

Capstone Project

NETFLIX MOVIES AND TV SHOWS CLUSTERING

Technical Documentation

Abhishek Jain

Klearpixeloff@gmail.com

Khushboo Chaurasiya

Sharmakhushboo771@gmail.com

Table of Content:-

1. Data Description
2. Problem Statement
3. Modelling
4. Evaluation
5. Conclusion

Problem Statement:-

This dataset consists of tv shows and movies available on Netflix as of 2019. The dataset is collected from Flexible which is a third-party Netflix search engine.

In 2018, they released an interesting report which shows that the number of TV shows on Netflix has nearly tripled since 2010. The streaming service's number of movies has decreased by more than 2,000 titles since 2010, while its number of TV shows has nearly tripled. It will be interesting to explore what all other insights can be obtained from the same dataset.

Integrating this dataset with other external datasets such as IMDB ratings, rotten tomatoes can also provide many interesting findings.

In this project, I have done

- Exploratory Data Analysis
- Understanding what type content is available in different countries
- Is Netflix increasingly focusing on TV rather than movies in recent years.
- Clustering similar content by matching text-based features

Data Description:-

The dataset consists of listings of all the movies and tv shows available on Netflix, along with details such as - cast, directors, ratings, release year, duration, etc

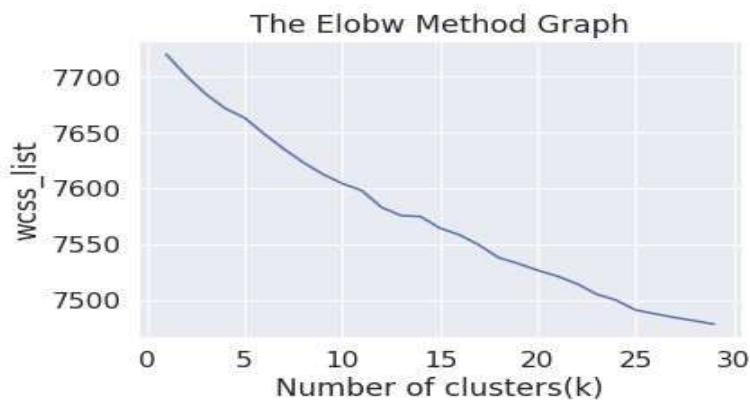
- show_id : Unique ID for every Movie / Tv Show
- type : Identifier - A Movie or TV Show
- title : Title of the Movie / Tv Show
- director : Director of the Movie
- cast : Actors involved in the movie / show
- country : Country where the movie / show was produced
- date_added : Date it was added on Netflix
- release_year : Actual Release Year of the movie / show
- rating : TV Rating of the movie / show
- duration : Total Duration - in minutes or number of seasons
- listed_in : Genre
- description: The Summary descript

Modelling :-

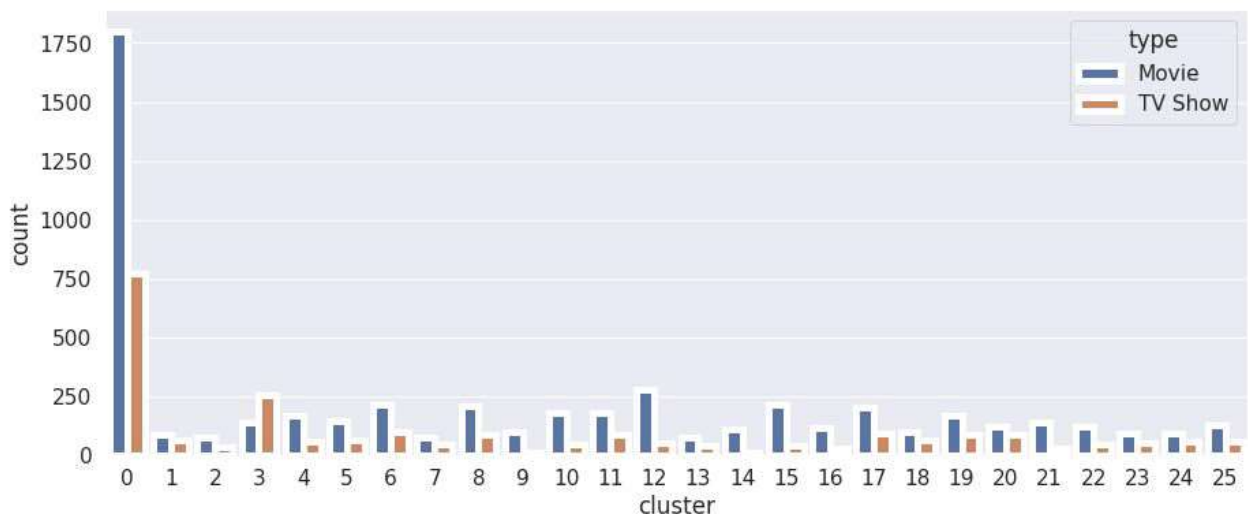
K-MEANS:

K-Means Clustering is an Unsupervised Learning algorithm which groups the unlabeled dataset into different clusters. Here K defines the number of predefined clusters that need to be created in the process, as if $K=2$, there will be two clusters, and for $K=3$, there will be three clusters, and so on.

