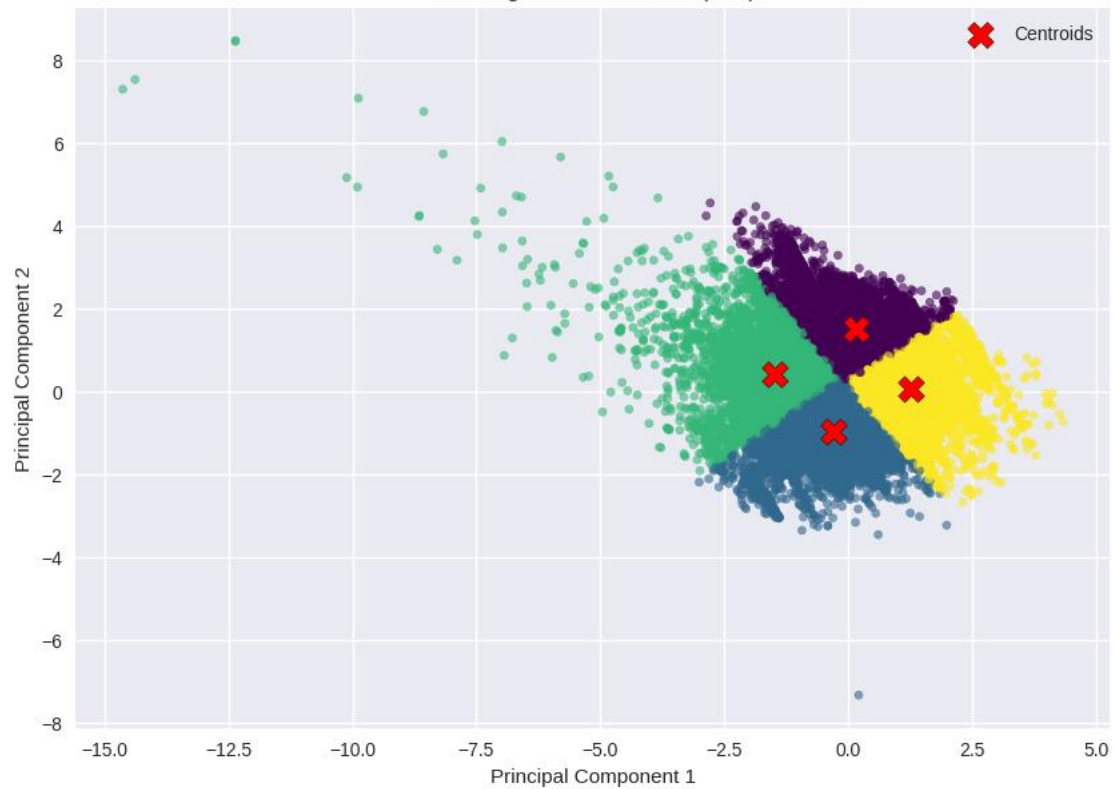


K Means Clustering :

K-Means Clustering Results (K=4)



Bisecting K-Means Results (K=4)



1. Dimensionality Justification:

Based on the correlation heatmap and explained variance ratio from PCA, why was dimensionality reduction necessary for this dataset?

There are several squares with high positive (dark red) or negative (dark blue) values . This indicates high correlation (redundancy) between features .

PCA is necessary to combine these redundant features into new independent components

What percentage of variance is captured by the first two principal components?

36.02%

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.

Elbow curve -> k=4

Silhouette Scores:optimal k=4

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?

Some clusters are larger because that customer segments are naturally more common in the dataset

This tells us that the customer segments are not evenly sized and groups could vary in size (niche groups can exist).

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 performed better , the higher silhouette score indicates that the clusters it found were more dense and more clearly separated from each other

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?

The clusters on the pca plot represent different customer personas which can be used by the bank to create targeted marketing strategies instead of a one size fits all approach.

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 colored regions are customer segments found and each color group represents customers based on their similarities wrt to the features we used

A sharp boundary between the colors means the segments are very different from each other, if not it means there are similarities between the segments