

Capstone Project 04

Netflix Movies and TV Shows Clustering

By: Prashant Gaikwad
(Individual Project)

Points for Discussion

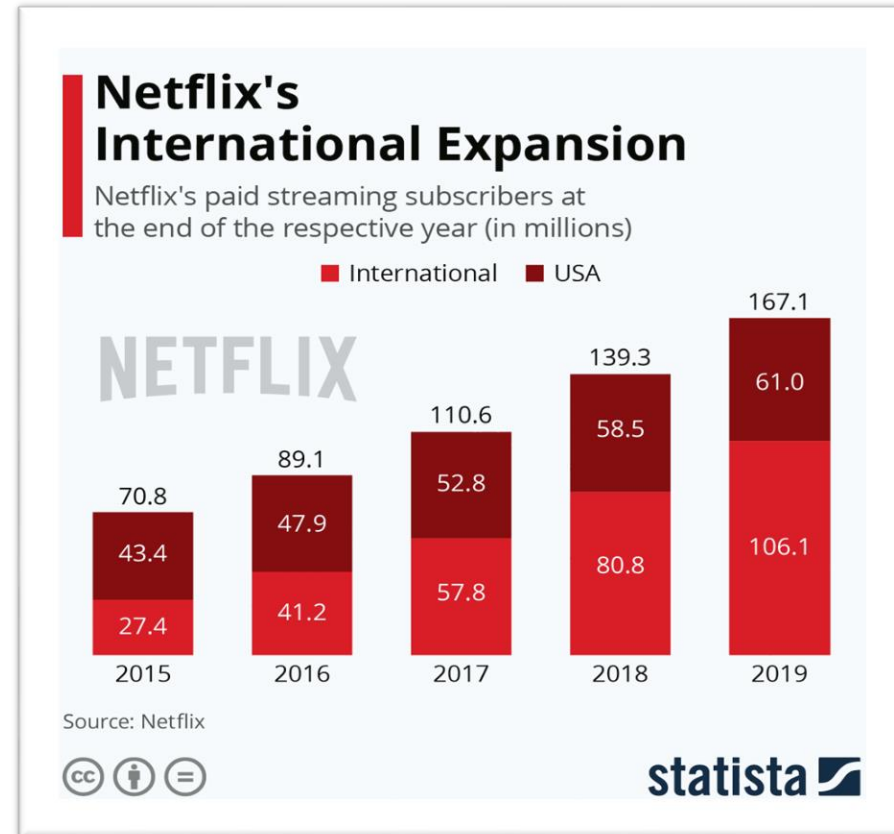
- ☐ Introduction
- ☐ Problem Statement
- ☐ Netflix Movies and TV Shows Dataset
- ☐ Data Exploration
- ☐ Data Preparation & Cleaning
- ☐ Exploratory Data Analysis (EDA)
- ☐ Feature Reduction & Feature Engineering
- ☐ Clustering Model Building
- ☐ Word Cloud
- ☐ Conclusion

INTRODUCTION

- Many online streaming services offer a large number of TV shows, which are at our disposal to watch, at the price of a subscription cost.
- Netflix is a popular entertainment service used by people around the world.
- We will explore the Netflix dataset through visualizations, graphs using python libraries and clustering model building.



- We will explore the Netflix dataset through visualizations like graphs using python libraries
- Understanding what type content is available in different countries
- Is Netflix has increasingly focusing on TV rather than movies in recent years.
- Clustering similar content by matching text-based features



PROBLEM STATEMENT

Netflix is an American technology and media services provider and production company headquartered in Los Gatos, California. Netflix was founded in 1997 by Reed Hastings and Marc Randolph in Scotts Valley, California. The company's primary business is its subscription-based streaming service, which offers online streaming of a library of films and television series, including those produced in-house.

Their most successful algorithm, Netflix Recommendation Engine (NRE), is made up of algorithms which filter content based on each individual user profile. The engine filters over 3,000 titles at a time using 1,300 recommendation clusters based on user preferences.

