

Capstone Project-4

NETFLIX MOVIES & TV SHOWS CLUSTERING

Team Member

Anupam Mishra

Kartike

Animesh Chakraborty

CONTENT

- Introduction
- Problem Statement
- Data Summary
- Data Cleaning
- EDA
- Text Preprocessing
- Unsupervised ML
 - K-Means Clustering
 - Agglomerative Clustering
- Conclusion

Introduction

- Netflix is a prominent OTT platform with a wide variety of content to view from a variety of nations and genres, so keep an eye on it. This dataset consists of tv shows and movies available on Netflix as of 2019. The dataset is collected from Flixable which is a third-party Netflix search engine.
- The idea of this project is to analyze and perform clustering to determine various patterns related to the content available in Netflix. Based on the attributes related to the Tv shows or movies, we will be implementing different clustering algorithms which comes under unsupervised Machine learning category.

PROBLEM STATEMENT

- 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.
- In this project, you are required to do :
 1. Exploratory Data Analysis
 2. Understanding what type of content is available in different countries
 3. Is Netflix has increase its focus on TV rather than movies in recent years.
 4. Clustering similar content by matching text-based features

DATA SUMMARY

The dataset has 7787 rows and 12 columns.

- **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 description

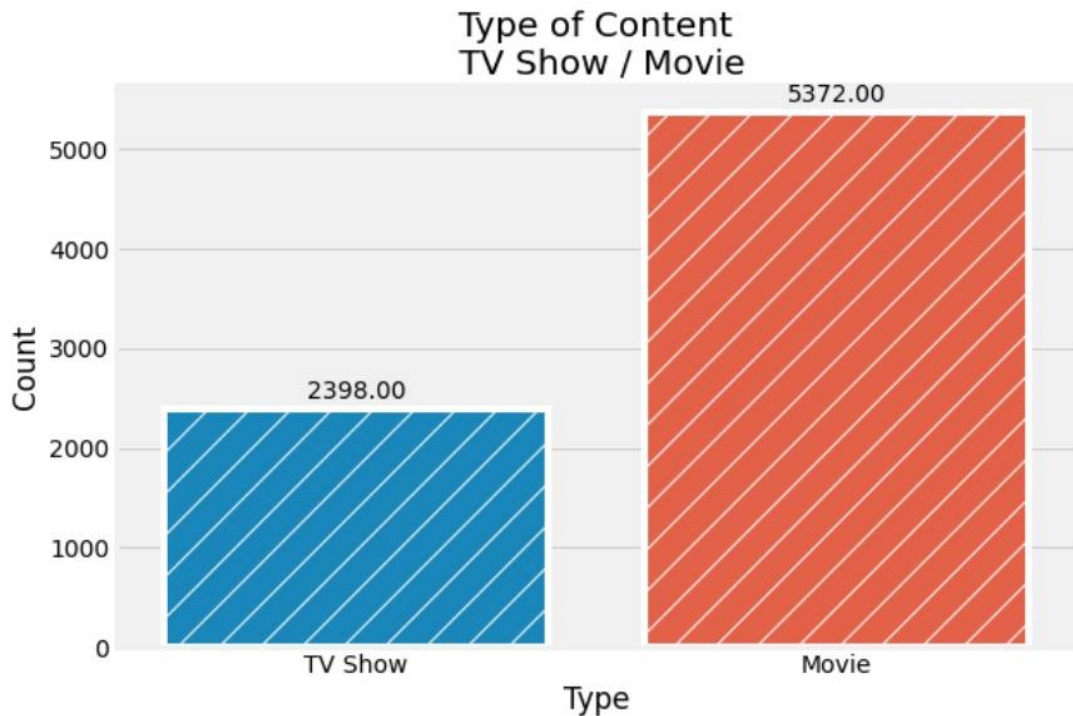
Data Cleaning

Null Value Treatment:

- **Director** feature have more than **30.68%** of null values. Filling null values by 'unknown'.
- **Country** feature have **6.51%** of null values. Filling null values by mode of feature.
- **Cast** feature have **9.22%** of null values. Filling null values by 'unknown'.
- **Rating** feature have **0.09%** of null values. Dropping rows corresponding to null values.
- **Date_added** feature have **0.13%** of null values. Dropping rows corresponding to null values.

