

Capstone Project: 04

Unsupervised Machine Learning Clustering

Netflix Movies and TV Shows Clustering

Author: Pankaj Beldar

Email id: pankajrbell@gmail.com

Abstract: Netflix is widely regarded as the leading over-the-top (OTT) platform because of its reputation for offering users a wide variety of high-quality streaming content. The reason why Netflix's services are so popular worldwide is that the company uses cutting-edge technology like artificial intelligence and machine learning to provide consumers with more appropriate and intuitive suggestions. This article explains how Netflix uses artificial intelligence, data science, and machine learning. Even now, over twenty years after it first launched, Netflix is still working to improve its service. Artificial intelligence has been put to use by Netflix to provide customers with the greatest possible service and experience.

1. Introduction:

It's fascinating how Netflix applies AI/Data Science/ML to running its operations, such as by implementing algorithms to provide movie recommendations and using AI to guarantee high-quality streaming even at reduced bandwidths. The following are some of the numerous

applications of AI, data science, and machine learning at Netflix. Improvement in Netflix's AI integration has made widespread individualization possible. Simply said, the AI engine keeps an eye on the flow of information and sometimes takes over so that it may make judgments and suggestions at predetermined moments. Netflix's AI considers your viewing habits and hobbies to provide Netflix recommendations. Users can take charge of their multimedia streaming and customize their interactions owing to the system's ability to compile and recommend content based on their preferences.

2. Data Overview:

There are a total 7787 entities and 12 features in our dataset. About 30.67% data is missing in director, 9.22% in cast, 6.51% in country and 0.0898 % in rating.



The attribute 'director', 'cast', 'country', 'date_added', 'rating' consists of missing values.

3. Data Cleaning:

• Handling Missing Value

The attribute 'director', 'cast', 'country', 'date_added', 'rating' consists of missing values.

To tackle missing values, we will replace 'country' and 'rating' missing values by the most frequent entity that is 'United States' and 'TV-MA' respectively.

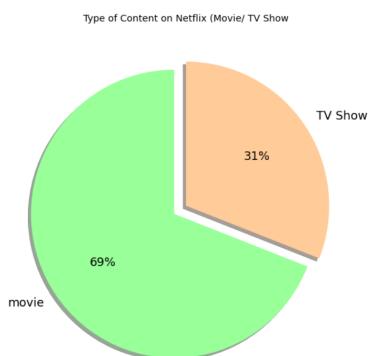
missing values in 'cast' by 'unknown'.

There are around 30.68 % values missing in 'director', hence we decide to drop it.



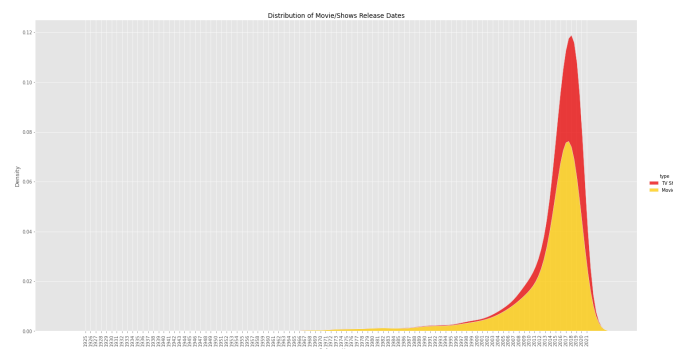
No duplicate values exist in the whole dataset.

4. Exploratory Data Analysis:

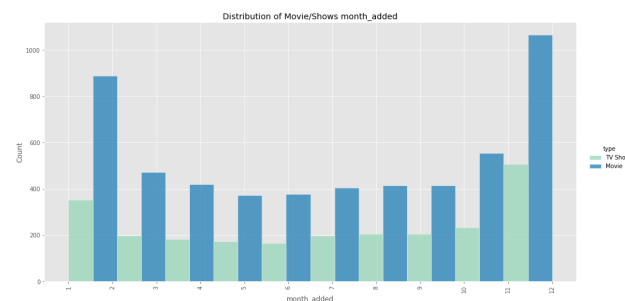


69% of the content available on Netflix are movies; the remaining 31% are TV Shows.