It's so accurate that 80% of Netflix viewer activity is driven by personalized recommendations from the engine.



Netflix Movies and TV Shows Dataset

1. show_id : Unique ID for every
Movie / Tv Show

2. type : Identifier - A Movie or
TV Show

3. title : Title of the Movie / Tv
Show

4. director : Director of the Movie

5. cast : Actors involved in the
movie / show

6. country : Country where the
movie / show was produced

7. date_added : Date it was
added on Netflix

8. release_year : Actual
Release year of the movie / show

9. rating : TV Rating of the
movie / show

10. duration : Total Duration - in
minutes or number of seasons

11. listed_in : Genre

12. description: The Summary
description

DATA EXPLORATION

- Data exploration refers to a data user being able to find his or her way through large amounts of data in order to gather necessary information.
- Explored and analyzed the data to discover key factors responsible for app engagement and success.



DATA PREPARATION & CLEANING

- Data preparation is the process of cleaning and transforming raw data prior to processing and analysis. It is an important step prior to processing and often involves reformatting data, making corrections to data, and the combining of data sets to enrich data.
- Data cleansing or data cleaning is the process of detecting and correcting (or removing) corrupt or inaccurate records from a record set, table, or database and refers to identifying incomplete, incorrect, inaccurate, or irrelevant parts of the data and then replacing, modifying, or deleting the dirty or coarse data.



- Dataset contains a large number of null values which might tend to disturb our analysis hence we dropped some unnecessary variables.
- Used “fillna()” method to filled some required variable null values from them at the beginning of project in order to get a better result.
- Encoded all categorical variable to numeric values



EXPLORATORY DATA ANALYSIS

- Exploratory Data Analysis is a process of examining or understanding the data and extracting insights or main characteristics of the data.
- **EDA** is very essential because it is a good practice to first understand the problem statement and the various relationships between the data features before getting hands dirty.



