

IEEE Machine Learning Hackathon at Nazarbayev University

“Smart Tech” team

Miras Mukhametkaziev, Almansur Kakimov

Desht Lab Case

Feb 23, 2025

1. Data Loading and Exploration

The dataset (WVS_Cross-National_Wave_7_csv_v6_0.csv) is loaded into a Pandas DataFrame.

The initial structure of the dataset is examined:

- Shape: (97220, 613) (97,220 rows and 613 columns)
- Data types: float64, int64, object
- Missing values: Checked per column, with some columns having significant missing data.
- Summary statistics (mean, std, min, max, quartiles) are computed.

2. Data Preprocessing

Missing Values Handling

- Columns with more than 50% missing values are dropped.
- Numerical missing values are filled with the median.
- Categorical missing values are replaced with the mode.

Feature Engineering

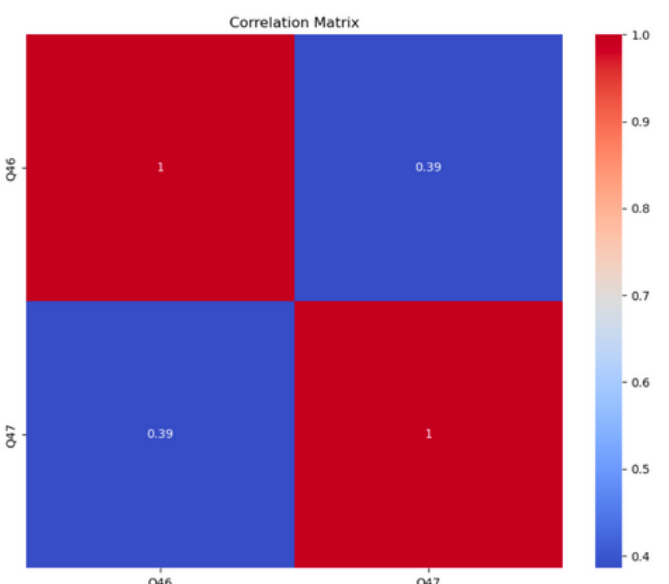
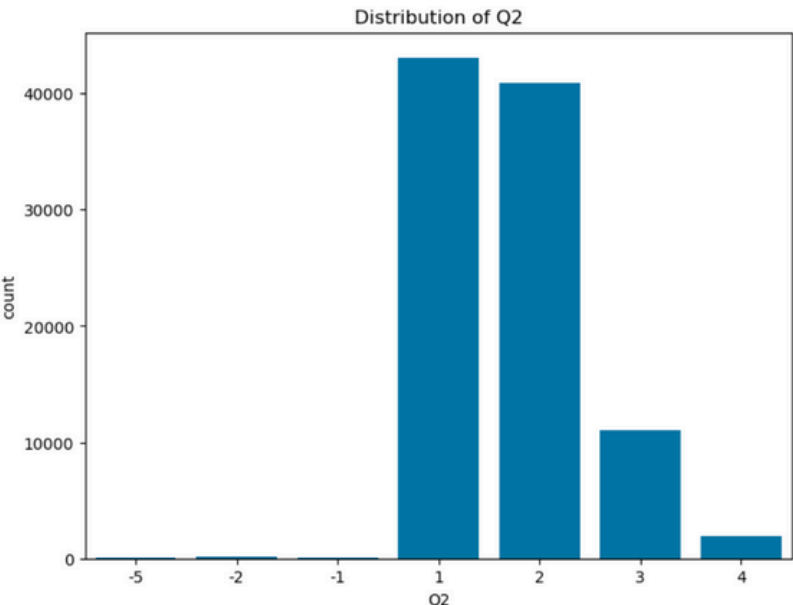
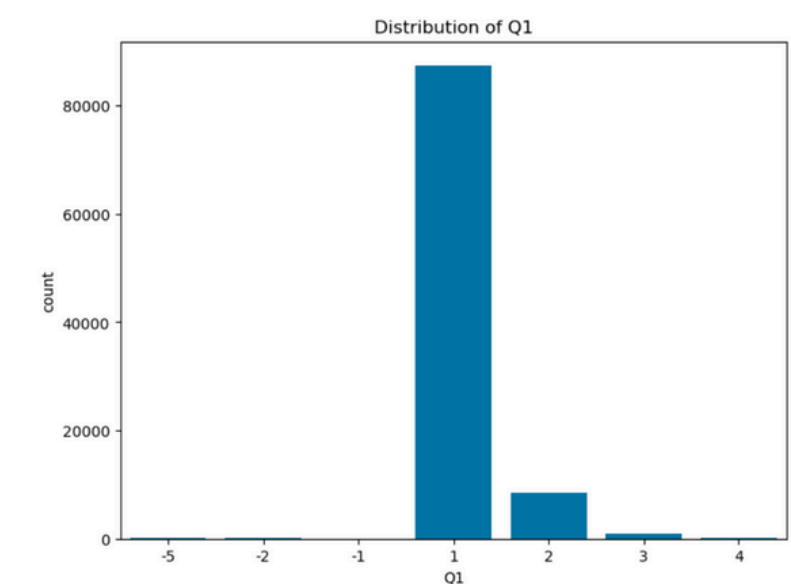
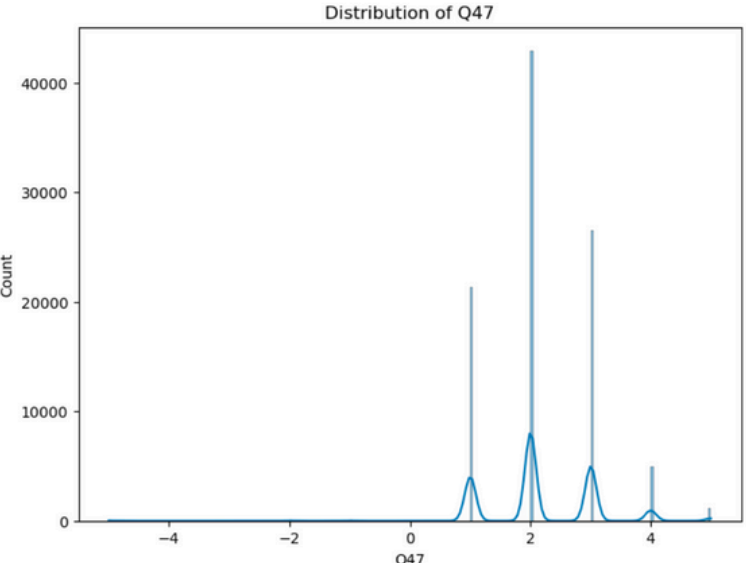
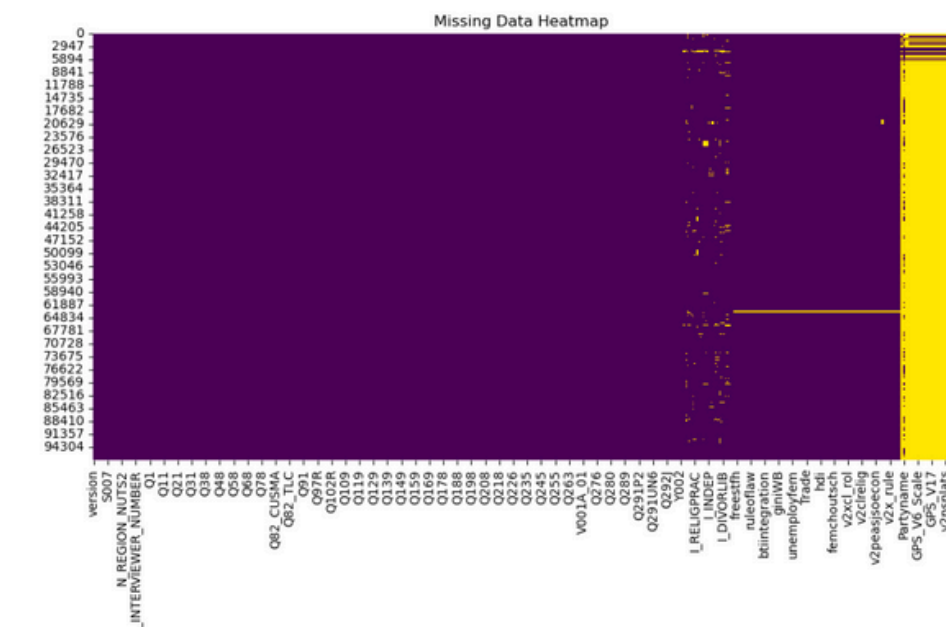
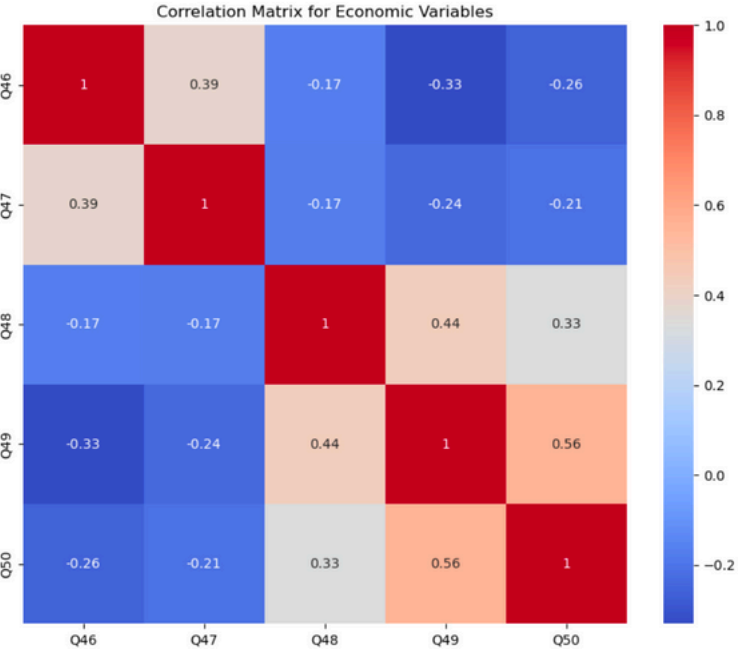
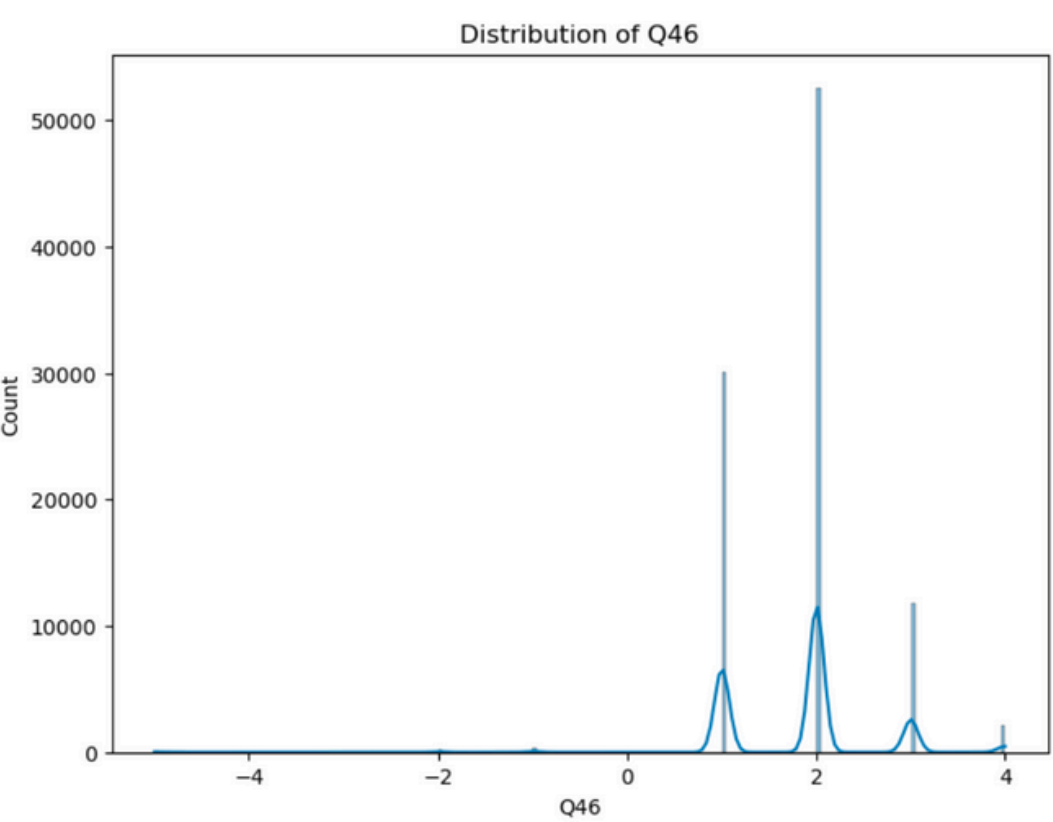
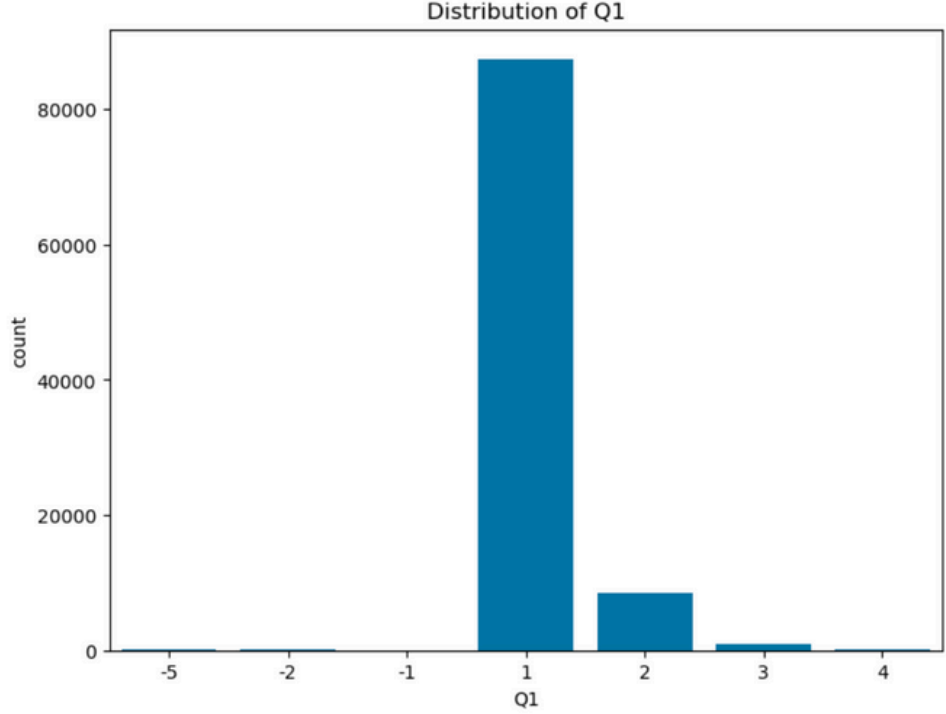
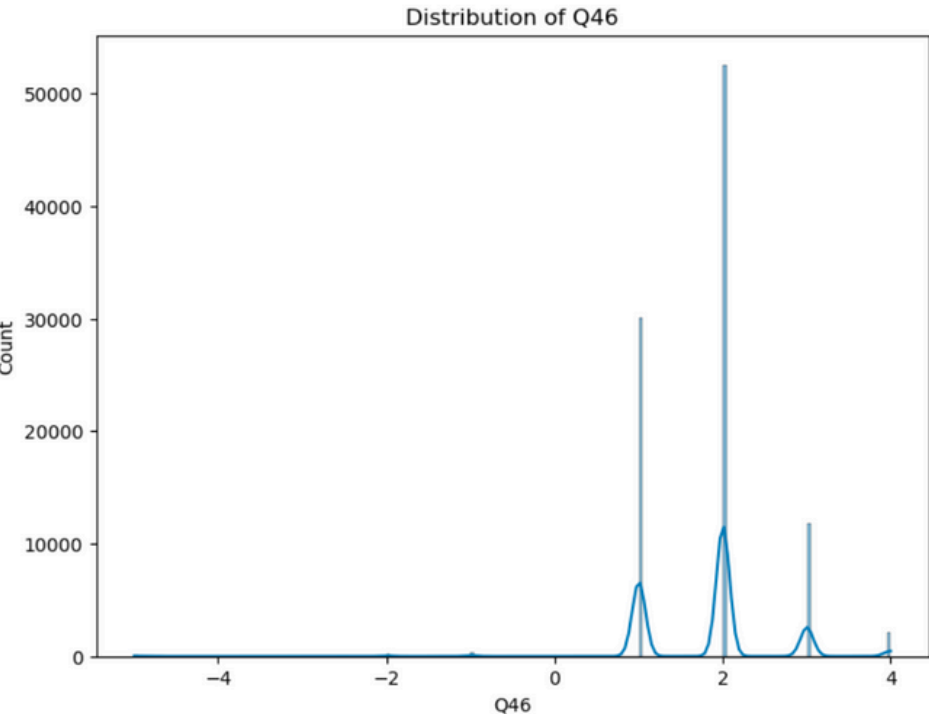
- Selected columns for economic and societal values analysis.
- Created histograms, bar plots, and correlation matrices.

