

Capstone Project

- Netflix Movies and TV Shows Clustering
- Unsupervised Machine Learning

By

Ankur

Barplot based on release month

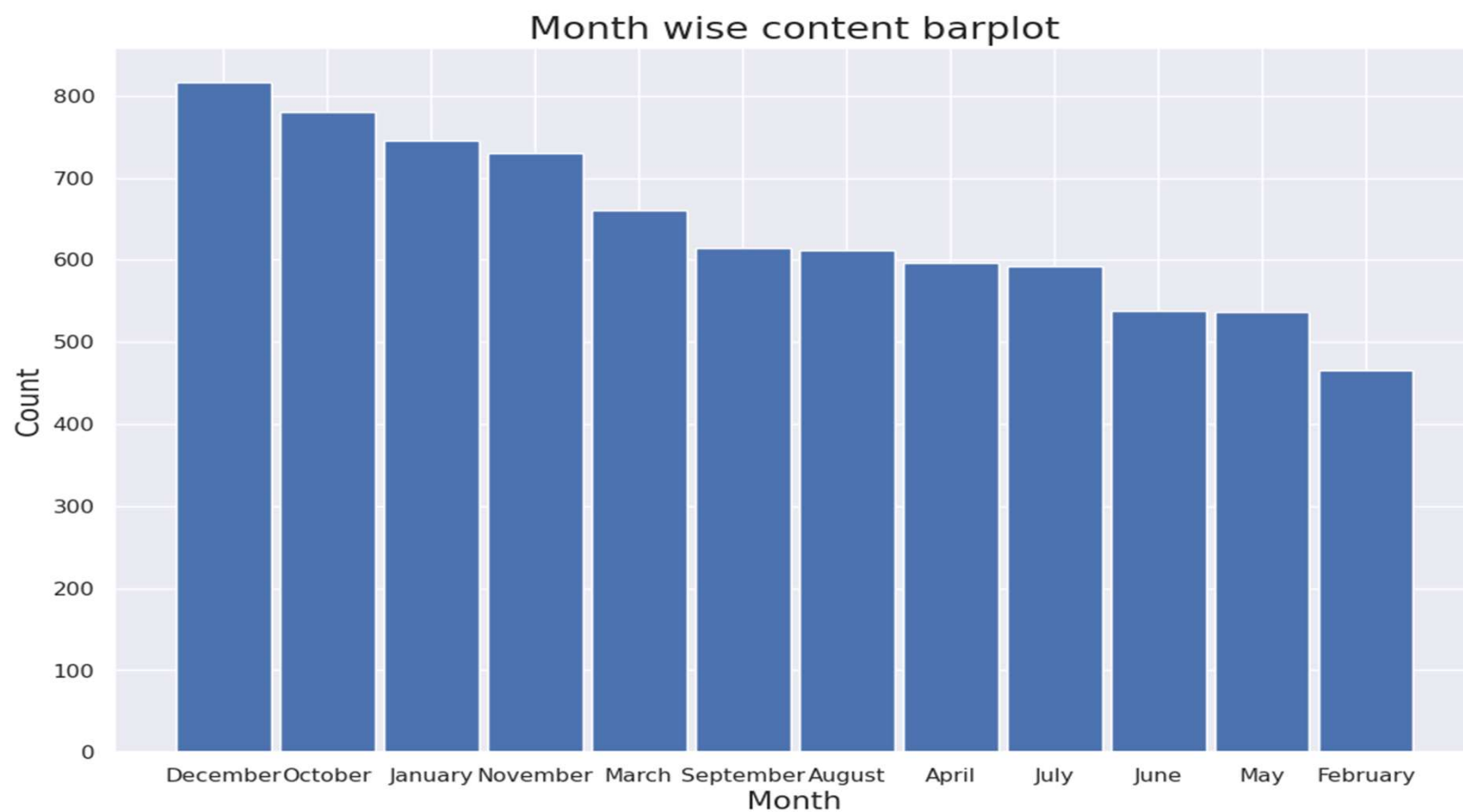


Table Of Contents

- 1. Defining problem statement
- 2. Data Cleaning & visualization
- 3. Data Preprocessing
- 4. Feature Selection
- 5. Applying different clustering methods
- 7. Applying Clustering Models 8. Conclusion

□ *Problem Statement*

- 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. 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



❑ *Data Summary*

- show_id : Unique ID for every Movie / Tv Show
- type : 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 : Genres
- description: The Summary description

□ Basic Data Exploration

- The dataset has 7787 observations and 12 features(columns).
- The dataset consists of eleven textual columns and one numeric column(release_year)
- No Duplicate values.

```
#   Column      Non-Null Count  Dtype
---  -
0   show_id    7787 non-null      object
1   type        7787 non-null      object
2   title       7787 non-null      object
3   director    5398 non-null      object
4   cast        7069 non-null      object
5   country     7280 non-null      object
6   date_added  7777 non-null      object
7   release_year 7787 non-null      int64
8   rating      7780 non-null      object
9   duration    7787 non-null      object
10  listed_in   7787 non-null      object
11  description  7787 non-null      object
dtypes: int64(1), object(11)
```