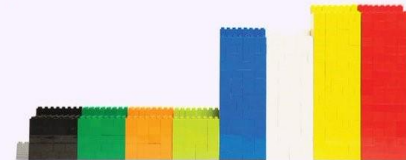
Data



Sorted



Arranged



Presented Visually

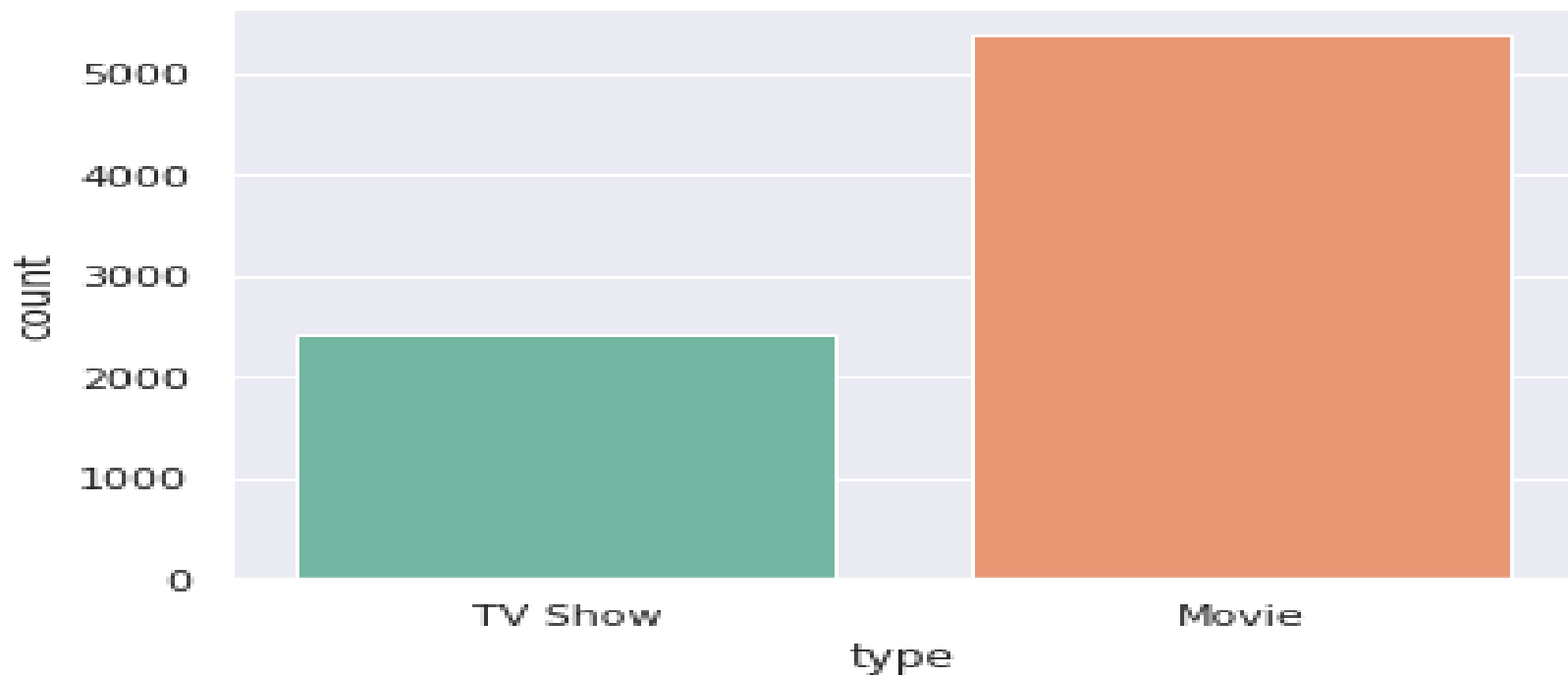
Types of EDA are:

- Univariate Analysis - analysis of a single variable
- Bivariate Analysis - analysis of exactly two variables
- Multivariate Analysis - analysis of dependent variable and multiple independent variables