Standardization

- Numerical columns are standardized using StandardScaler.

One-Hot Encoding

- Categorical variables are converted into numerical format via `pd.get_dummies()`.



3. Clustering Methods

K-Means Clustering

MiniBatchKMeans is used for clustering:

- A sample of 10,000 observations is taken for efficiency.

The Elbow Method (inertia) and Silhouette Score are used to determine the optimal number of clusters (K).

K=4 is selected based on the Elbow and Silhouette analysis.

Final K-Means Model: the dataset is clustered into 4 groups.

Distribution:

- Cluster 1: 41,970 samples
- Cluster 3: 26,916 samples
- Cluster 0: 19,522 samples
- Cluster 2: 8,812 samples

Cluster labels are added to the dataset.

DBSCAN Clustering

DBSCAN (Density-Based Spatial Clustering) is also applied:

Parameters: $\text{eps}=0.5$, $\text{min_samples}=5$

Many samples (-1) are classified as noise, meaning DBSCAN is not very effective on this dataset.

4. Clustering Evaluation

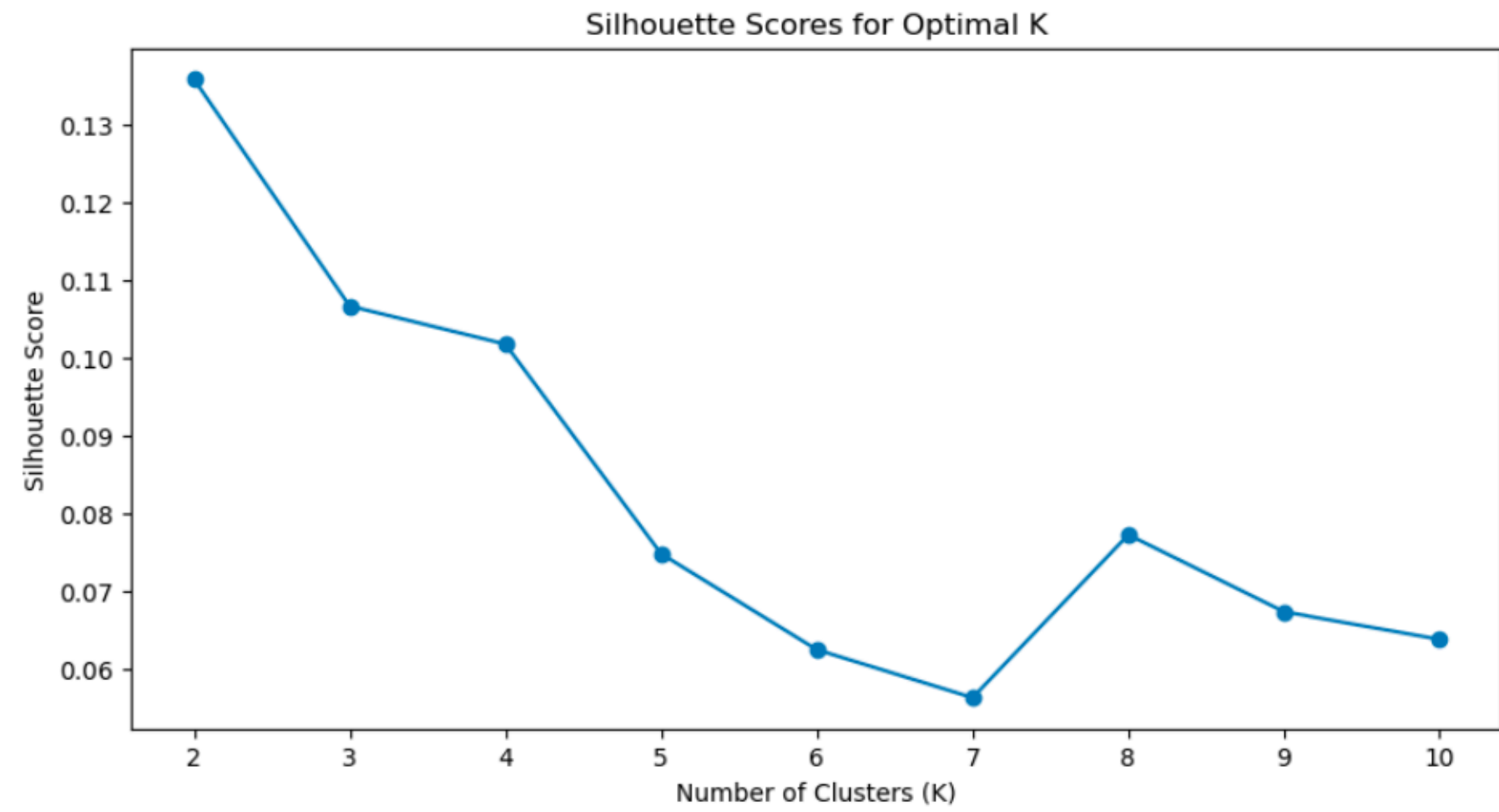
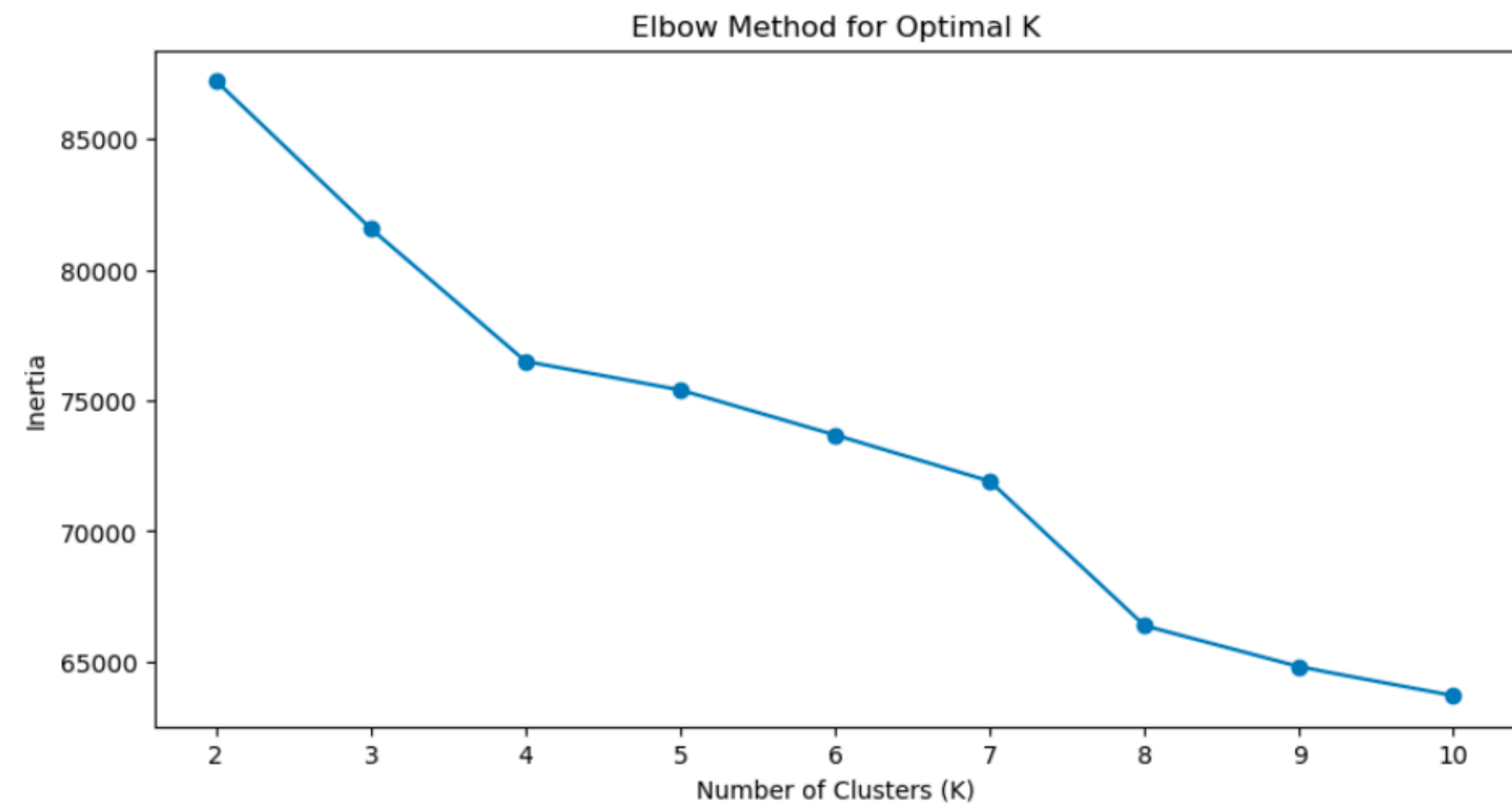
Silhouette Score

- K-Means: 0.120 (low but acceptable)
- DBSCAN: -0.220 (poor clustering)

Davies-Bouldin Index

- K-Means: 2.195 (lower is better)
- DBSCAN: 1.539

DBSCAN performs poorly compared to K-Means, likely due to the dataset's high dimensionality and sparse density regions.



5. Dimensionality Reduction and Visualization

PCA (Principal Component Analysis)

- Reduced dataset to 2 components for visualization.
- Scatter plot of clusters.

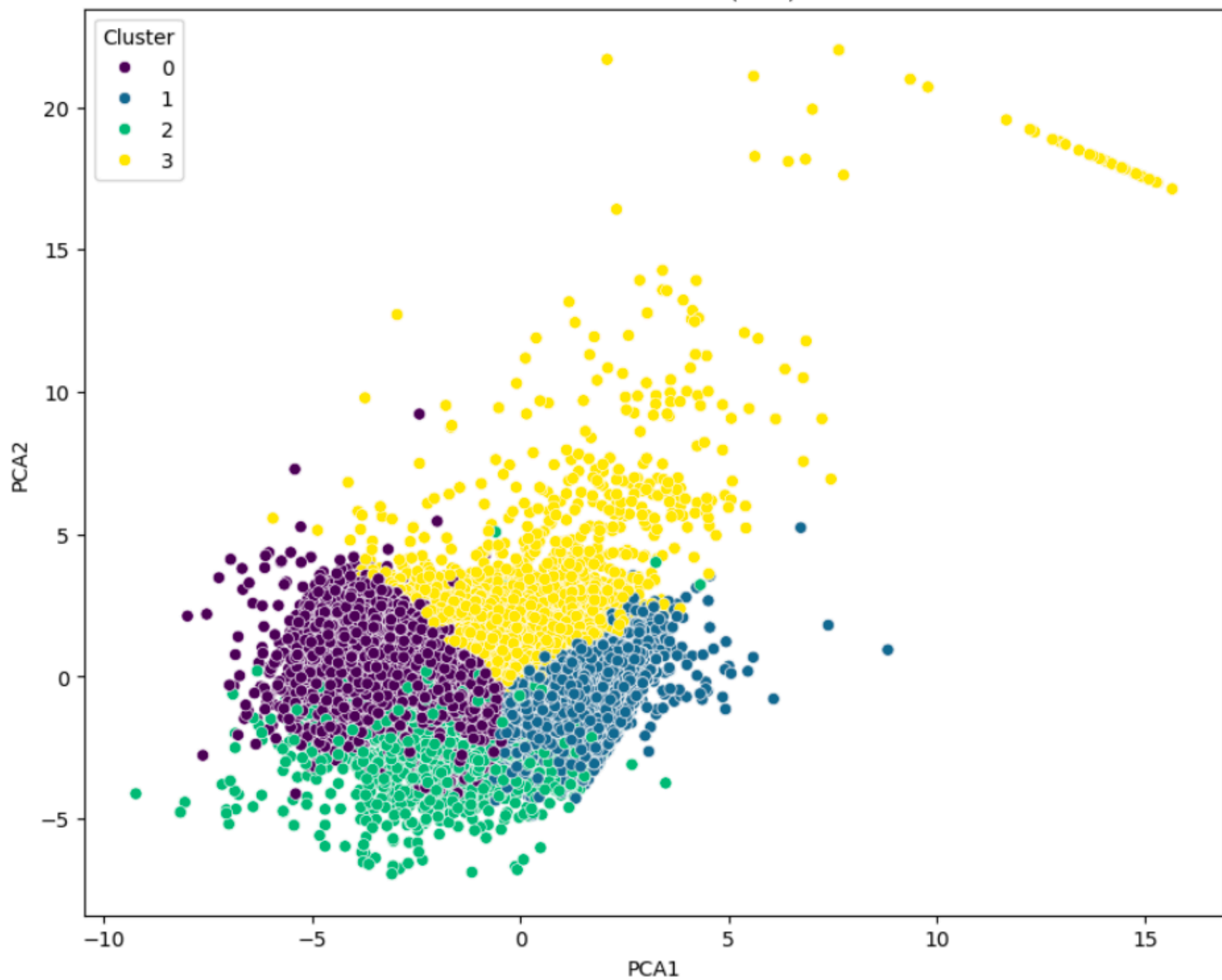
t-SNE (t-Distributed Stochastic Neighbor Embedding)