❖ *Checking NaN values*

☐ Null values present in this columns

- director
- cast
- country
- Rating

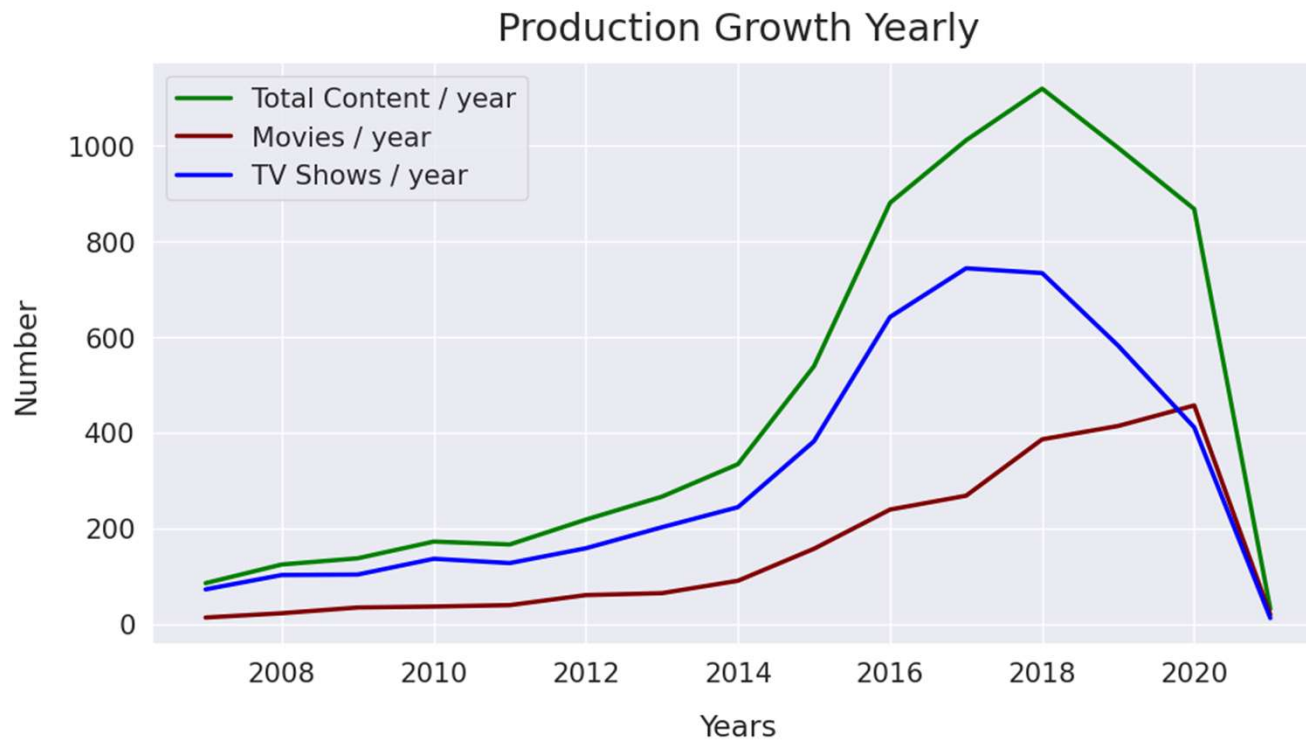
☐ No missing value present in this columns

- show_id
- type
- title
- date_added
- release_year
- duration
- listed_in
- description