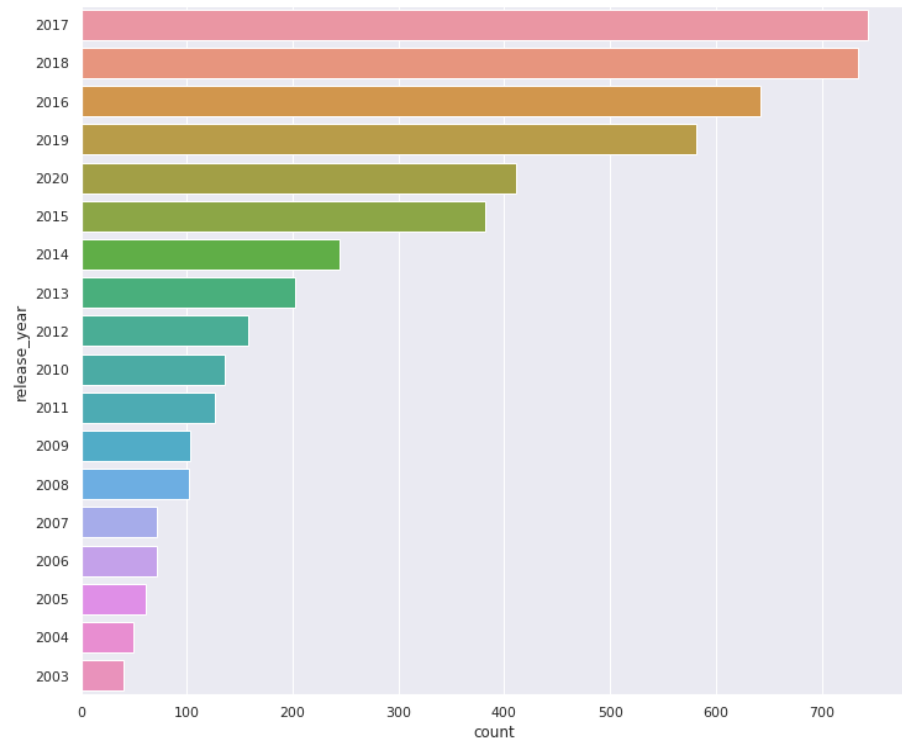
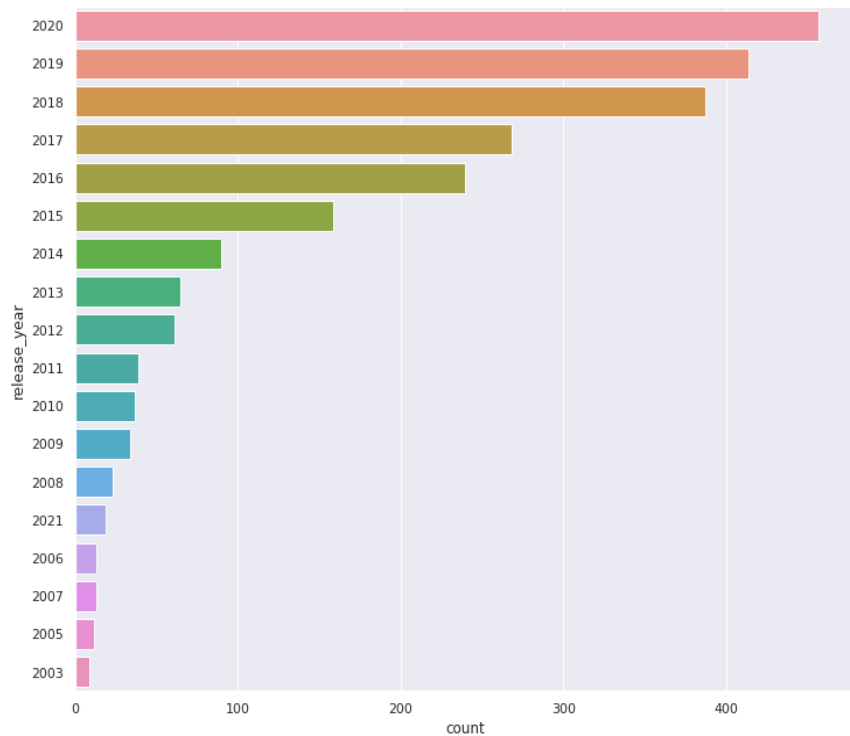
The primary motive of EDA is to

- Examine the data distribution
- Handling missing values of the dataset
(a most common issue with every dataset)
- Handling the outliers
- Removing duplicate data
- Encoding the categorical variables
- Normalizing and Scaling