- Another technique to reduce dimensions.
- Scatter plot of clusters with better separability.

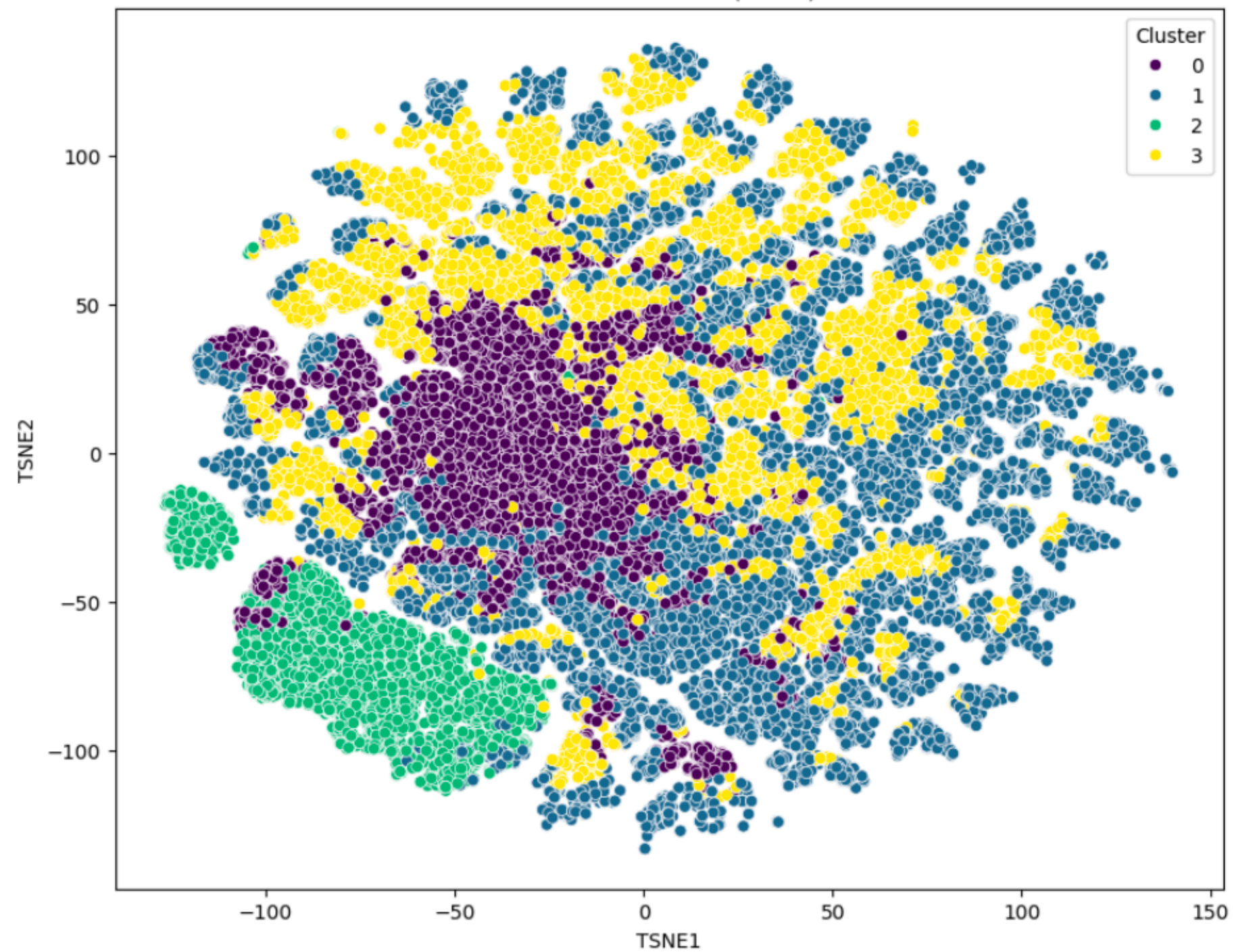
Parallel Coordinates Plot

- Used to visualize feature differences across clusters.

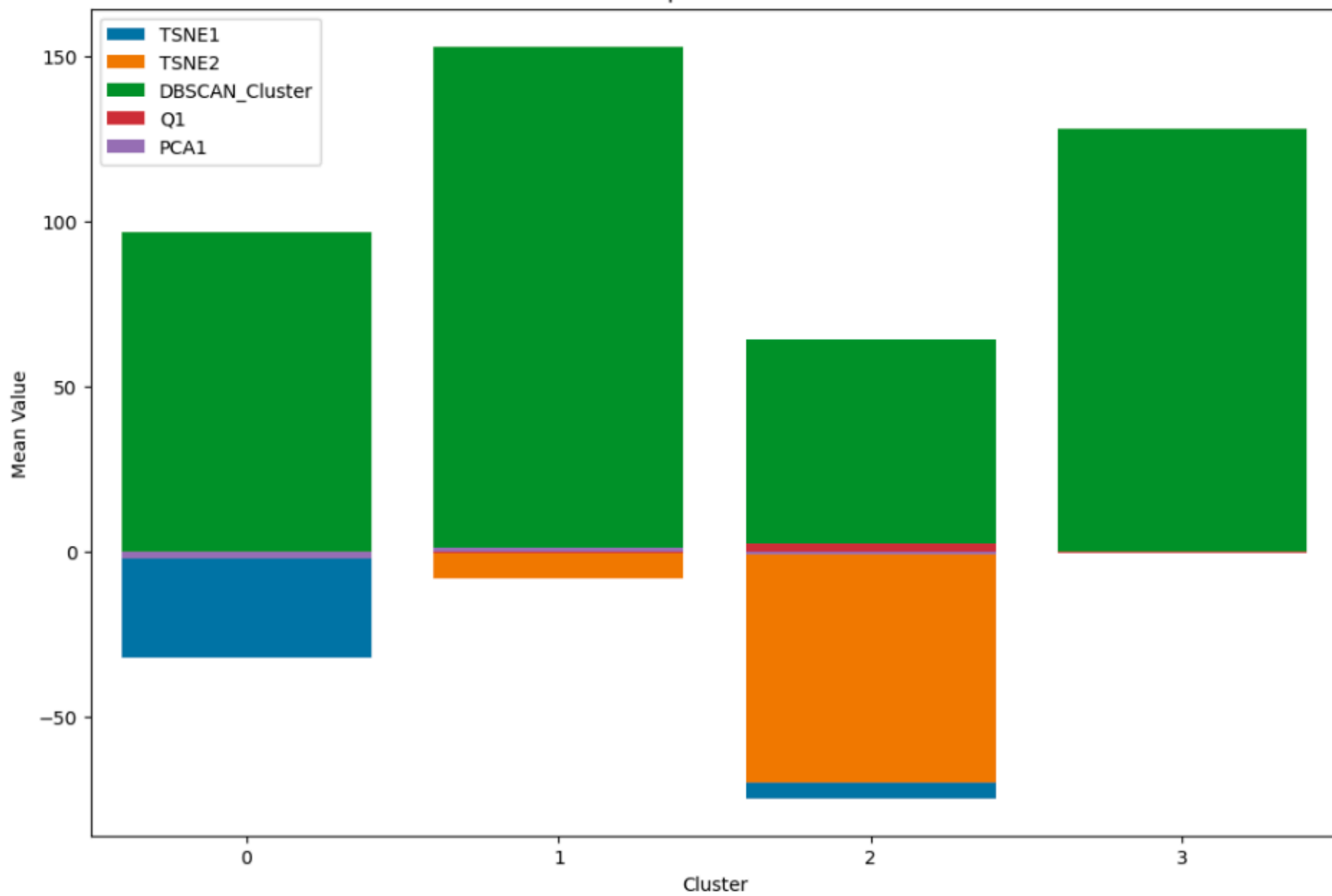
K-Means Clusters (PCA)