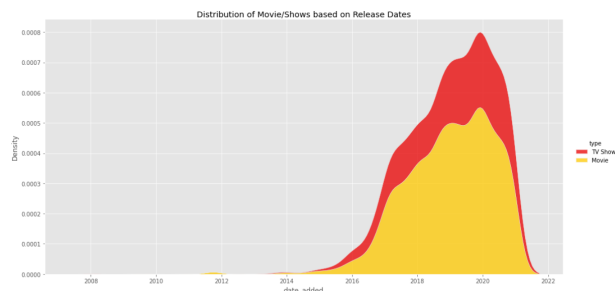
Netflix has 5377 movies, which is more than double the quantity of TV shows.



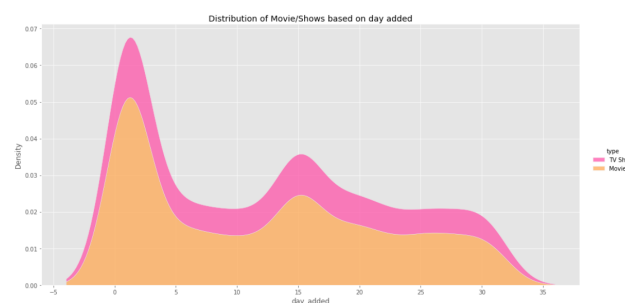
In recent years more TV Shows are released as compared to Movies on Netflix. Less number of TV shows and Movies were released in 2020-2021 due to coronavirus pandemic



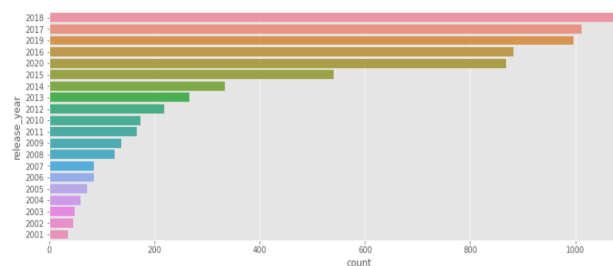
Most of the Movies/TV Shows were added in the month of December and January. Number of Movies added on Netflix are more as compared to TV Shows throughout the year.



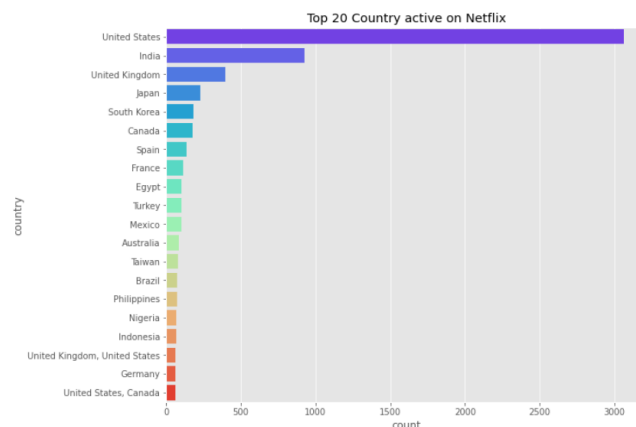
In recent few year more number of TV Shows were added on NetFlix as compared to Movies , We can say Netflix is more focusing on TV Shows than Movies.



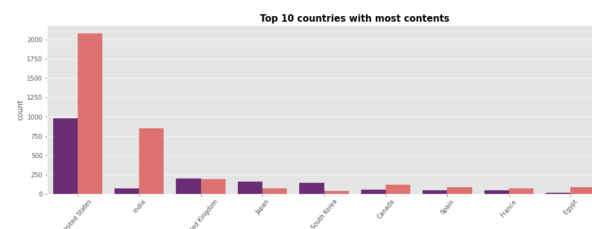
Most of the TV Shows/Movies are added on early 5 days of the month, then middle of the month. Again Number of TV Shows added is more as compared to Movies.