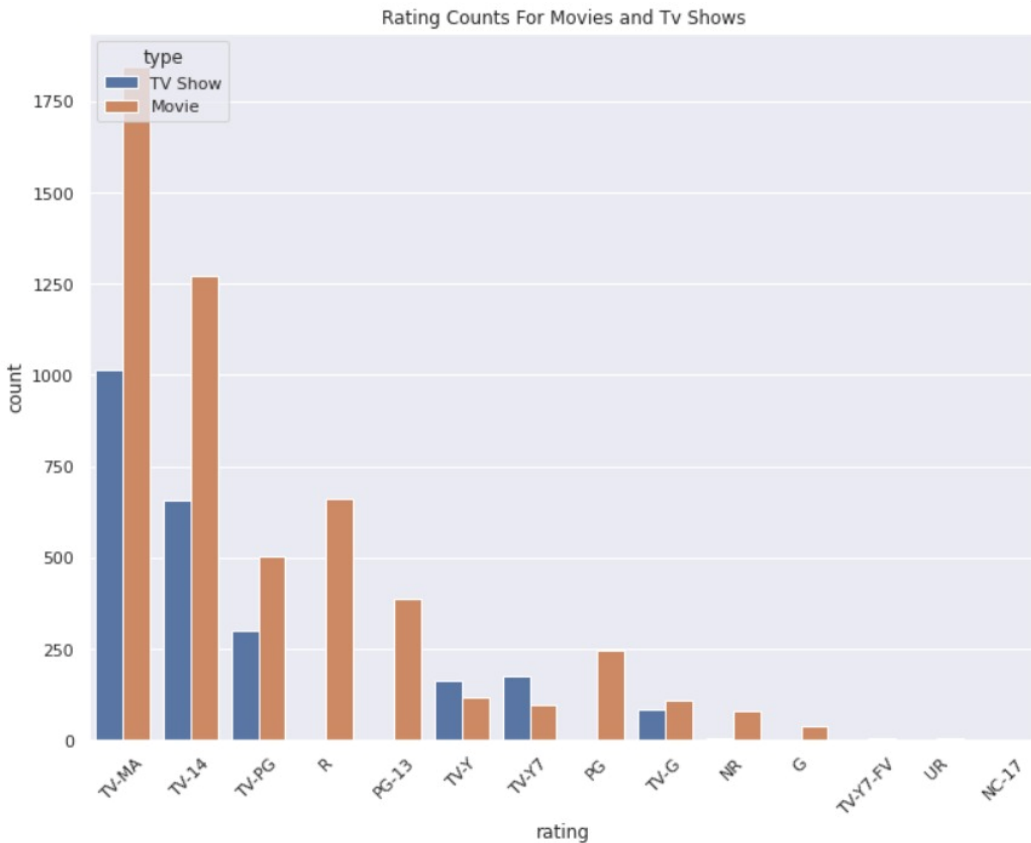
Exploratory Data Analysis(EDA)

Type of content available on Netflix



- It is evident that there are more movies on Netflix than TV shows.
- Netflix has 5372 movies, which is more than double the quantity of TV shows.

Movie ratings analysis

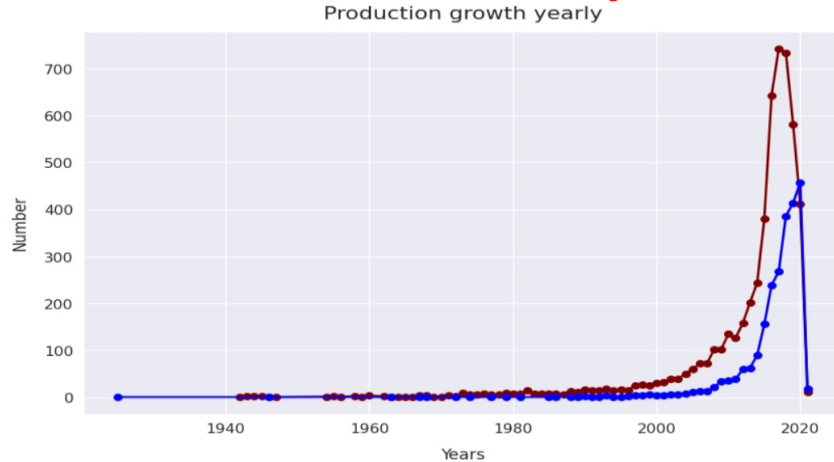


The 'TV-MA' rating is used in the majority of the film. The TV Parental Guidelines provide a "TV-MA" classification to a television programme that is intended solely for mature audiences.

The second largest is 'TV-14,' which stands for content that may be inappropriate for minors under the age of 14.

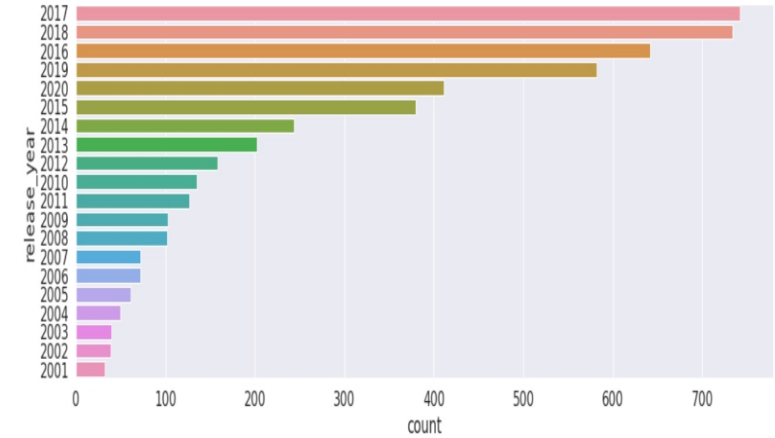
The third most common is the extremely popular 'R' rating. The Motion Picture Association of America defines an R-rated film as one that contains material that may be inappropriate for children under the age of 17; the MPAA states that "Under 17 requires accompanying parent or adult guardian."

