

Netflix Movies and TV Shows Clustering – Project Report

Objective

The goal of this project is to apply clustering techniques on the Netflix Movies and TV Shows dataset to identify meaningful groups of content based on numerical and textual features.

Dataset Summary

Total records: 7,787

Key feature used:

- **Numerical:** release_year, duration_num, genre_count, content_age
- **Categorical:** type, rating (encoded)
- **Text:** description, listed_in (genres)

Data Preprocessing

- Missing values were handled using simple strategies (mode / derived values).
- Duration was converted into a numeric format.
- New features created:
 - **Content age** = current year – release year
 - **Genre count** = number of genres per title

Text columns were converted into numerical form using **TF-IDF Vectorization**.

Feature Engineering & Reduction

- TF-IDF was applied to `description` and `listed_in`.
- All features were combined into a single feature matrix.
- **Truncated SVD** was used to reduce dimensionality and improve clustering performance.
- Features were scaled using **StandardScaler**.

Clustering Approach

Algorithm Used

- **K-Means Clustering**
- Initialization: **k-means++**
- Number of initializations: **n_init = 10**

Choosing Number of Clusters

- **Elbow Method** showed a clear bend at **k = 3**.

- **Silhouette Score** was highest at **k = 3**.

Hence, **3 clusters** were selected.

Clustering Evaluation Metrics

Best Model (Selected)

- **Algorithm:** K-Means
- **Initialization:** k-means++
- **Number of clusters (k):** 3
- **Number of initializations:** 10

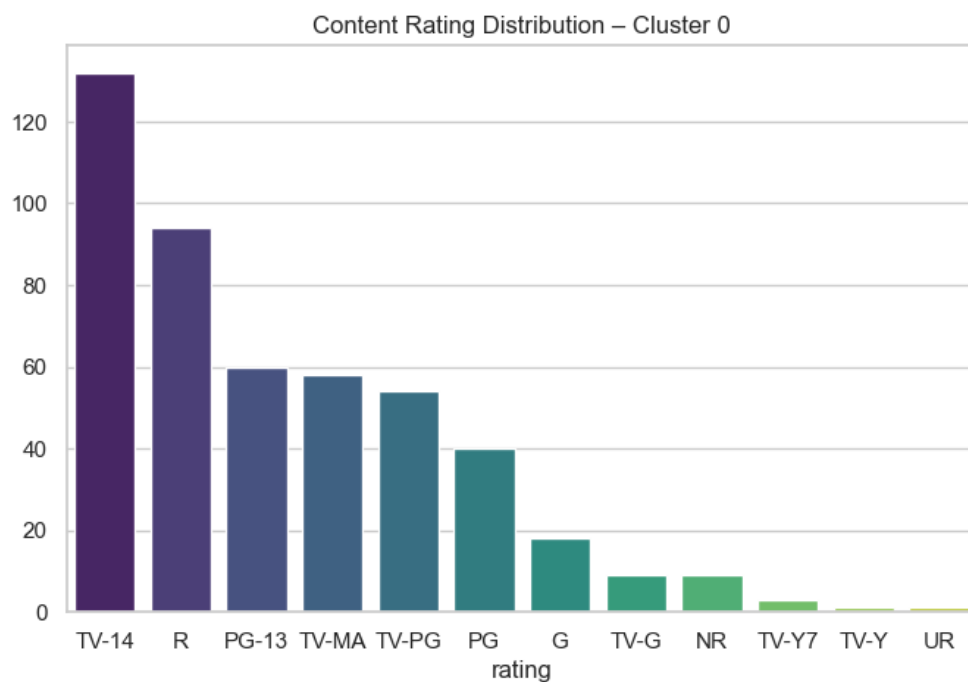
Metric Values

- **Silhouette Score: 0.36** → Good cluster separation
- **Inertia: 26,944.79** → Compact clusters without overfitting
- **Calinski–Harabasz Index: 3,663.24** → Strong cluster structure
- **Davies–Bouldin Index: 1.06** → Lower value indicates better clustering