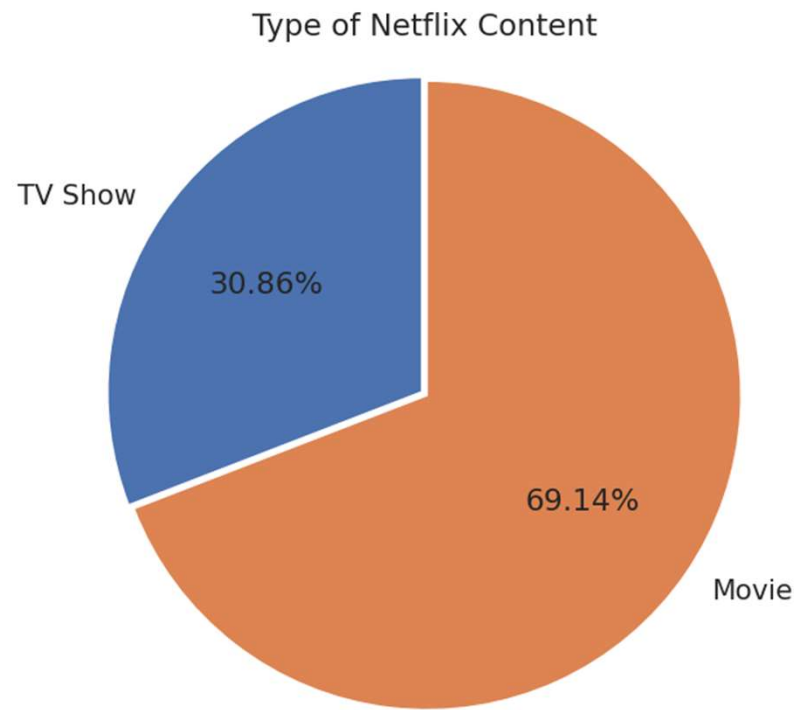
❑ *Data Cleaning*

- Removing unnecessary columns like 'director', 'cast'
- Dropping all the NaN containing date_added observations(Only 10 observations was there)
- Created 4 new columns
- No_of_categories based on listed_in
- Date_added_month based on date_added