Growth in content over the years

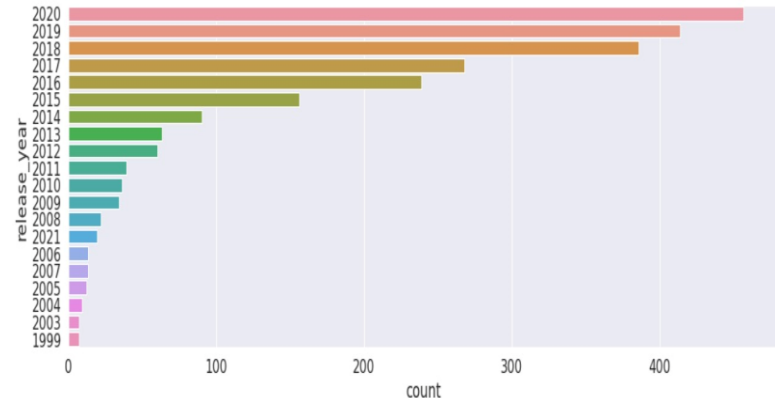


- * We saw a huge increase in the number of movies and TV Shows after 2015.
- * Highest number of Movies released in 2020 and 2019.
- * Highest number of TV Shows released in 2017 and 2018.
- * The number of movies on Netflix is growing significantly faster than the number of TV shows.
- * It appears that Netflix has focused more attention on increasing Movie content than TV Shows. Movies have increased much more dramatically than TV shows.

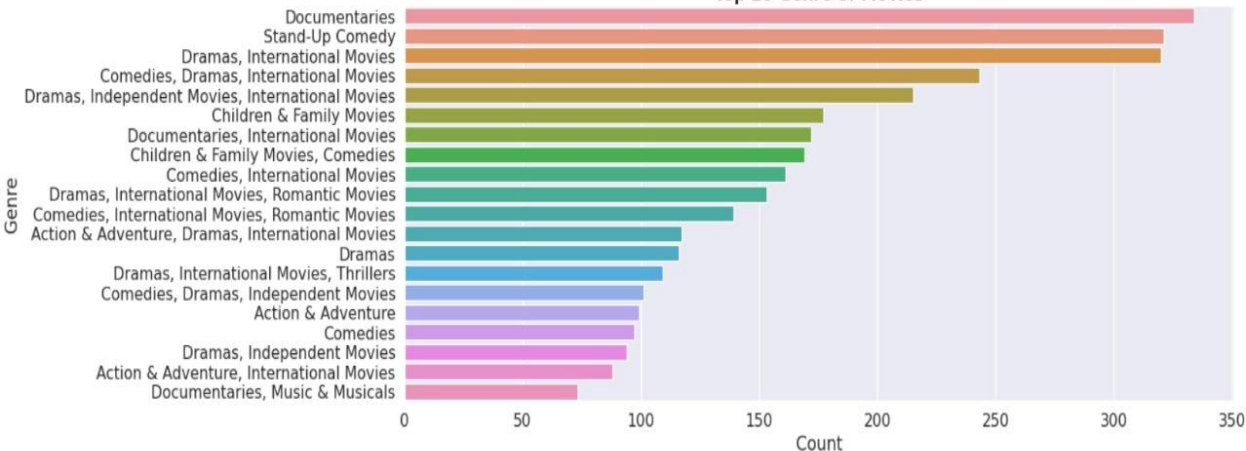
tv_shows



movies

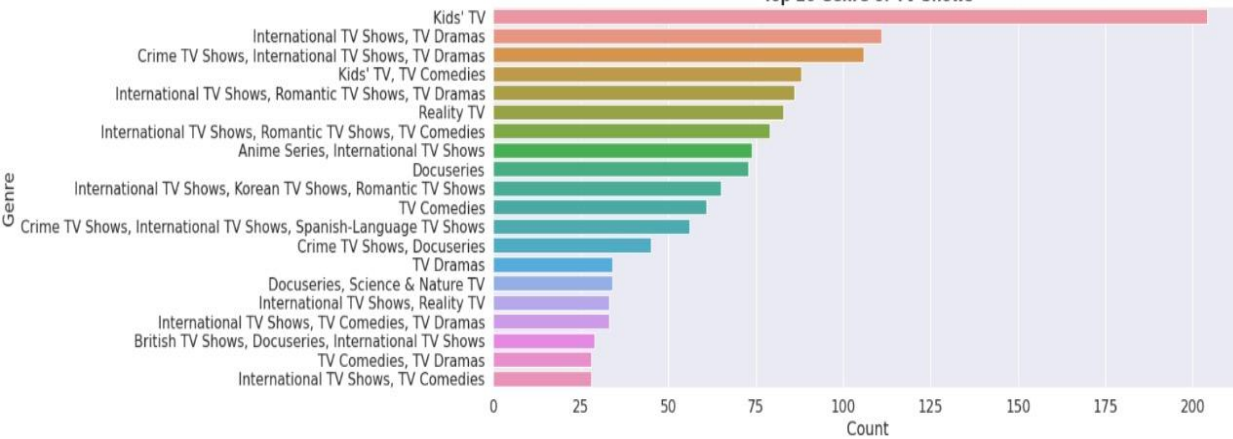


Top 20 Genre of Movies



• "Documentaries" are the top most movies genre in Netflix which is followed by "Stand-Up comedy" and "Dramas, International movies".