Type of content

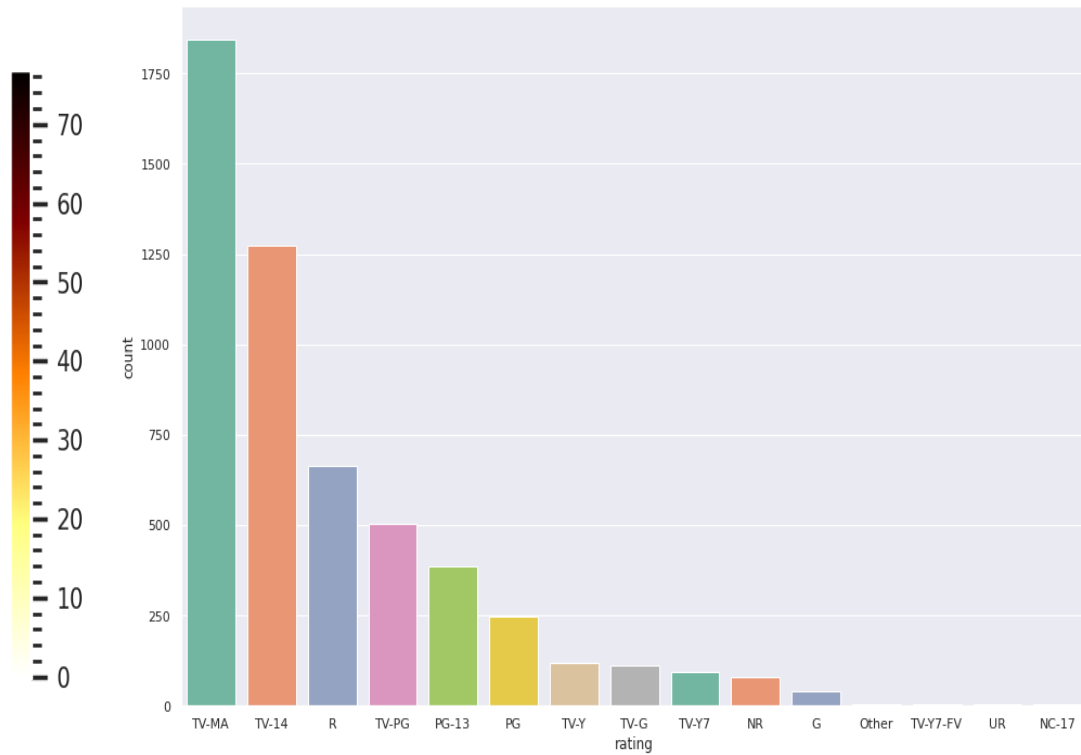
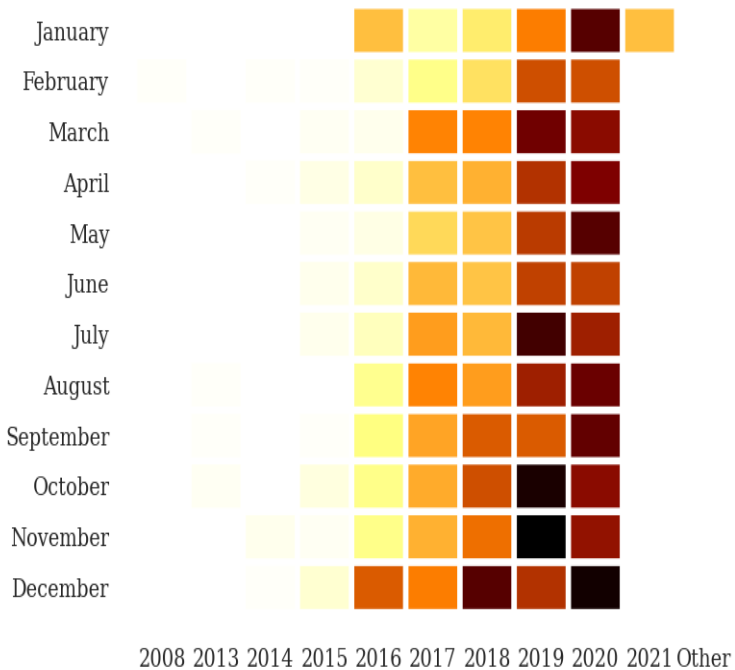