❑ *Production Yearly Growth*

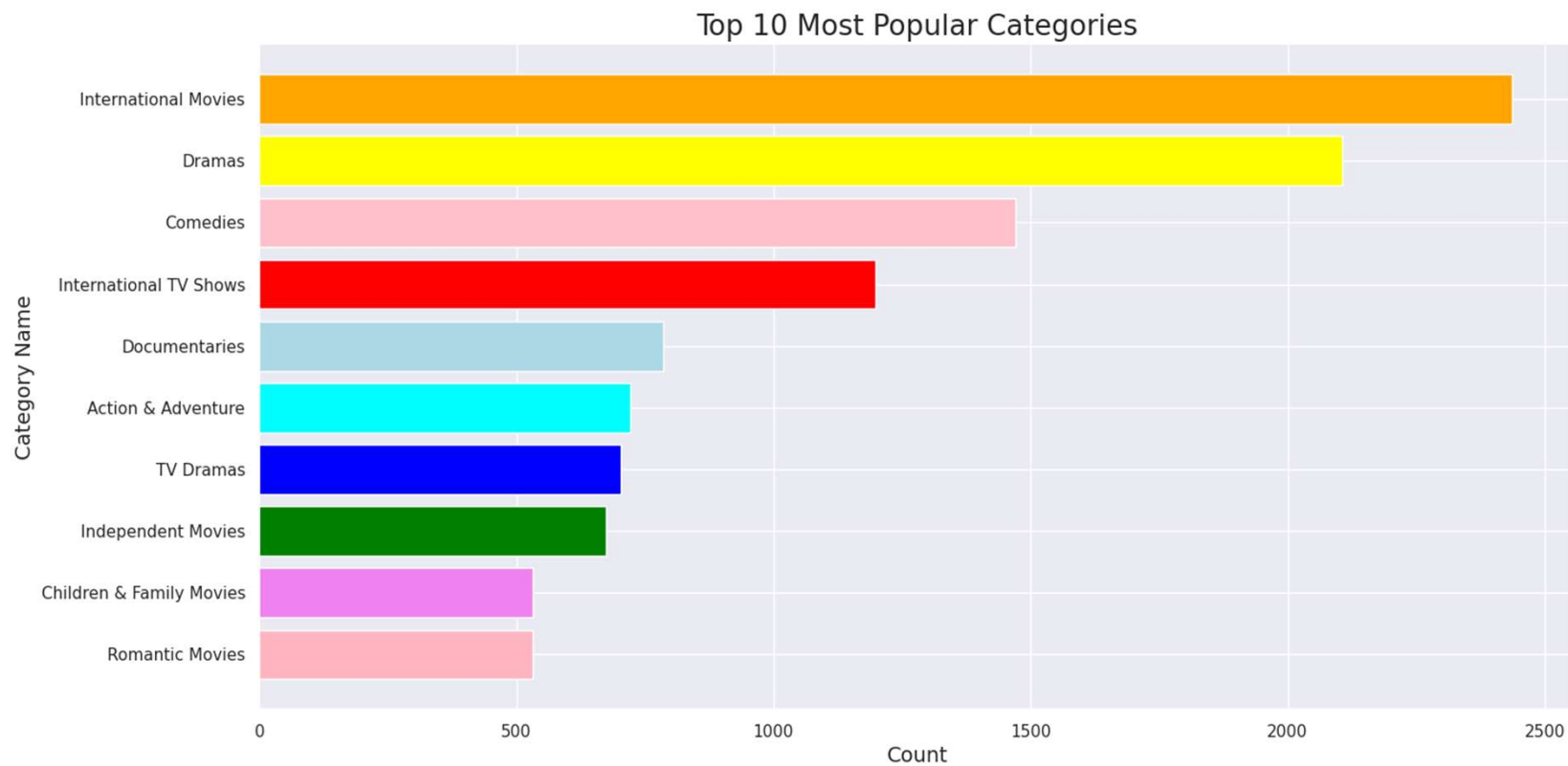


Tv shows or Movies

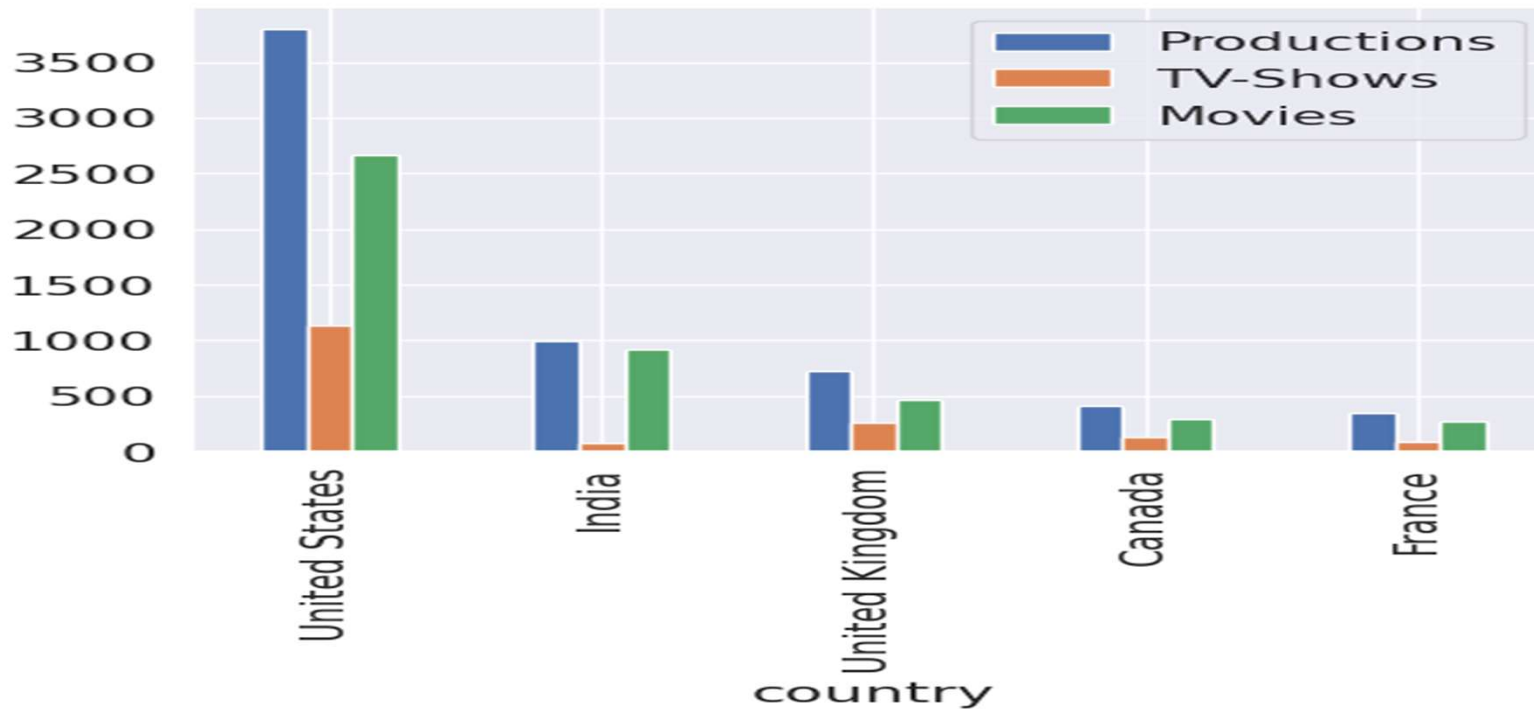


- Most of the contents are Movies
- Less than $\frac{1}{3}$ content are Tv Shows

TOP 10 Most Occurred Category By Count