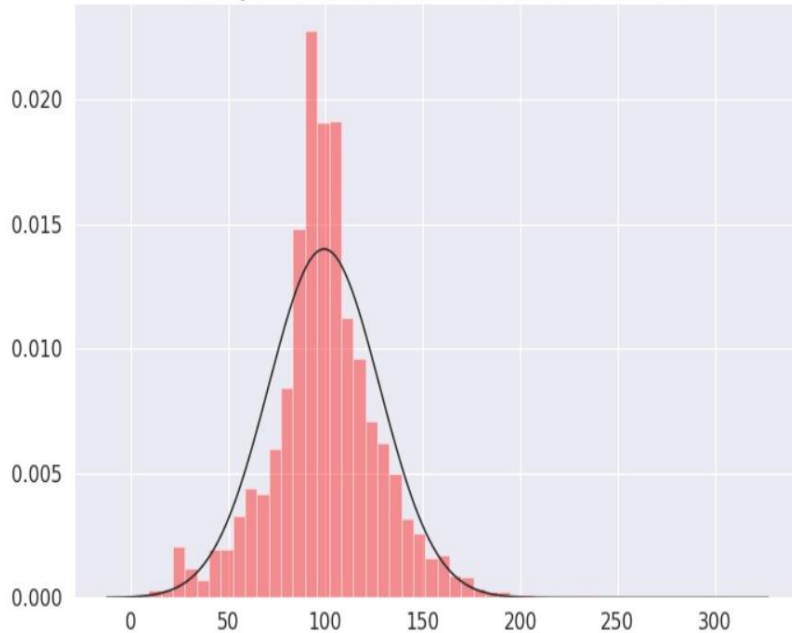
Top 20 Genre of TV Shows



• "Kids' TV" are the top most TV Shows genre in Netflix which is followed by "International TV Shows, TV Dramas" and "Crime TV Shows, International TV Shows, TV Dramas".