Elbow Method



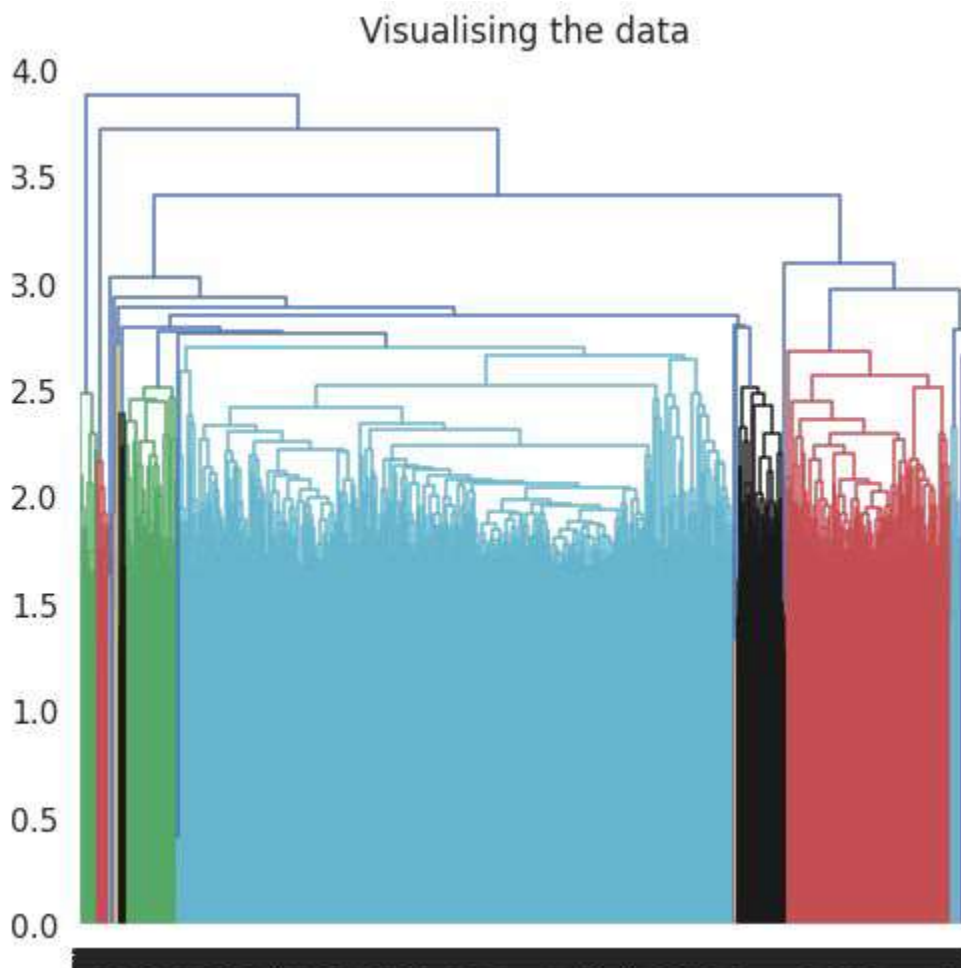
- From elbow method generating 26 clusters



- cluster 0 has the highest number of datapoints and evenly distributed for other cluster

Agglomerative Clustering

- In agglomerative clustering no need to give the value of k beforehand
- The agglomerative hierarchical clustering algorithm is a popular example of HCA
- Here I used ward linkage



- the optimal number of clusters is 4 using the Dendrogram

Evaluation

The Davies-Bouldin index (DBI).It is most commonly used to evaluate the goodness of split by a K-Means clustering algorithm for a given number of clusters.

Davies_bouldin_score is 9.05605194948868(K-Means)

Davies_bouldin_score is 13.979382092977453 (Agglomerative Clustering)

Conclusion

- from elbow and silhouette score ,optimal of 26 clusters formed , K Means is best for identification than Hierarchical as the evaluation metrics also indicates the same.in kmean cluster 0 has the highest number of datapoints and evenly distributed for other cluster
- Netflix has 5372 movies and 2398 TV shows, there are more number movies on Netflix than TV shows.
- TV-MA has the highest number of ratings for tv shows i.e adult ratings
- highest number of movies released in 2017 and 2018 highest number of movies released in 2020 The number of movies on Netflix is growing significantly faster than the number of TV shows. We saw a huge increase in the number of movies and television episodes after 2015. there is a significant drop in the number of movies and television episodes produced after 2020. It appears that Netflix has focused more attention on increasing Movie content than TV Shows. Movies have increased much more dramatically than TV shows
- the most content is added to Netflix from october to january
- Documentaries are the top most genre in netflix which is followed by standup comedy and Drams and international movies

- kids tv is the top most TV show genre in netflix
- most of the movies have duration of between 50 to 150 highest number of tv_shows consistig of single season Those movies that have a rating of NC-17 have the longest average duration. When it comes to movies having a TV-Y rating, they have the shortest runtime on average
- usa has the highest number of content on the netflix ,followed by india