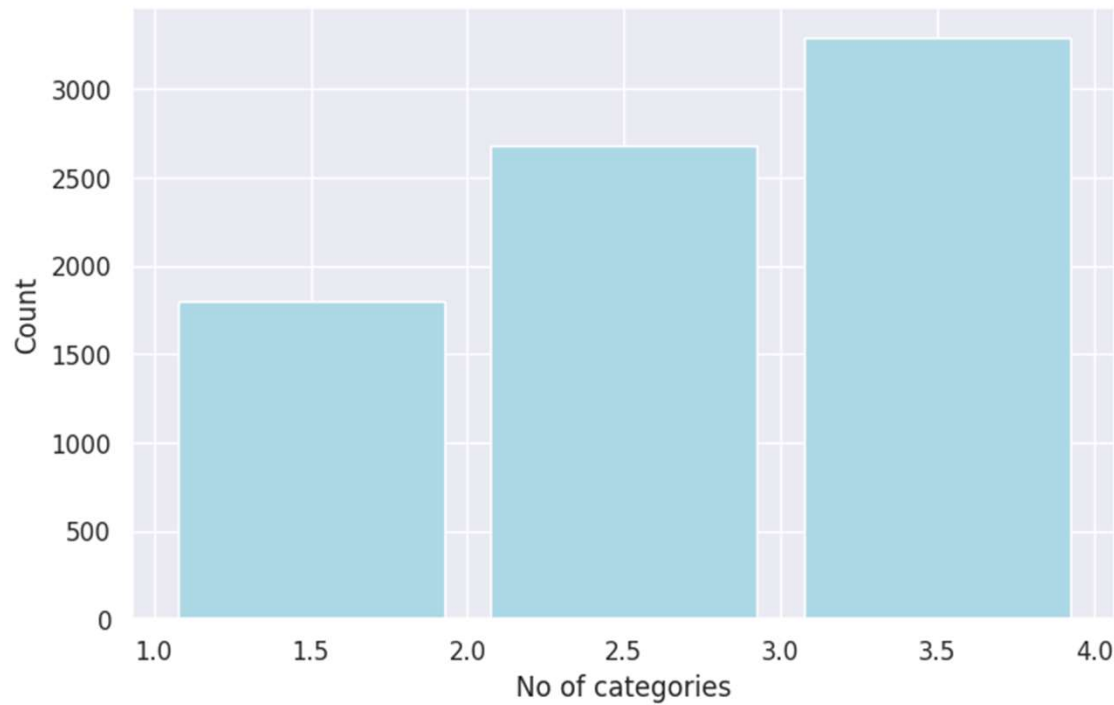


Countries producing most no of contents

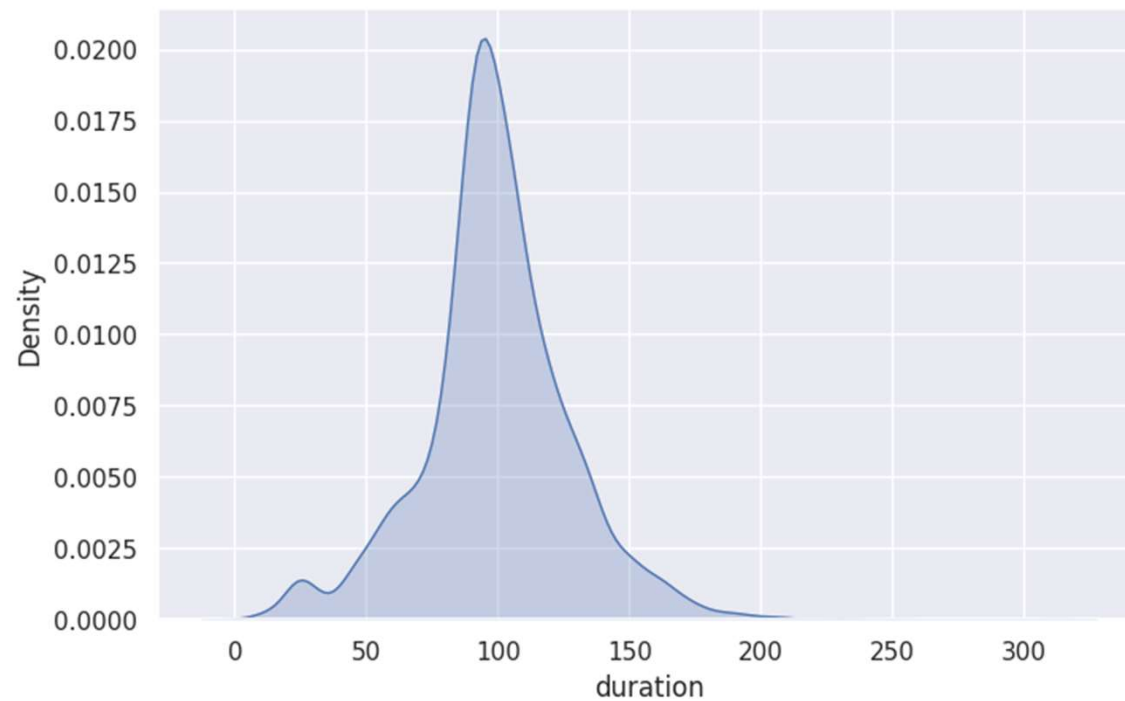


- United states have the most number of content and then india and so on
- We can conclude that except Japan other countries are producing movies more than TV-Show