Results

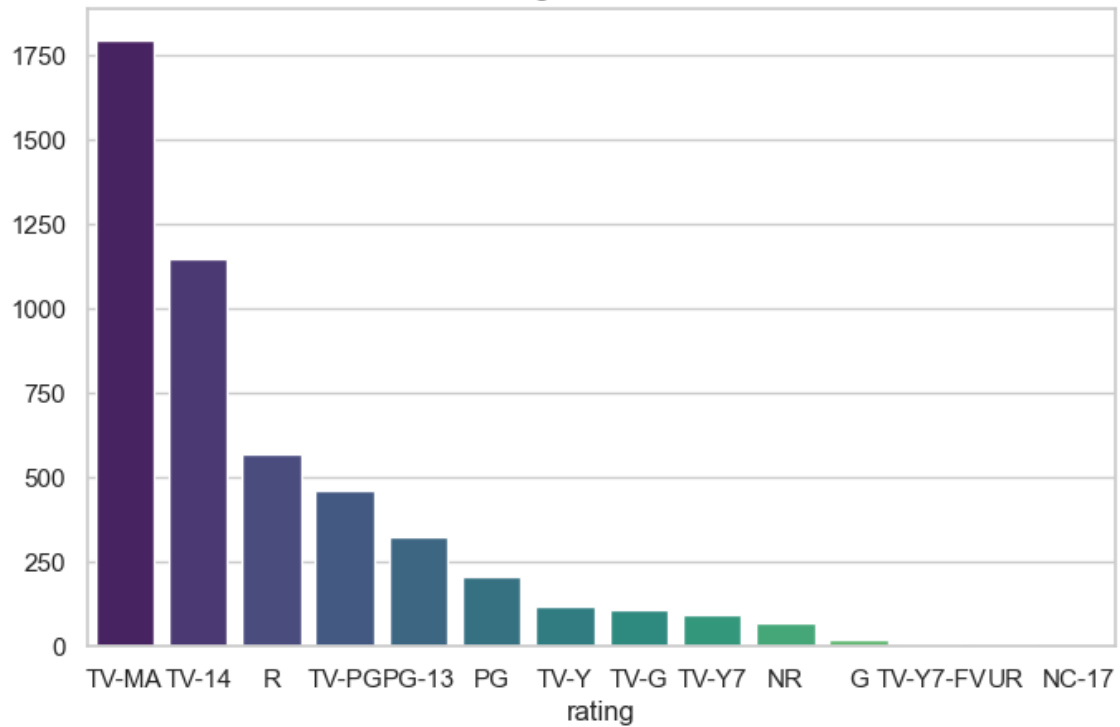
- The chosen model produced well-separated clusters.
- PCA visualization confirmed clear cluster separation.
- Evaluation metrics supported the selected configuration.

Cluster 0: Older, long-duration content



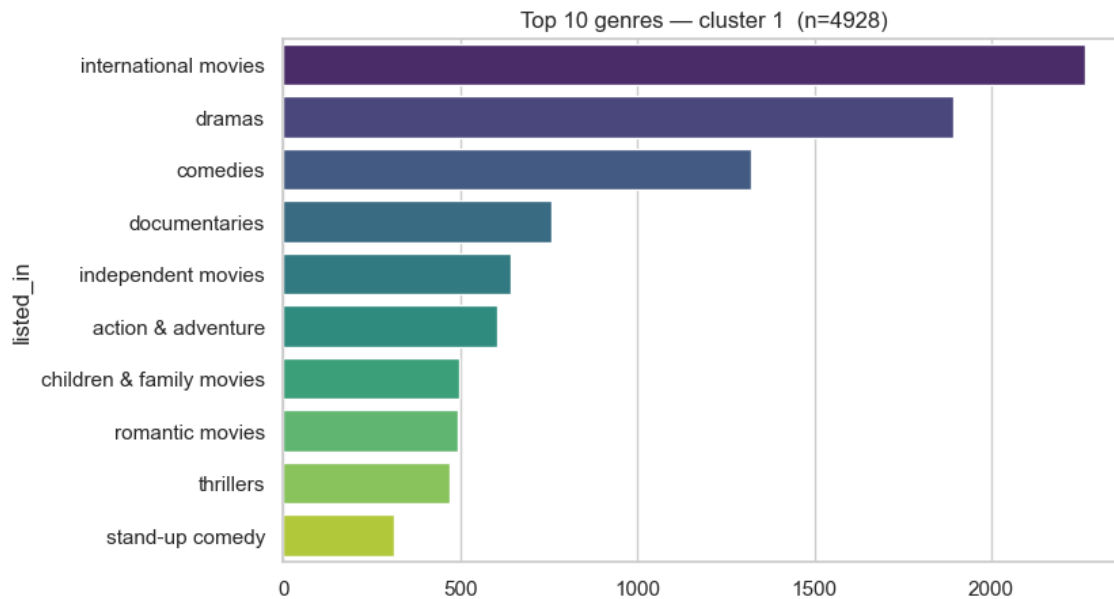
Cluster 1: Modern mainstream Netflix content