Release year of Shows and Movies



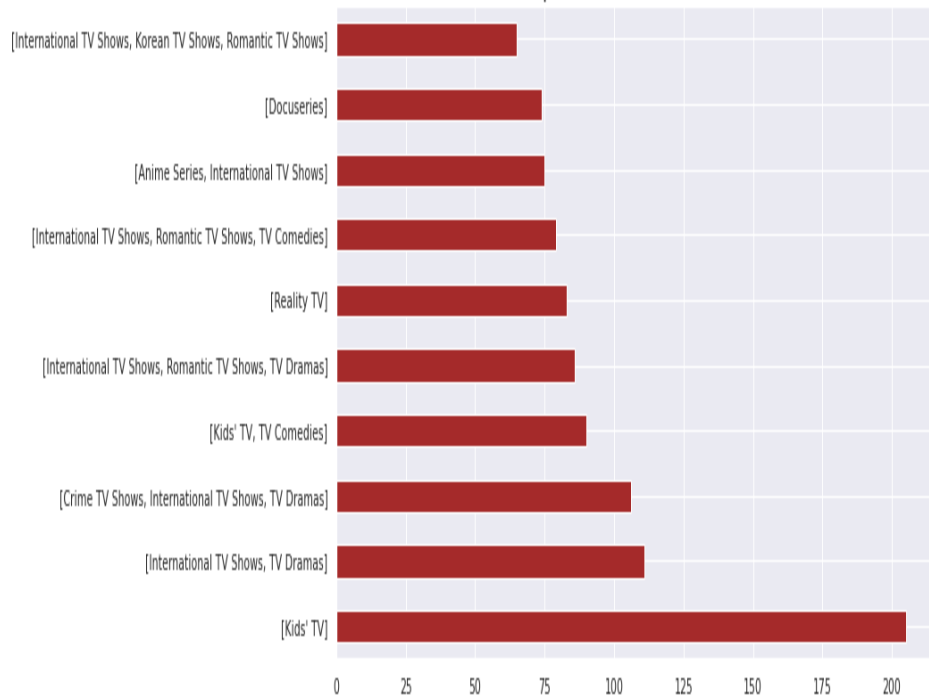
EDA (continued)