How many no of categories are present there in each content

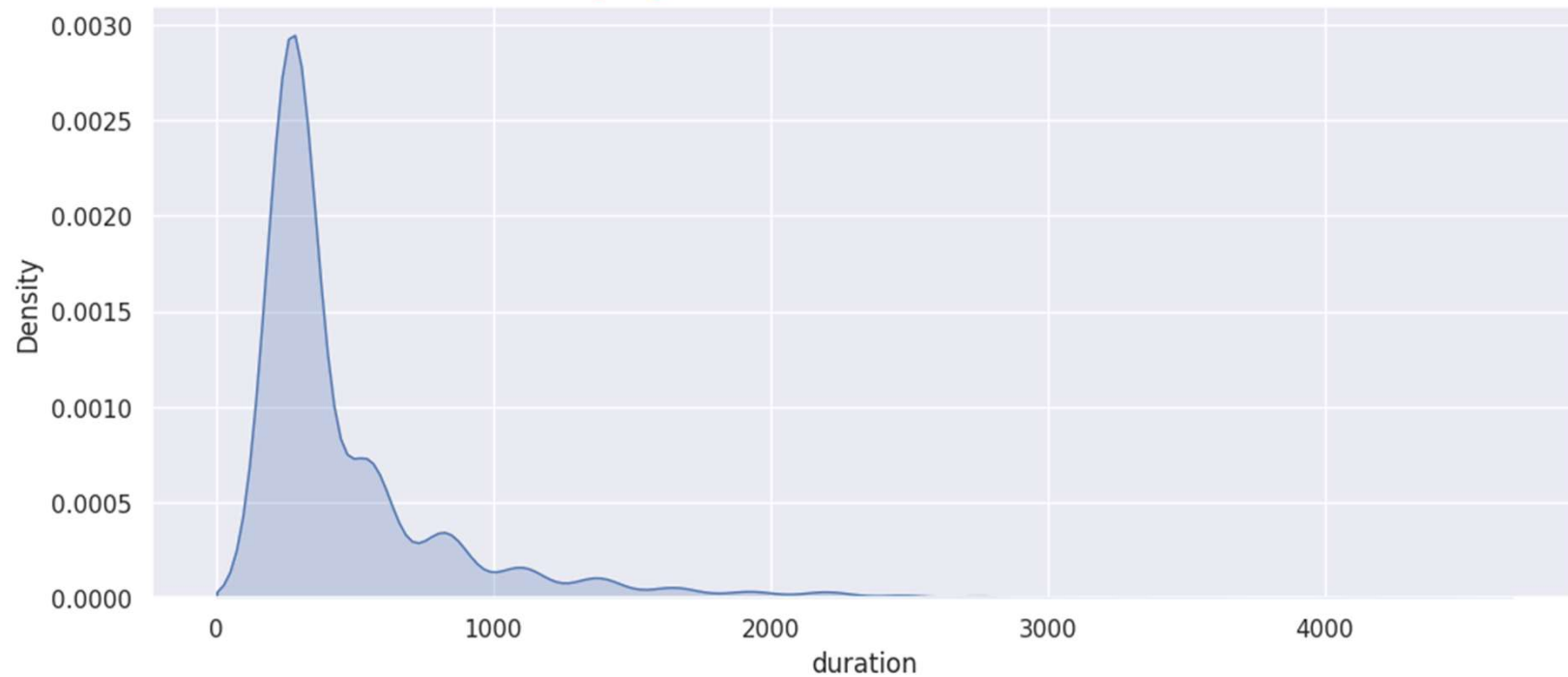


Movie wise density plot



- ❑ Most movies are about 70 to 120 min duration for movie

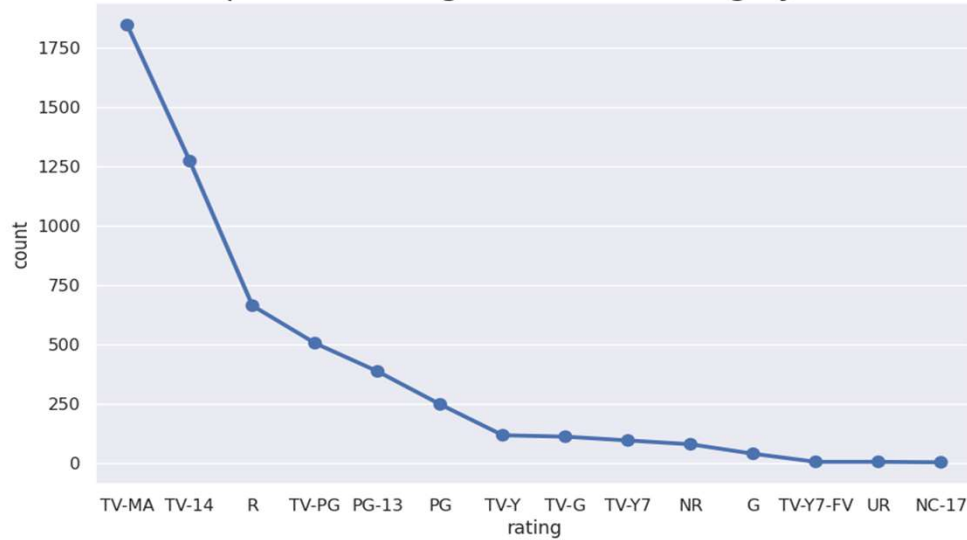
TV-Shows wise density plot



- Most contents are about 0 to 750 min duration for movies
- There are very few shows which is having more than 1000 mins. (may be the no of episodes/ seasons are more