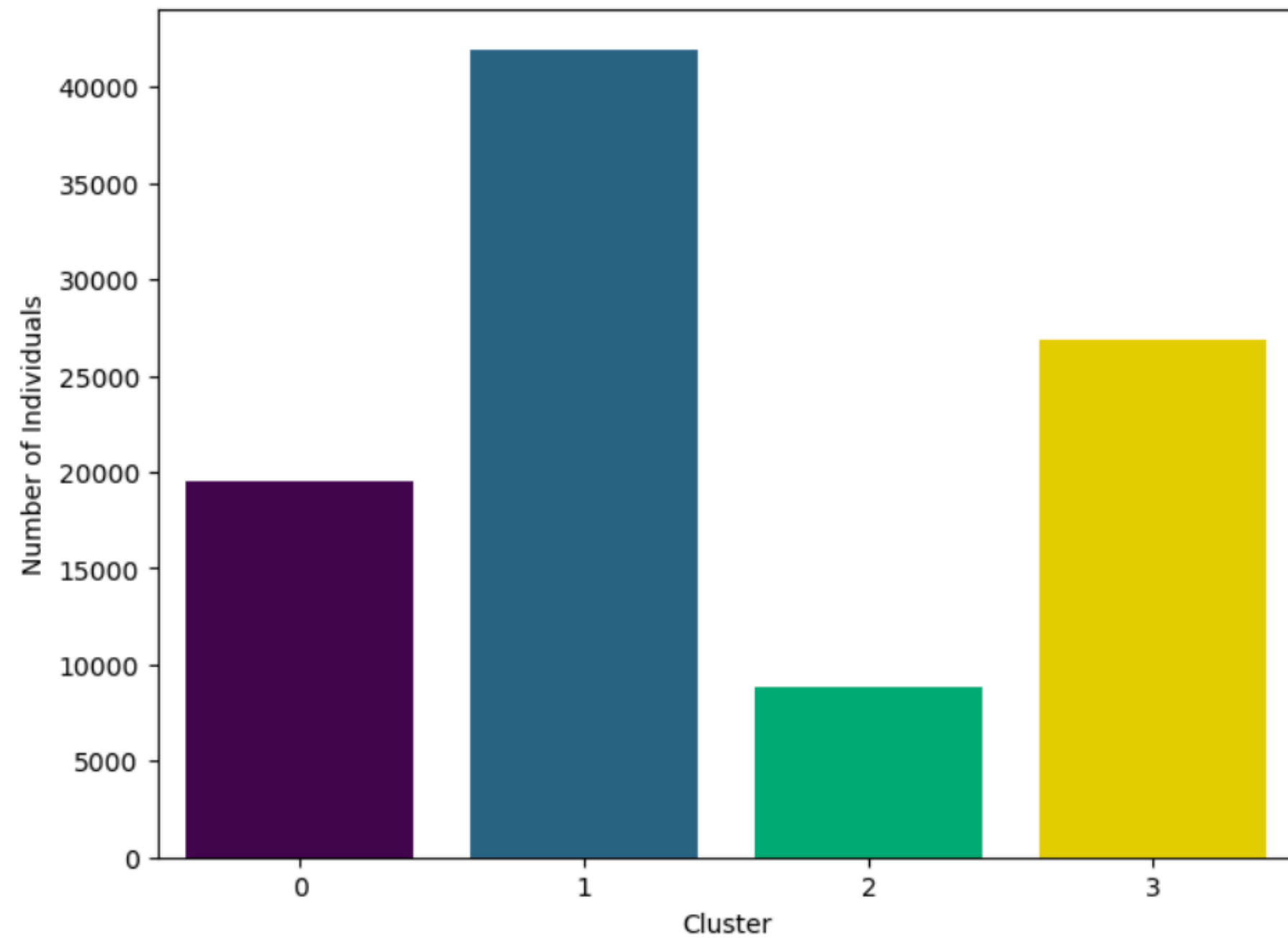
K-Means Clusters (t-SNE)