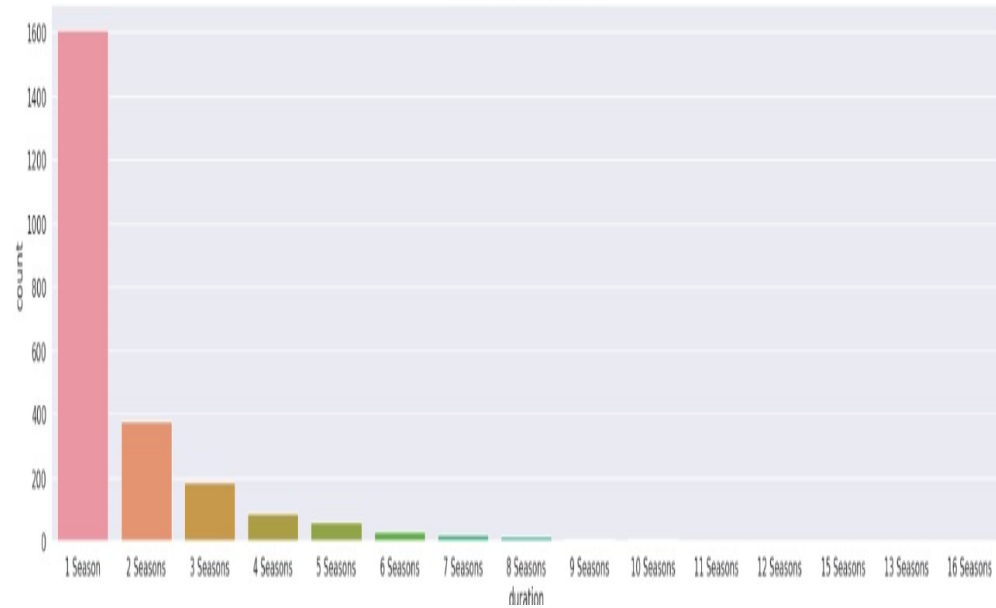
Duration

Distplot with Normal distribution for Movies



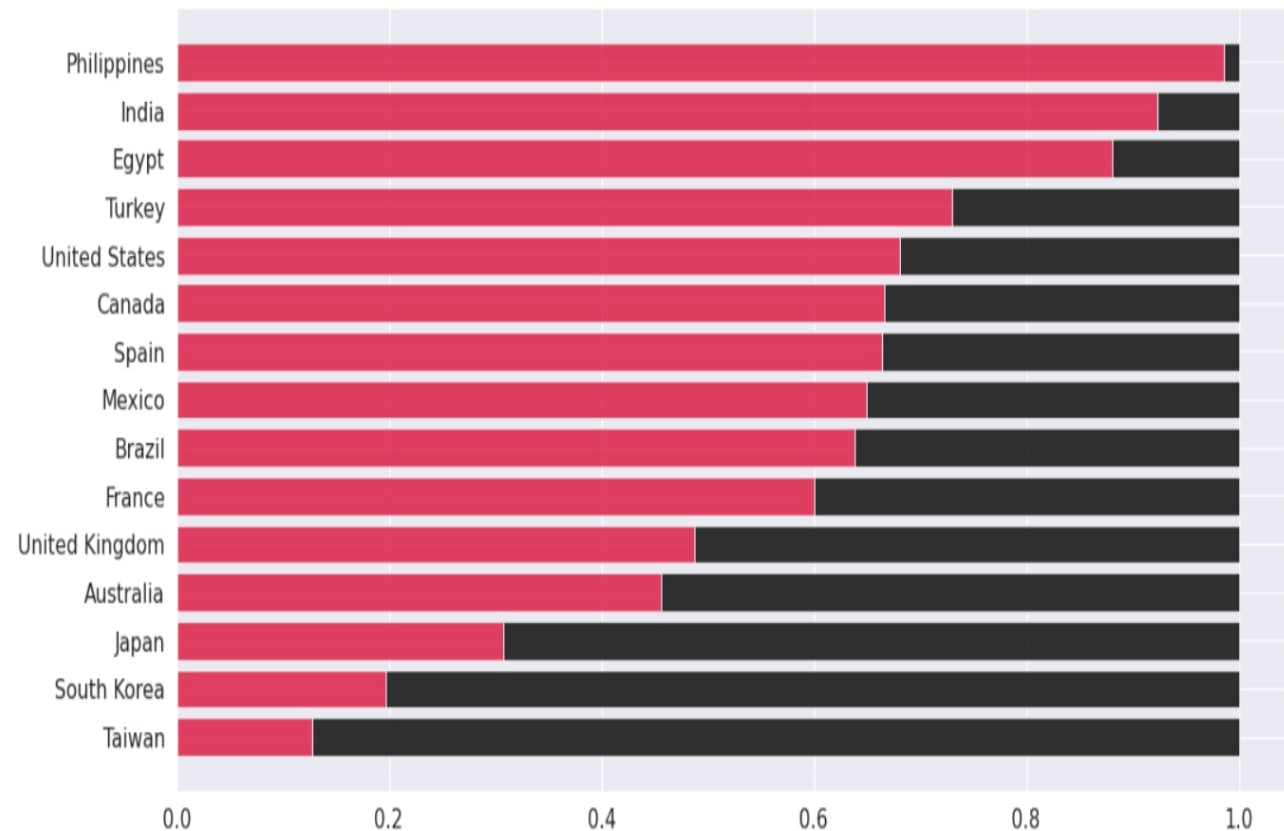
- most of the movies have duration of between 50 to 150

Distribution of TV Shows duration



- highest number of tv_shows consist of single season

How does content differ by country in the top fifteen lists

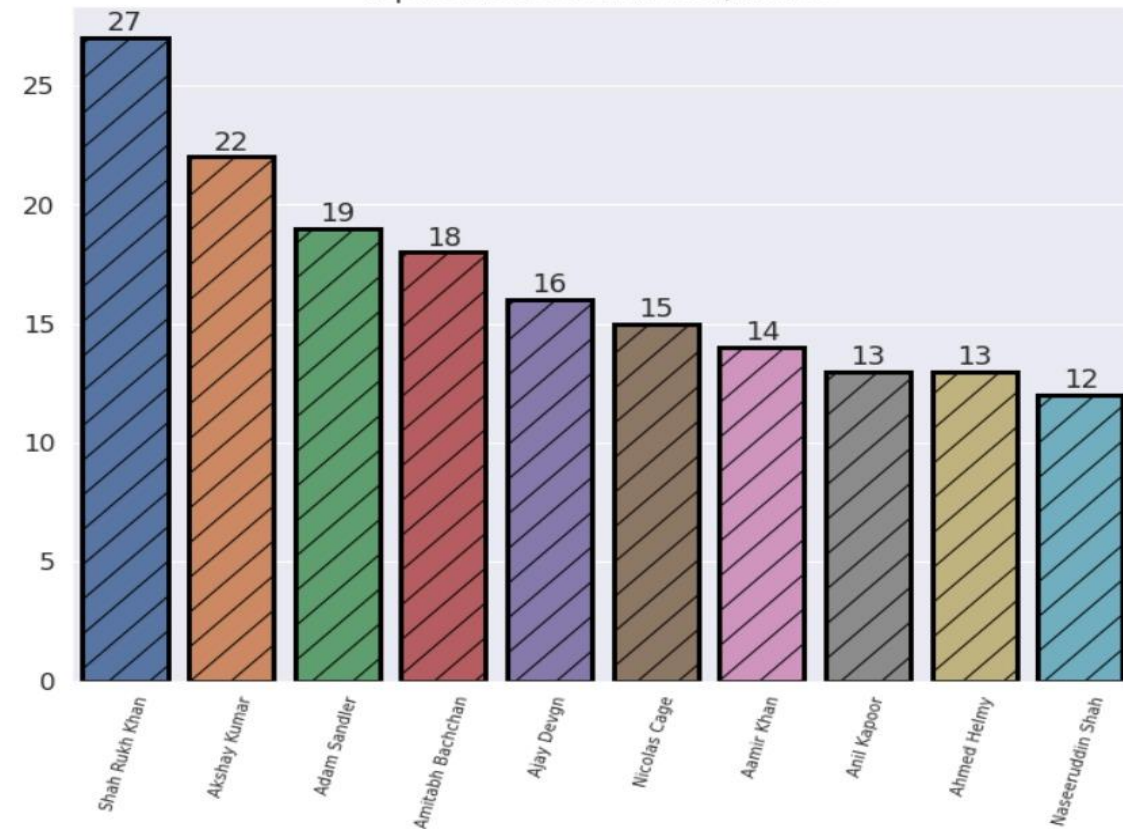


- Philippines has highest number of Movies percentage compared to TV Shows in Netflix.