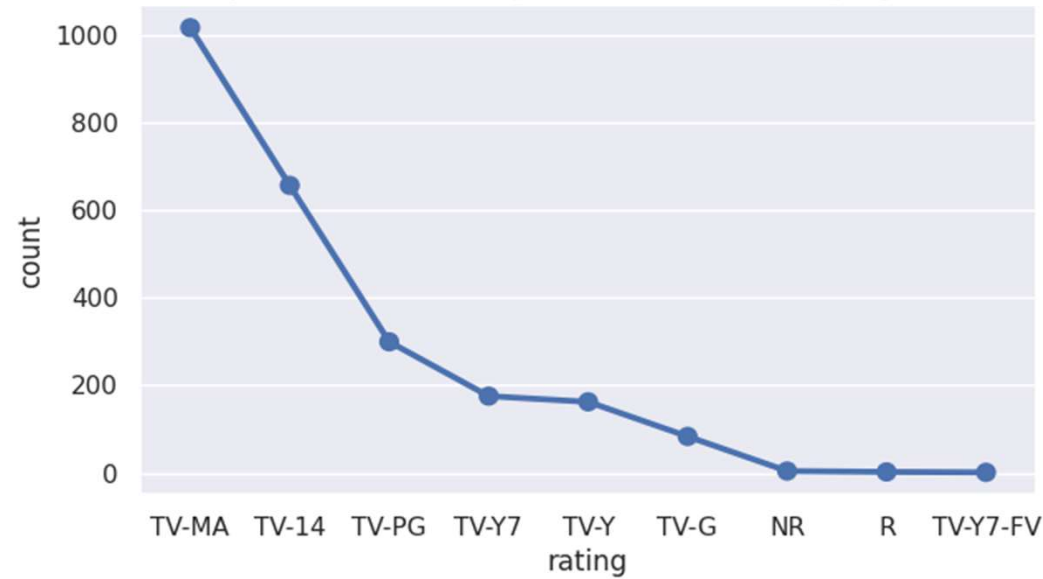
TOP Content Based On Rating

Top Movie Ratings Based On Rating System



- Most of the contents got ratings like
- TV-MA (For Mature Audiences)
- TV-14 (May be unsuitable for children under 14)
- TV-PG (Parental Guidance Suggested)
- NR (Not Rated)

Top TV Show Ratings Based On Rating System



Feature Selection & ML algo used

- ❑ Only selected 3 features , to do clustering

1. no_of_category.
2. Length(description)
3. Length(listed-in)