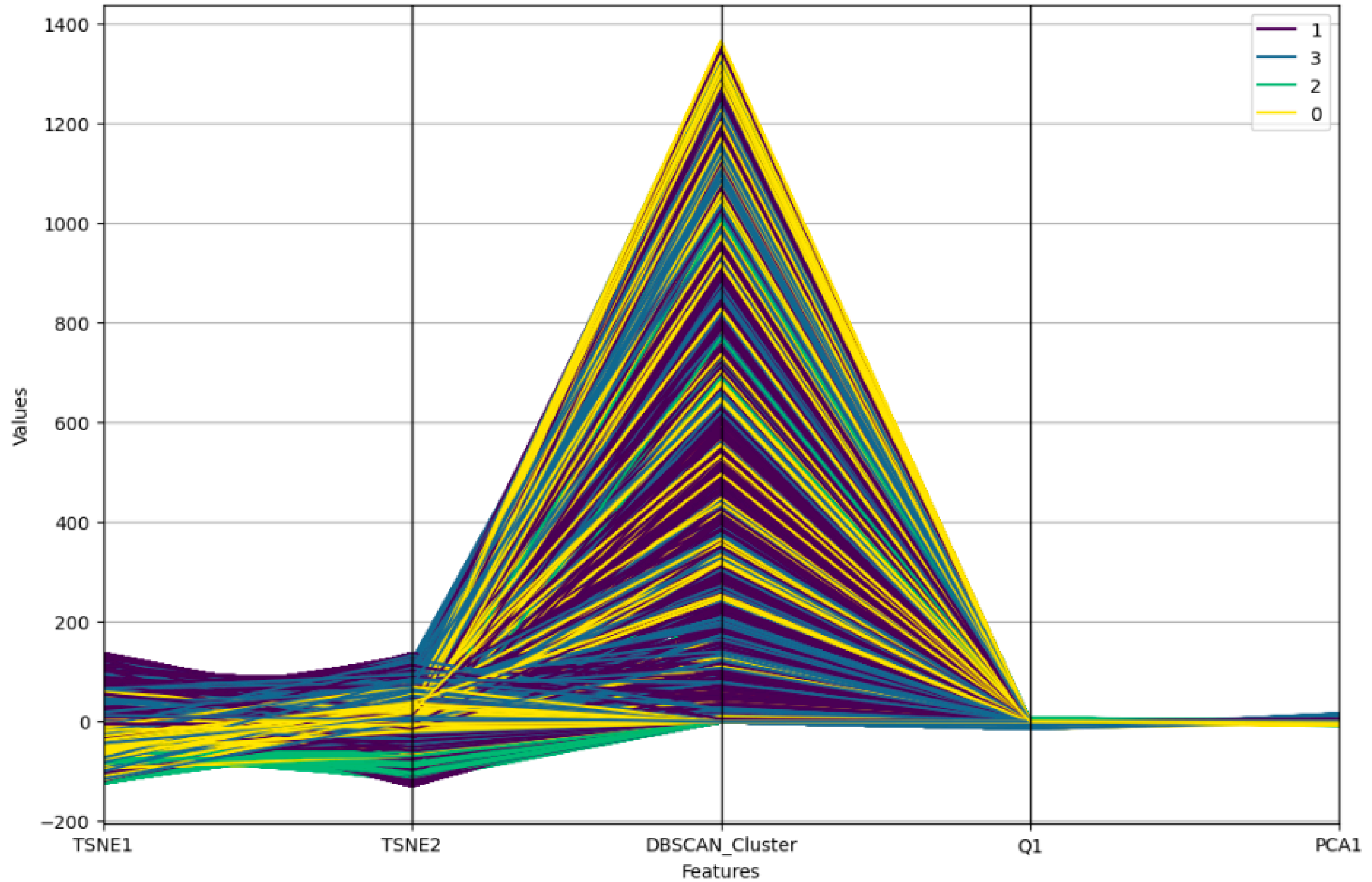
Mean Values of Top Features Across Clusters



Distribution of Individuals Across Clusters



Parallel Coordinates Plot of Top Features Across Clusters



6. Interpretation of Clusters

Cluster 0 ("Rural Poor")

- High religiosity
- Low income and education
- Mostly rural population.

Cluster 1 ("Urban Middle Class")

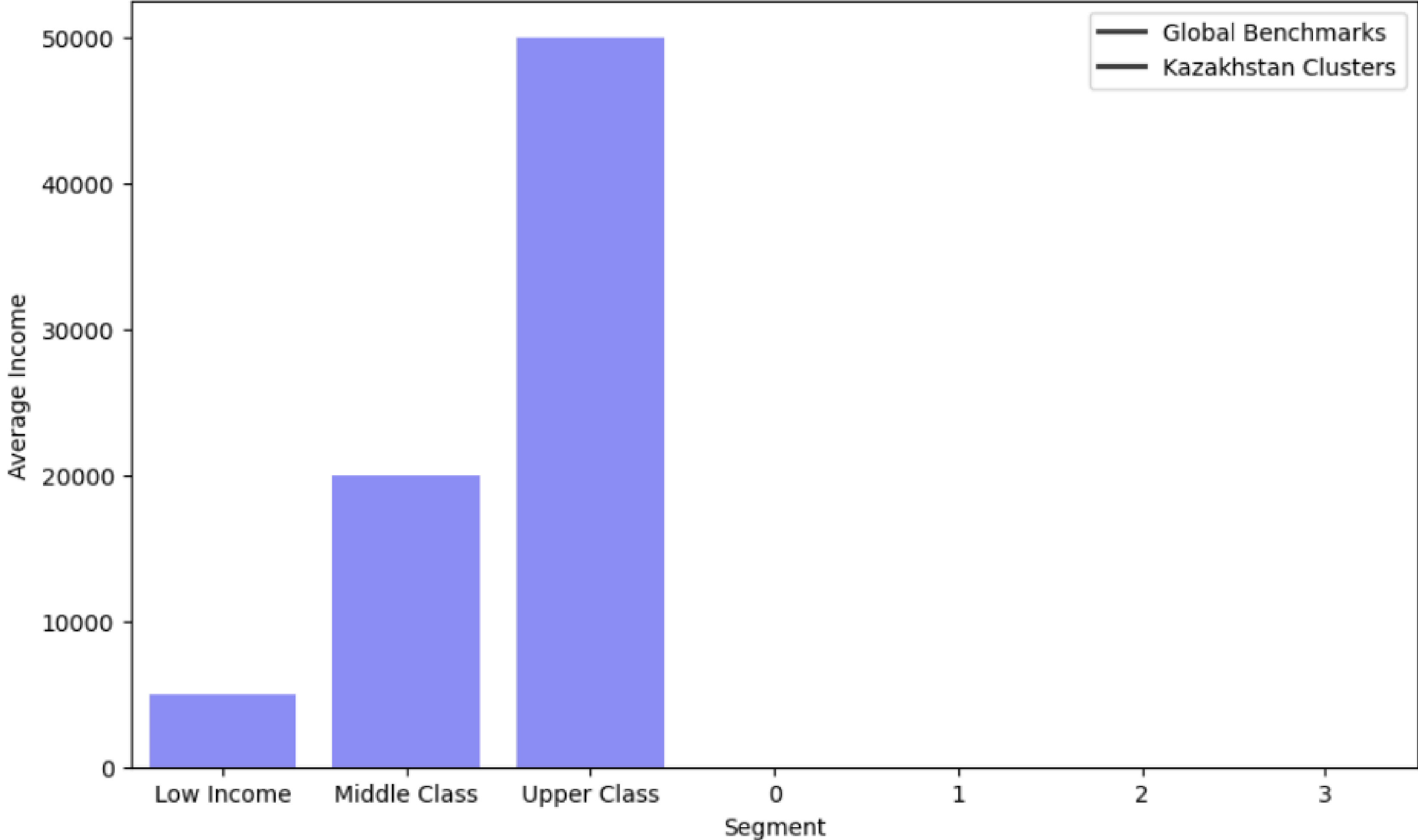
- Moderate religiosity
- Higher income and education
- Mostly urban population.