Content Rating Distribution – Cluster 1



WordCloud — cluster 1



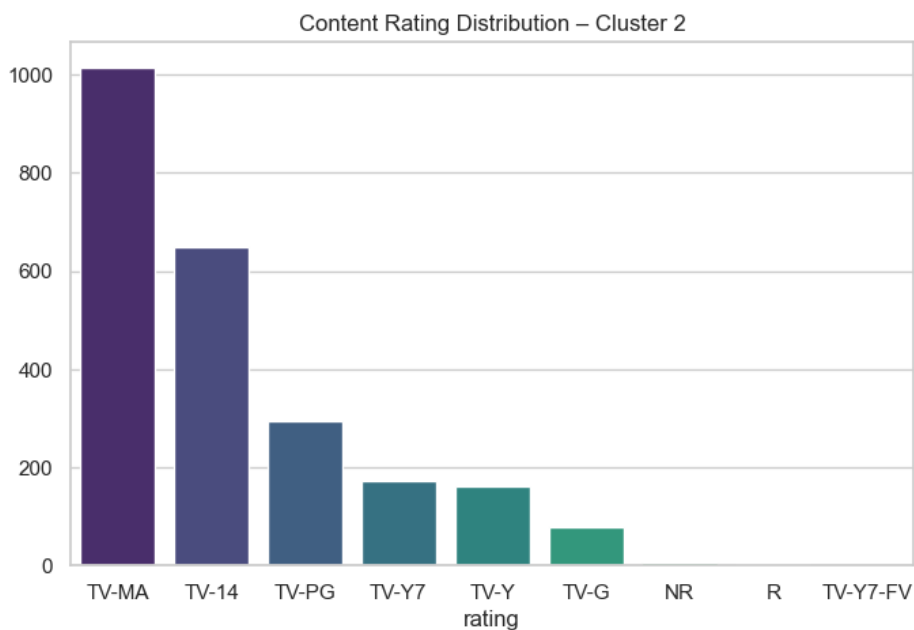


Cluster 1 represents the main Netflix catalog, dominated by international movies, dramas, and comedies.

The WordCloud shows everyday and relationship-based themes such as life, family, friends, love, and youth, indicating general audience content.

Most content is rated TV-MA and TV-14, showing this cluster targets teen and adult viewers with a wide variety of content types.