- Taiwan has highest number of TV Shows percentage compared to Movies in Netflix.

Top 10 Actors who appear in the majority of films

Top 10 Movies Lead Actor/Actress



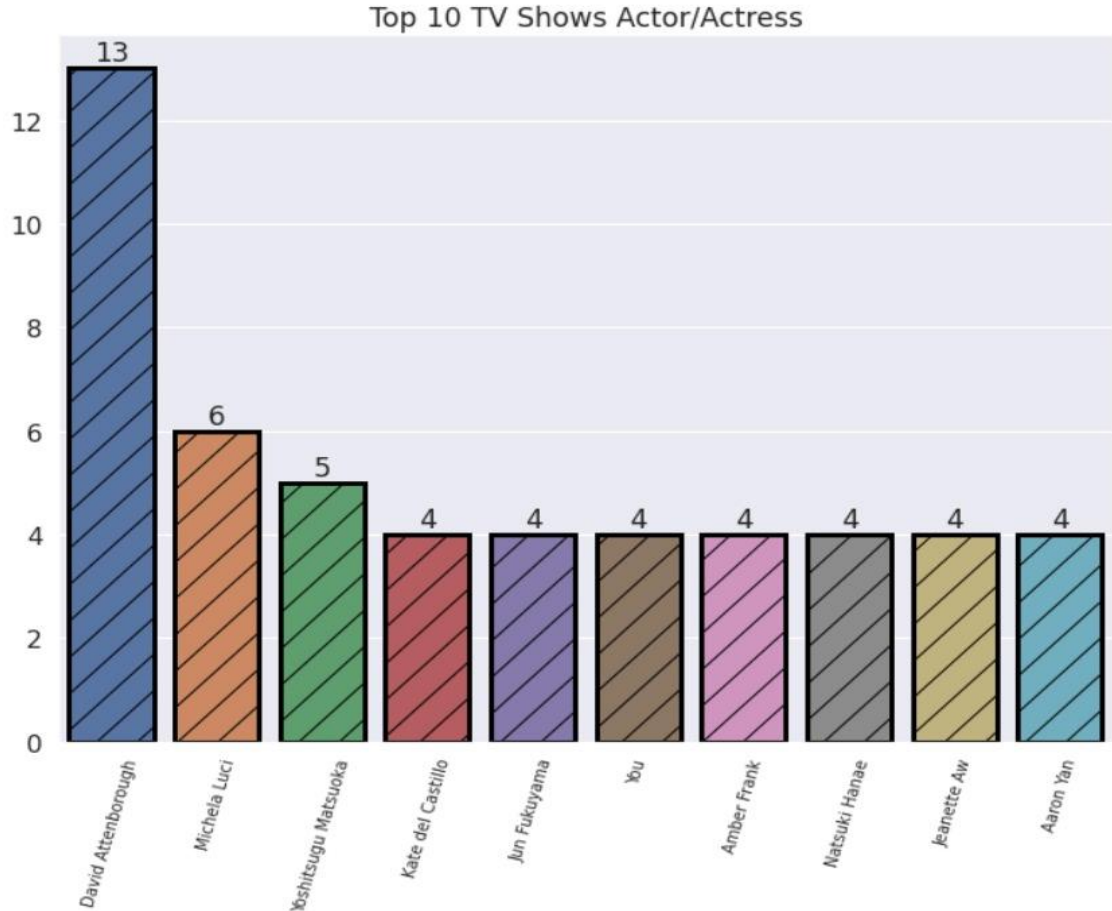
* According to the above barplot, Shah Rukh Khan has worked in over 27 (Lead Actor) films.

* After Shah Rukh Khan, Akshay Kumar (Lead Actor) is ranked second, with 22 films under his belt.

Top 10 Actors who appear in the majority of TV Shows

* According to the above barplot, David Attenborough has worked in over 23 (Lead Actor) TV Shows.

* After David Attenborough, Michela Luci (Lead Actress) is ranked second, with 6 TV Shows under his belt.



Netflix Content for different age groups in top 15 countries

* It is also interesting to see parallels between culturally comparable nations -

the US and UK are closely aligned with their Netflix target ages, but radically different from, example, India or Japan!

* Also, Mexico and Spain have similar content on Netflix for different age groups.



Text Pre-processing for Clustering

1. Removing Punctuation:

- Punctuations does not carry any meaning in clustering.
- So, removing punctuations helps to get rid of unhelpful parts of the data, or noise.

2. Removing Stopwords:

- Stopwords are basically a set of commonly used words in any language, not just in English.
- If we remove the words that are very commonly used in a given language, we can focus on the important words instead.

3. Stemming:

1. • Stemming is the process of removing a part of a word, or reducing a word to its stem or root.
- Applying stemming to reduce words to their basic form or stem, which may or may not be a legitimate word in the language.

1.K-Means Clustering

K-means algorithm is an iterative algorithm that tries to partition the dataset into K pre defined distinct non overlapping subgroups where each data point belongs to only one group

1.. Vectorization:

- Here we have textual data
- Clustering algorithms cannot understand textual data
- So, we use vectorization technique to convert textual data to numerical vectors.