Netflix Contents Update

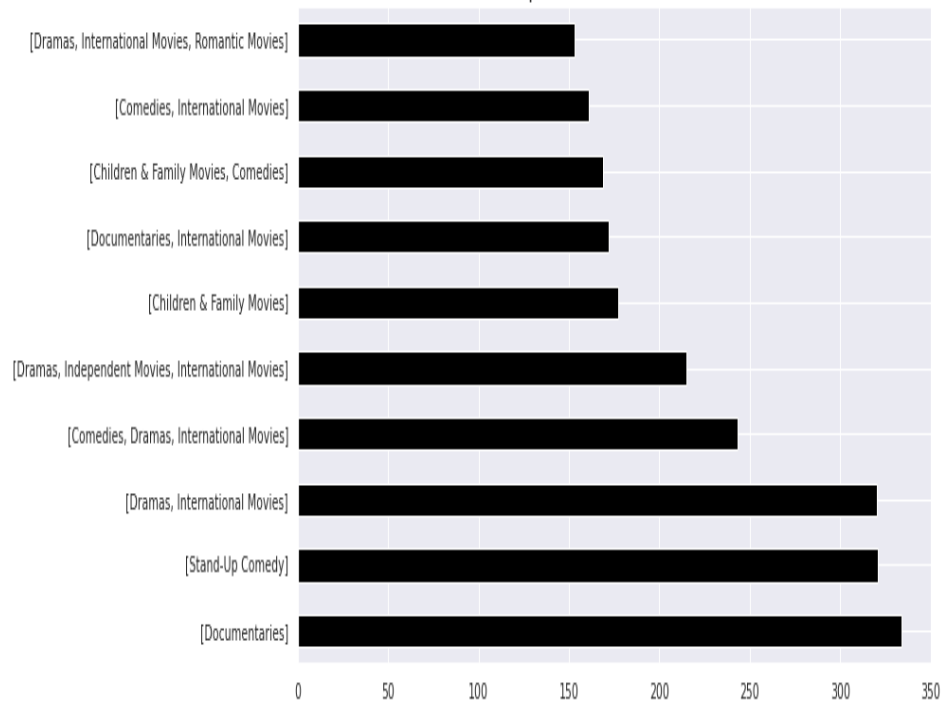


TOP 10 GENRES

Top 10 Genres of TV Shows

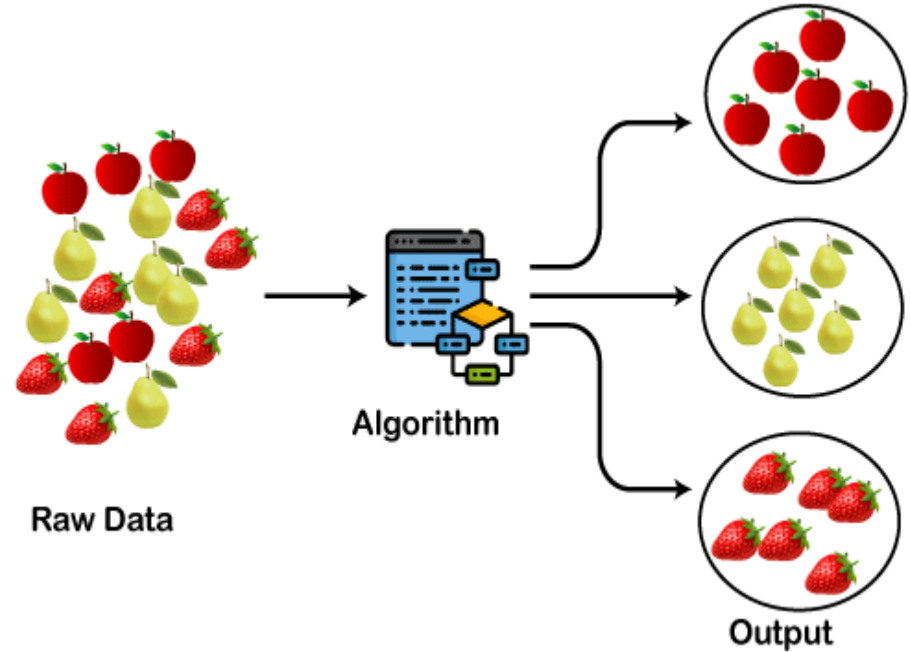


Top 10 Genres of Movies



Cluster Model Building

- Clustering :
- Clustering is the task of dividing the population or data points into a number of groups such that data points in the same groups are more similar to other data points in the same group than those in other groups.
- In simple words, the aim is to segregate groups with similar traits and assign them into clusters.
-



K-Means Clustering:

The objective of K-means is simple: group similar data points together and discover underlying patterns. To achieve this objective, K-means looks for a fixed number (k) of clusters in a dataset.

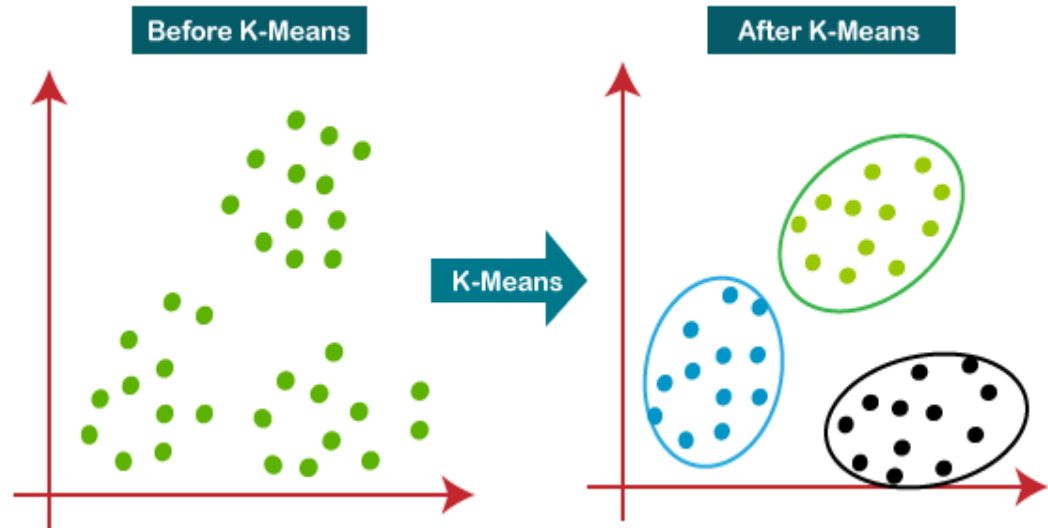
The 'means' in the K-means refers to averaging of the data; that is, finding the centroid.