Cluster 2: Very recent content, short duration

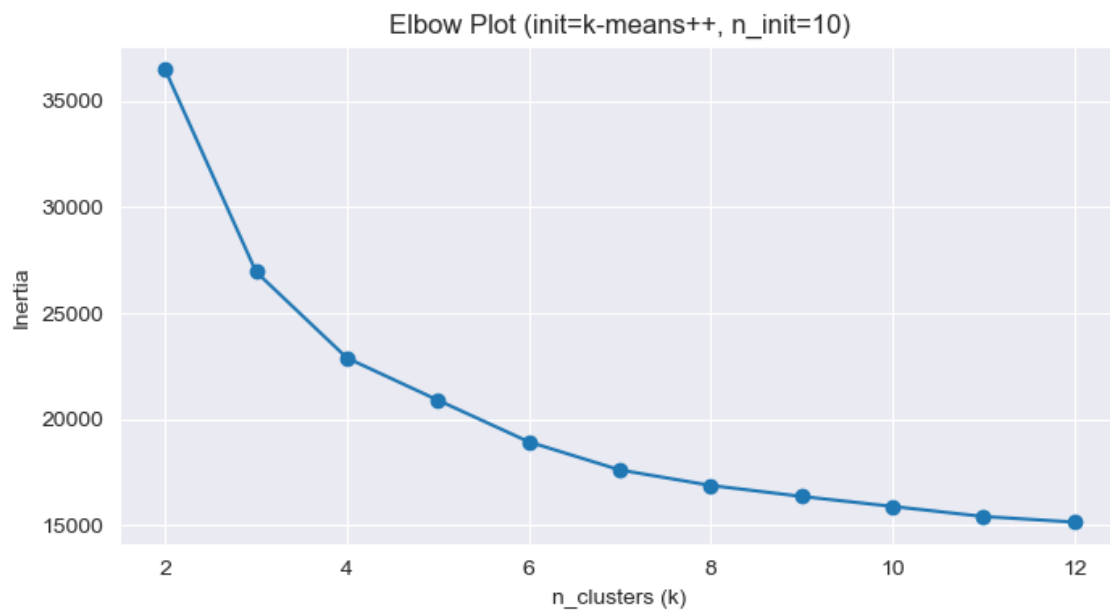


Key Observations

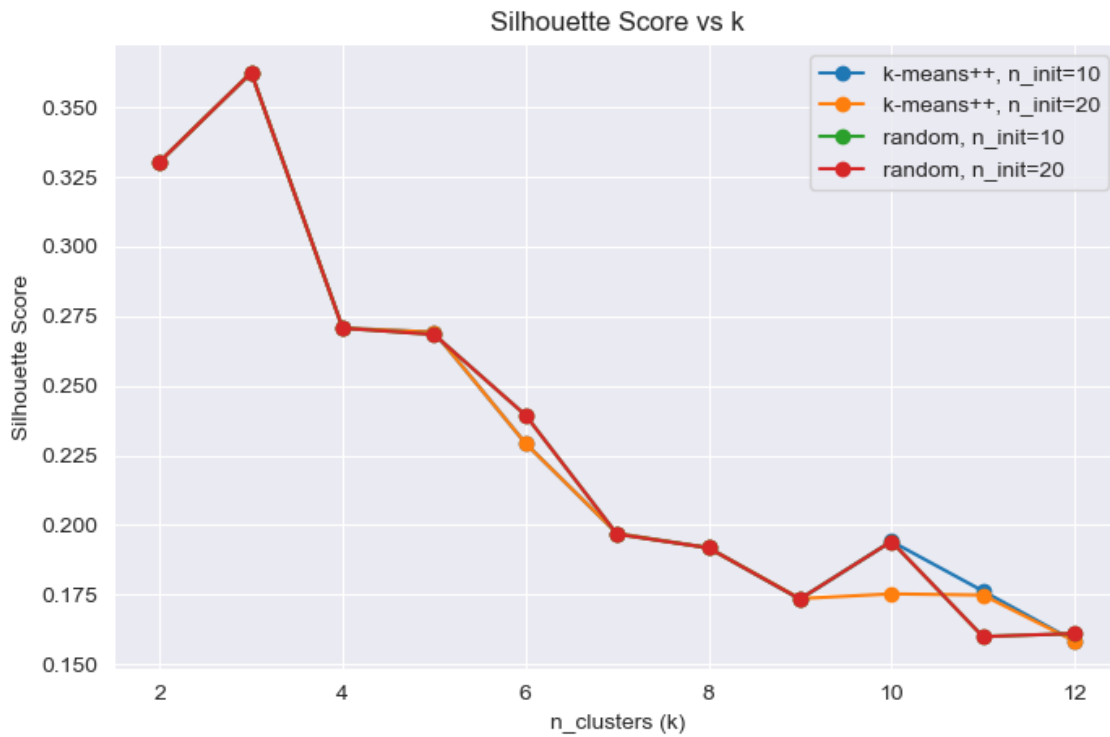
- Release year and content age strongly influence clustering.
- Duration helps separate long-form and short-form content.