So, we use vectorization technique to convert textual data to numerical vectors.

2. Elbow Curve:

- The Elbow Curve is one of the most popular methods to determine this optimal value of k .
- The elbow curve uses the sum of squared distance (SSE) to choose an ideal value of k based on the distance between the data points and their assigned clusters.

3. Silhouette score :

- Silhouette score is used to evaluate the quality of clusters created
- using clustering algorithms such as K Means in terms of how well
- samples are clustered with other samples that are similar to each
- other.

1. Silhouette Score

Silhouette Coefficient Formula

$$S = \frac{(b-a)}{\max(a,b)}$$

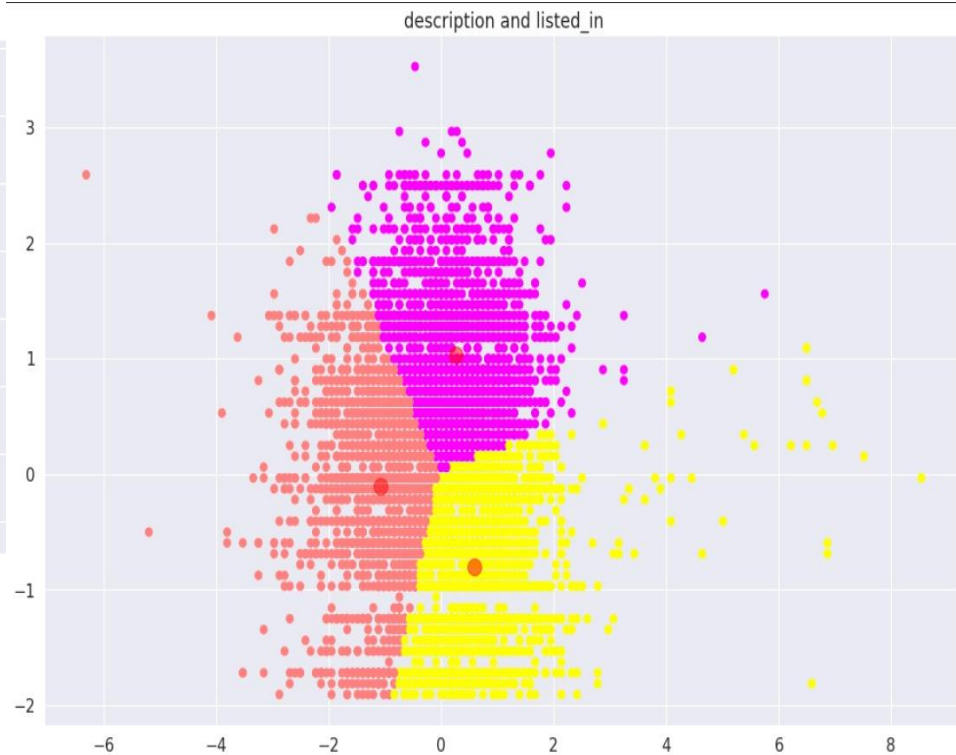
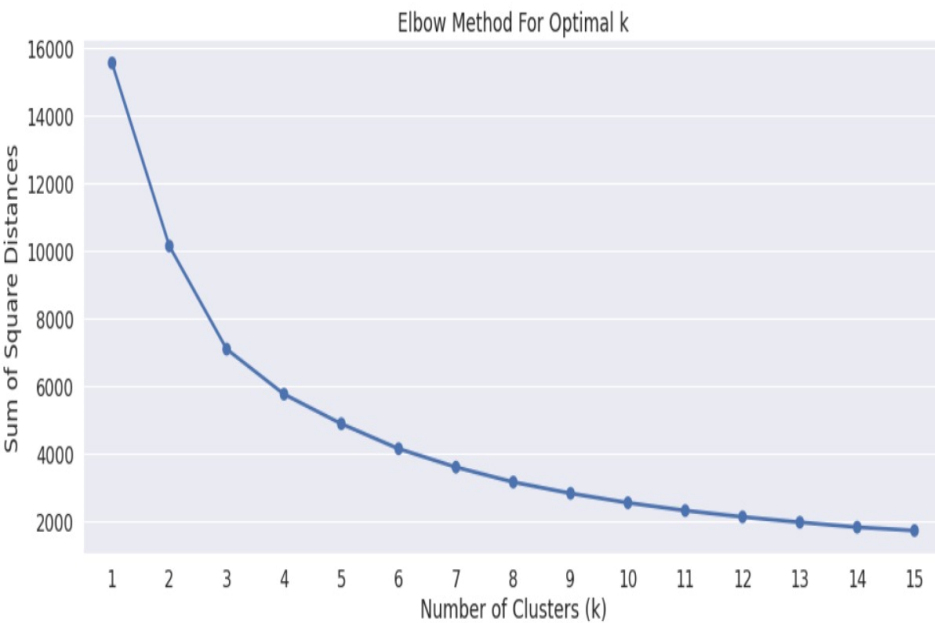
- **mean intra-cluster distance (a)** :- Mean distance between the observation and all other data points in the same cluster.
- **mean nearest-cluster distance (b)** :- Mean distance between the observation and all other data points of the next nearest cluster. This distance can also be called b.