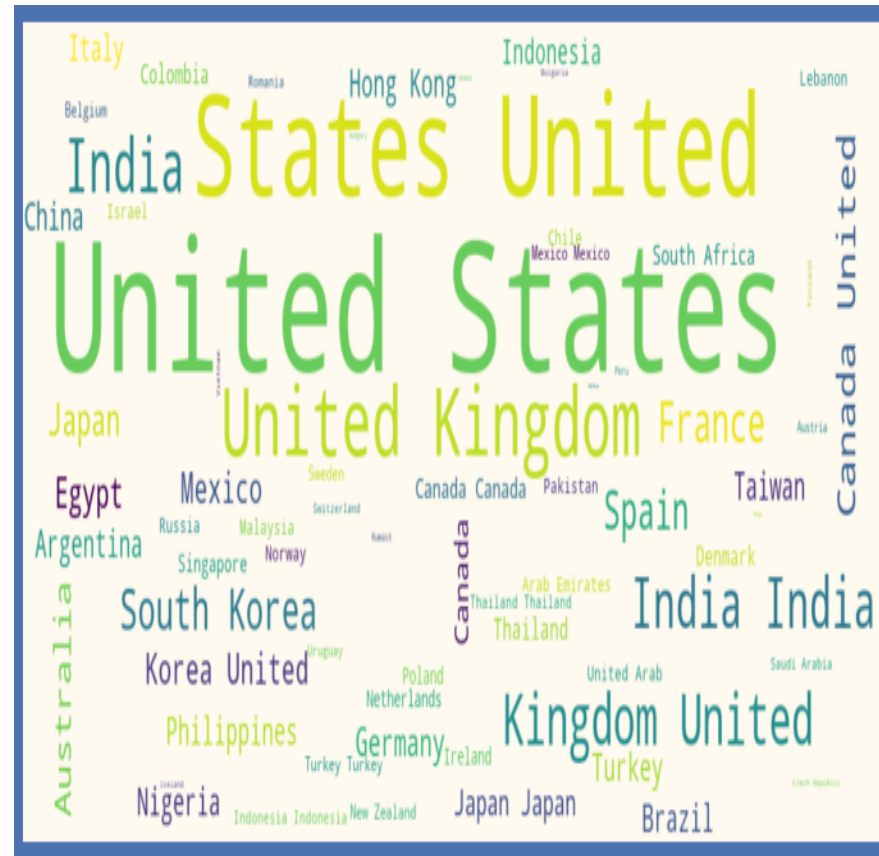
Advantages

Recommendation engines

Market segmentation

Social network analysis





CONCLUSION

- 1) The largest count of movies are made with the 'TV-MA' rating. "TV-MA" is a rating assigned by the TV Parental Guidelines to a television program that was designed for mature audiences only.
- 2) Highest number of movies released in Netflix was in the year 2017 and 2018
- 3) Highest number of Tv shows released in Netflix was in the years 2018, 2019, and 2020.
- 4) The months of October, November, December and January had the largest number of films and television series released.
- 5) In the given years we have seen in the month of Feb, May less movies was released so producer can release on this months
- 6) There are about 70% movies and 30% TV shows on Netflix.

7) The United States has the highest number of content on Netflix by a huge margin followed by India.

8) Raul Campos and Jan Sulter collectively have directed the most content on Netflix.

9) James, Lee, Michael, Daniel, John, etc have casted in most of the content on Netflix.

10) Anupam Kher has acted in the highest number of films on Netflix. Drama is the most popular genre followed by comedy.

11) More of the content is released in holiday season - October, November, December and January.

12) The number of releases have significantly increased after 2015 and have dropped in 2021 because of Covid 19.

13) The most frequent words used in description are LIFE, FAMILY, WOMEN, FATHER, LOVE, FRIEND, etc.

THANK YOU