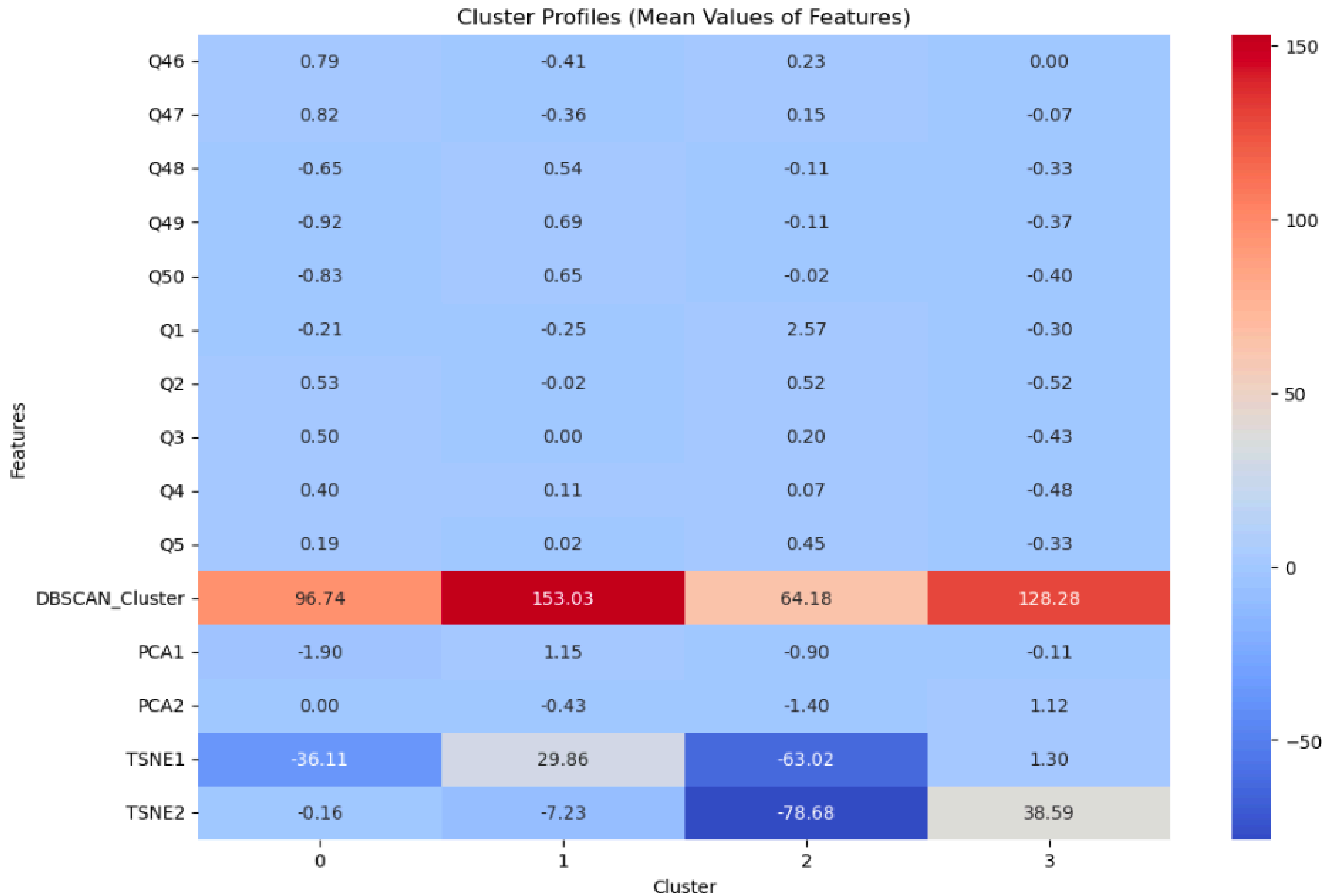
Cluster 2 ("Elite")

- High income and education
- Low religiosity
- Urban-based.

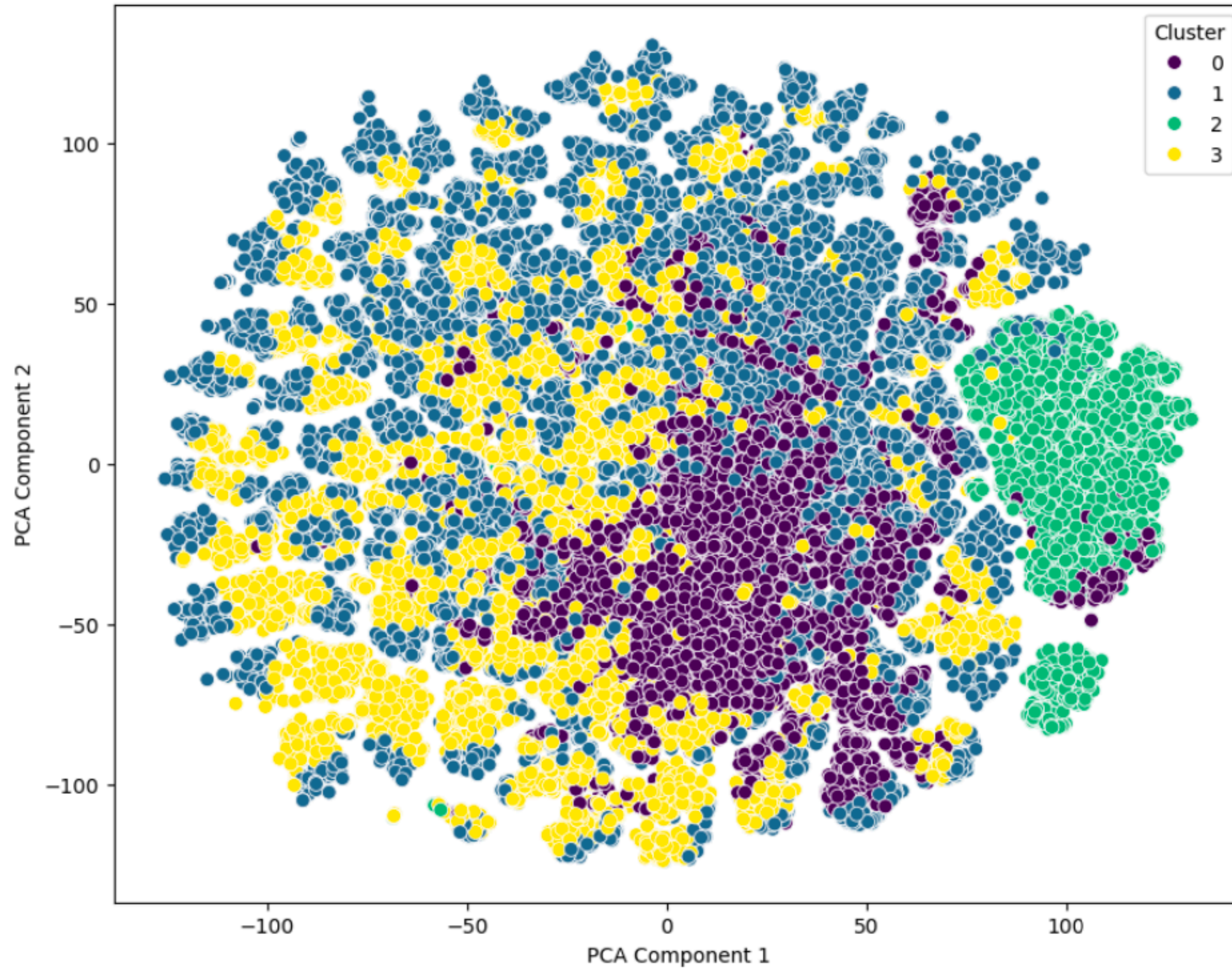
Comparison of Income Levels

- Global Benchmarks
- Kazakhstan Clusters





K-Means Clusters (PCA)



7. Saving Results

Cleaned dataset (cleaned_wvs_data.csv)

Clustered dataset (clustered_wvs_data.csv)

Cluster summary (cluster_summary.csv)

Important feature list (important_features.csv)

Conclusion

- K-Means performed better than DBSCAN, as the latter produced too many noise points.
- 4 distinct clusters were identified, with clear socio-economic differences.
- PCA and t-SNE helped visualize cluster separability.
- The results offer insights into different socio-economic groups based on values, income, and education.