The value of the silhouette coefficient is between [-1, 1]

- If score is **1** denotes the **best** meaning that the data point is very compact within the cluster to which it belongs and far away from the other clusters.
- The worst value is -1
- If score is 0 denotes overlapping clusters

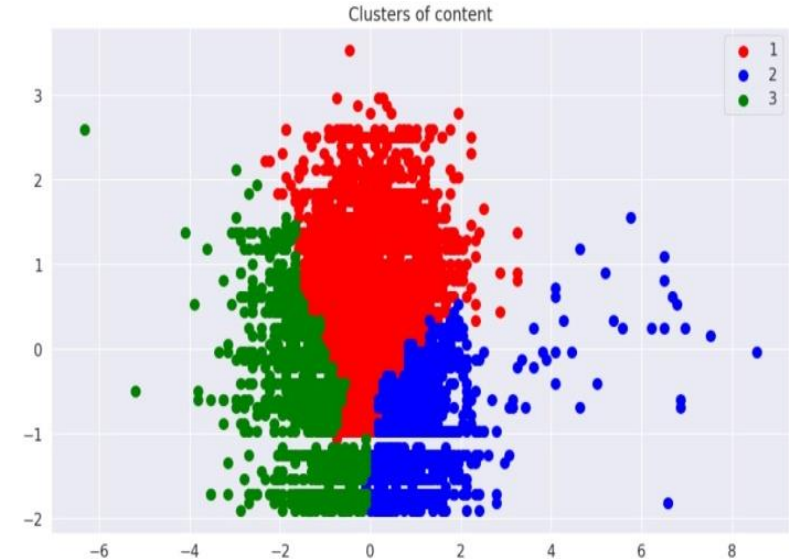
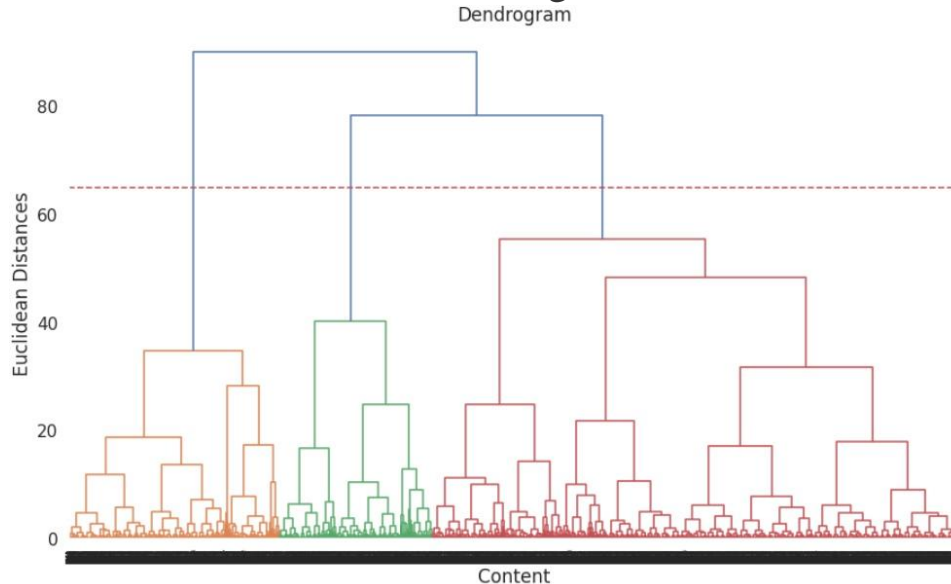
	n clusters	silhouette score
1	3	0.348
0	2	0.337
12	14	0.332
5	7	0.330
11	13	0.329
10	12	0.328
13	15	0.326
9	11	0.324
8	10	0.323
7	9	0.322
2	4	0.320
4	6	0.320
6	8	0.316
3	5	0.308

2. Elbow Method



2. 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 3 using the Dendrogram

Conclusion

- Data set contains 7787 rows and 12 columns in that cast and director features contains large number of missing values so we can drop it and we have 10 features for the further implementation
- We have two types of content TV shows and Movies (30.86% contains TV shows and 69.14% contains Movies)
- By analysing the content added over years we get to know that in recent years netflix is focusing on movies than TV shows (movies is increased by 80% and TV shows is increased by 73% compare to 2016 data)
- The most number of the movies and TV shows release in 2017 and 2020 respectively and United States have the maximum content on netflix

- On Netflix, Dramas genre contains the maximum content among all of the genres and the most of the content added in december month and less content in February.
- By applying the silhouette score method for n range clusters on dataset we got best score which is 0.348 for 3 clusters it means content explained well on their own clusters, by using elbow method after $k = 3$ curve gets linear it means $k = 3$ will be the best cluster.
- Applied different clustering models Kmeans, hierarchical, Agglomerative clustering on data we got the best cluster arrangements.
- By applying different clustering algorithms to our dataset ,we get the optimal number of cluster is equal to 3.

❖ Challenges

- ❑ A huge amount of data needed to be deal while doing the project which is quite an important task and also even small inferences need to be kept in mind.
- ❑ As dataset was quite big enough which led more computation time.
- ❑ Handling the numerical and categorical data to build high accuracy model.



Thank you