- ❑ Using Standard Scaler

- ❑ Used 5 algo to find out best k value

1. Silhouette score
2. Elbow Method
3. DBSCAN
4. Dendrogram
5. Agglomerative Clustering

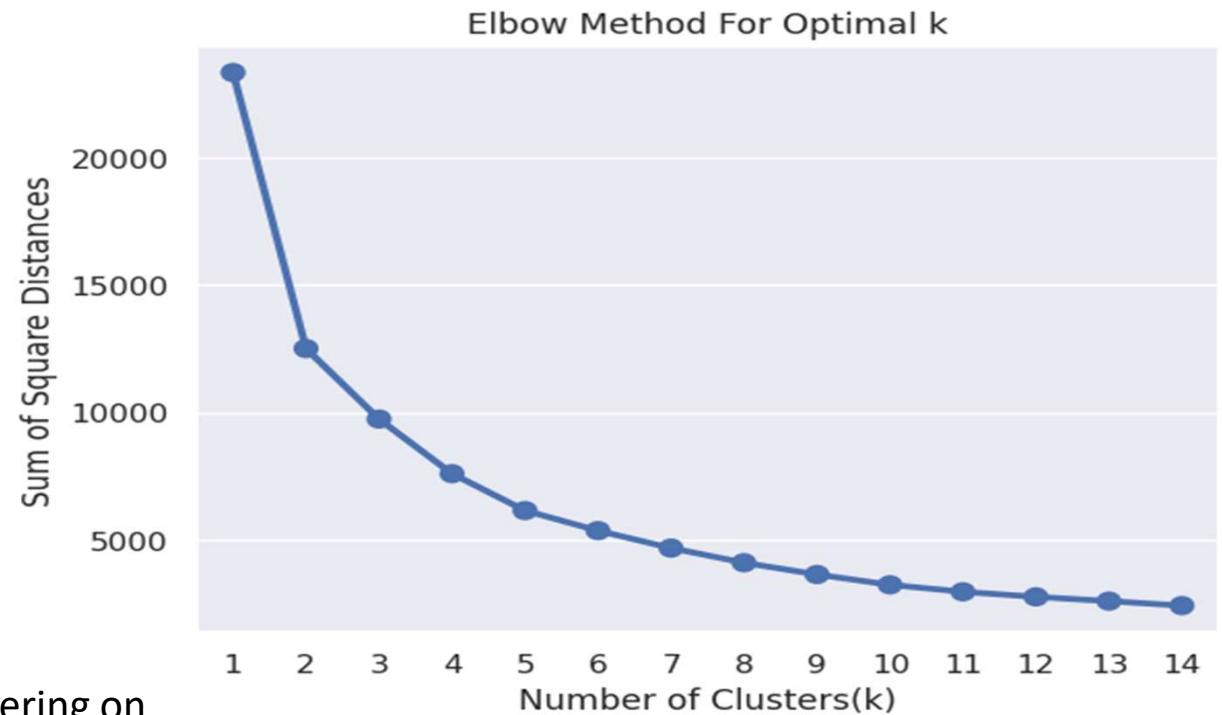
□ *Silhouette Score*

- Silhouette Coefficient Formula

$$S = (b - a) / \max(a, b).$$

- 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 a.

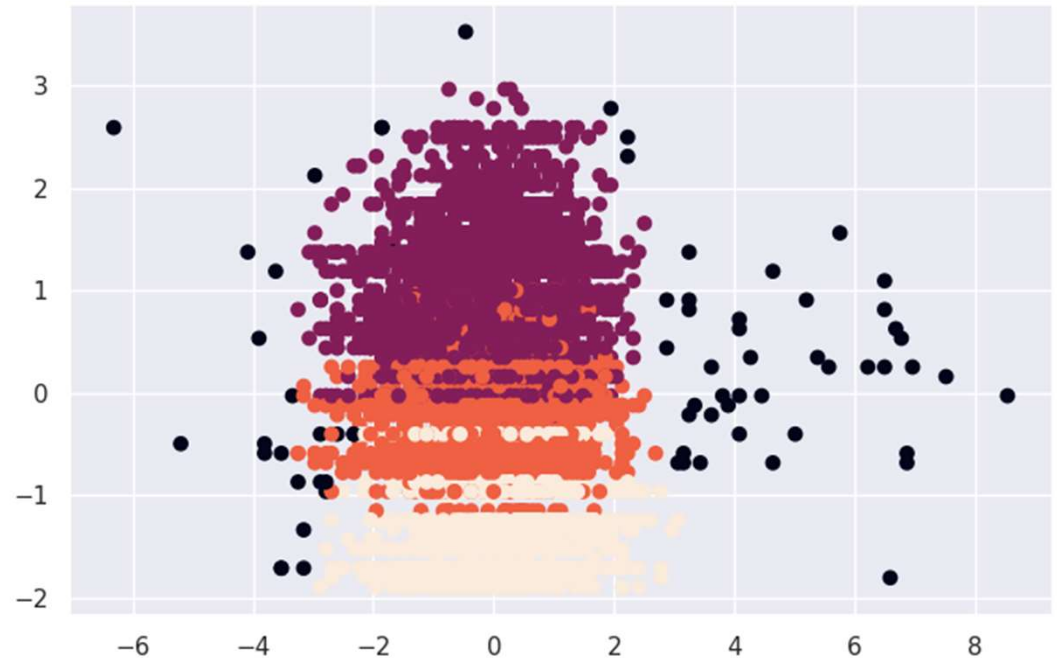
❑ *Elbow Method*



The elbow method runs k-means clustering on the dataset for a range of values for k (say from 1-15) and then for each value of k computes WCSS value . By default, the distortion score is computed, the sum of square distances from each point to its assigned center.

Agglomerative Clustering

Steps: - 1. Each data point is assigned as a single cluster. 2. Determine the distance measurement and calculate the distance matrix. 3. Determine the linkage criteria to merge the clusters. 4. Update the distance matrix. 5. Repeat the process until every data point become one cluster



Conclusion

1. Director and cast contains a large number of null values so we will drop these 2 columns .
2. In this dataset there are two types of contents where 30.86% includes TV shows and the remaining 69.14% carries Movies.
3. We have reached a conclusion from our analysis from the content added over years that Netflix is focusing movies and TV shows (From 2016 data we get to know that Movies is increased by 80% and TV shows is increased by 73% compare)
4. From the dataset insights we can conclude that the most number of TV Shows released in 2017 and for Movies it is 2020
5. On Netflix USA has the largest number of contents. And most of the countries preferred to produce movies more than TV shows.
6. Most of the movies are belonging to 3 categories
7. TOP 3 content categories are International movies , dramas , comedies.
8. In text analysis (NLP) I used stop words, removed punctuations , stemming & TF-IDF vectorizer and other functions of NLP.
9. Applied different clustering models like K means, hierarchical, Agglomerative clustering, DBSCAN on data we got the best cluster arrangements.
10. By applying different clustering algorithms to our dataset .we get the optimal number of cluster is equal to 3.