Visual Analysis

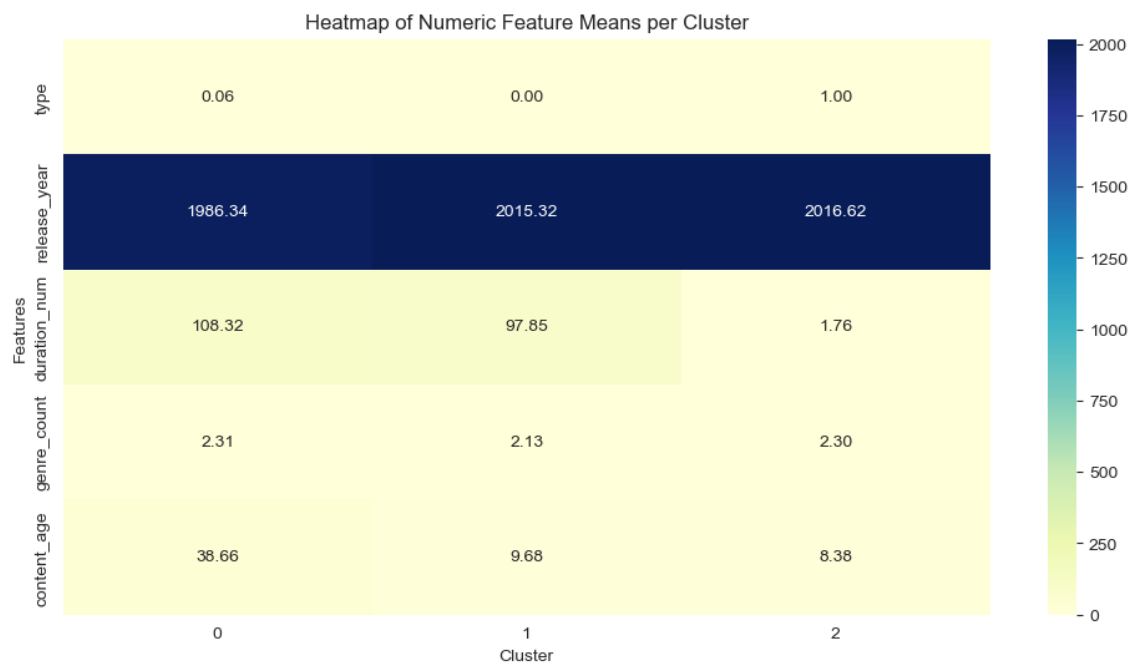
Elbow plot: Helped identify optimal k



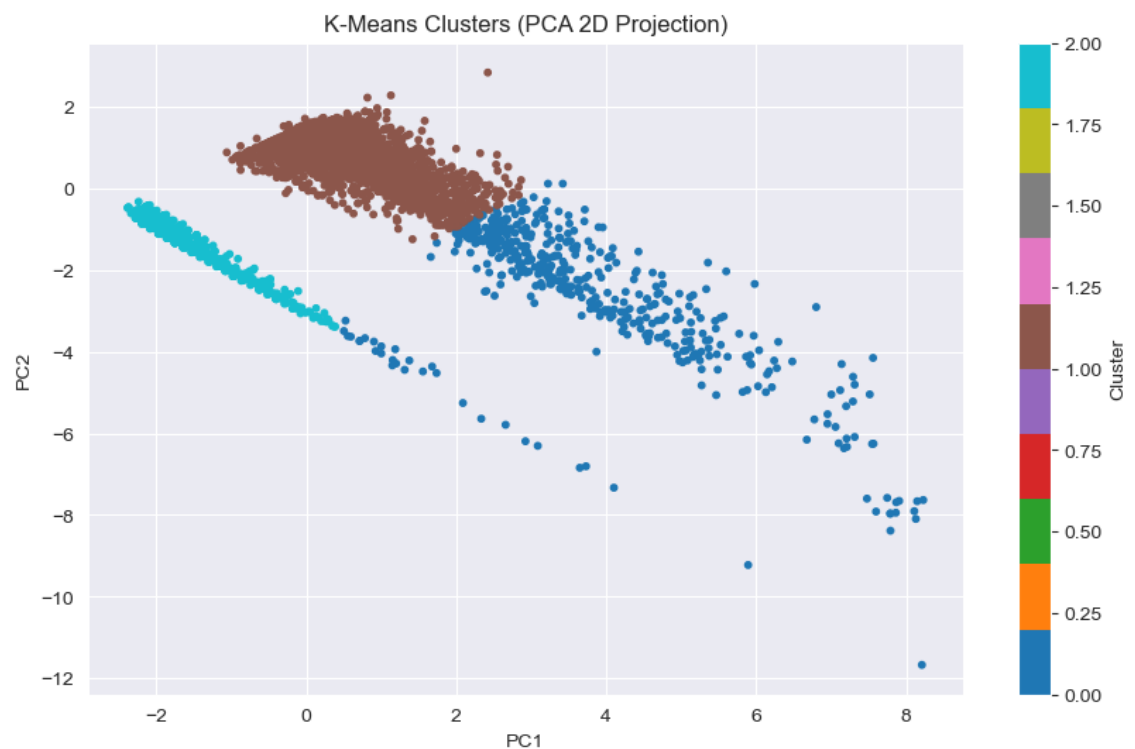
Silhouette plot: Validated cluster quality



Heatmap: Showed numerical difference across clusters



PCA scatter plots: Visualized cluster separation in 2D



Conclusion

The clustering approach successfully grouped Netflix content into three meaningful clusters, providing insights into content age, duration, and structure.