Number of Movies/TV Shows added on Netflix had drastically increased after 2014. Maximum content was uploaded on Netflix in the year of 2018, after that due to Covid, The count is decreasing slightly.



The United States, India, United Kingdom, Japan, South Korea, Canada, Spain, France, Egypt and Turkey are the Top 10 countries which produce most of the content on Netflix.



The United States produced most of the content on Netflix. Also, the number of movies released are more than TV Shows in the United States. In India, Canada, Spain, France, Egypt and Turkey , Most of the content on Netflix is Movies. The United Kingdom has almost equal production of Movies and TV Shows. In Japan and South Korea, Number of TV Shows are available on Netflix.

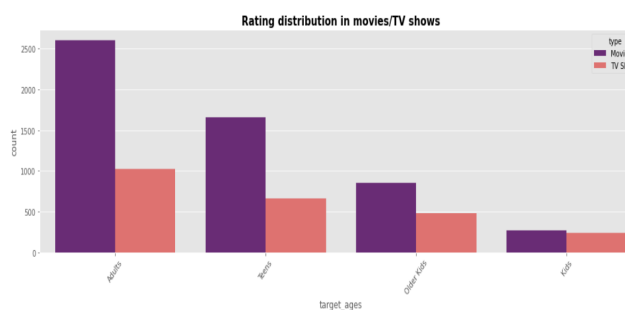
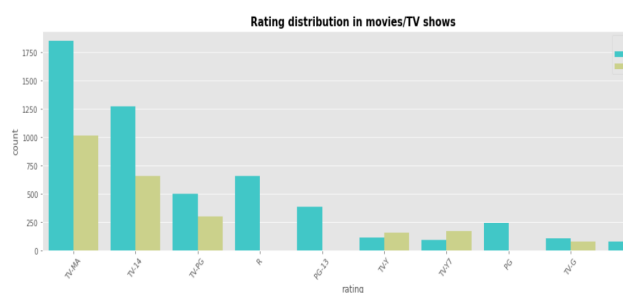
What Ages are Netflix Parental Controls Good For?

Maturity ratings for each piece of content on Netflix are determined

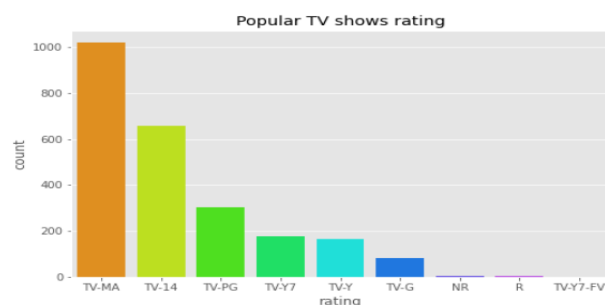
by a local standards organization or by Netflix considering the “frequency and impact of mature content in a TV show or movie”. These ratings are put in place to help your family make informed decisions about the content you are viewing.

To get a better understanding of how Netflix breaks down its ratings categories:

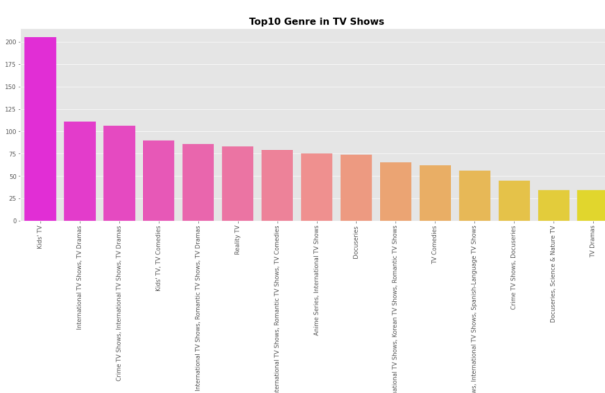
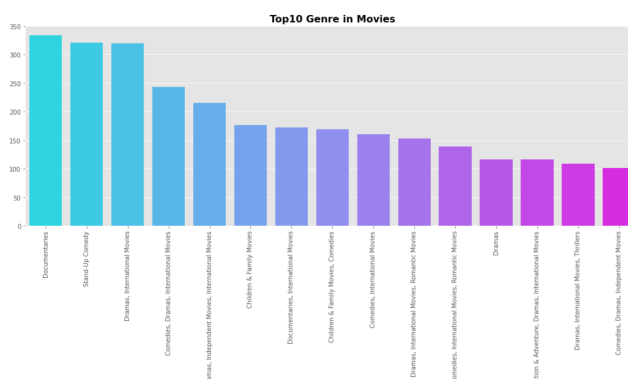
1. **Little Kids: G, TV-Y, TV-G**
2. **Older Kids: PG, TV-Y7, TV-Y7-FV, TV-PG**
3. **Teens: PG-13, TV-14**
4. **Adults: R, NC-17, TV-MA**



It is observed that, in each category, Quantity of Movies is more than the Quantity of TV Shows. The Availability of the Adult Content is more on Netflix and Least for the Kids.



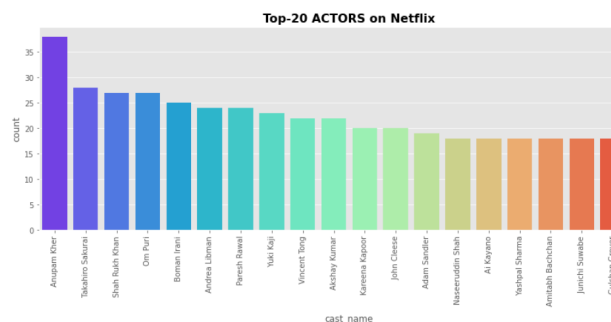
Popular Movies ratings are TV-MA, TV-14, R, TV-PG, PG-14 and PG. It is observed that Adults and Teens are mostly active on Netflix. Popular TV Shows ratings are TV-MA, TV-14, R, TV-PG, PG-14 and PG.



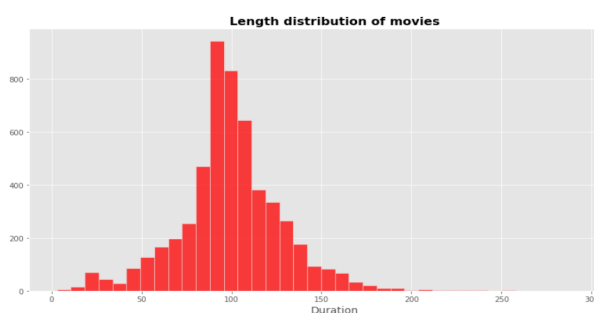
Top 5 Genres in 'TV Shows' are Kid's TV, TV Dramas, TV Crime Shows, TV Comedies, TV Romantic. Top 5 Genres in 'Movies' are Documentaries, Stand up Comedy, Dramas and International Movies,



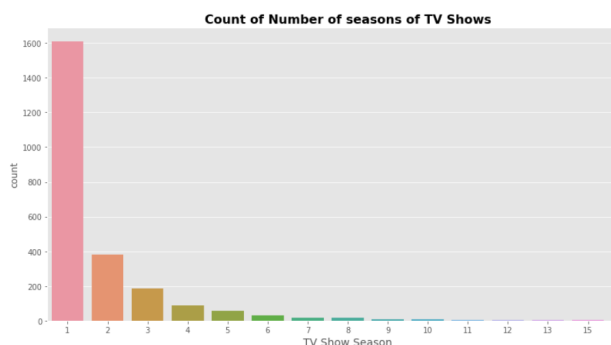
Comedies and Independent Movies.



Famous Actors on Netflix based on the Frequency of their occurrence on screen are Anupam Kher, Takahiro Sakurai, Shah Rukh Khan, Om Puri and Boman Irani and so on.

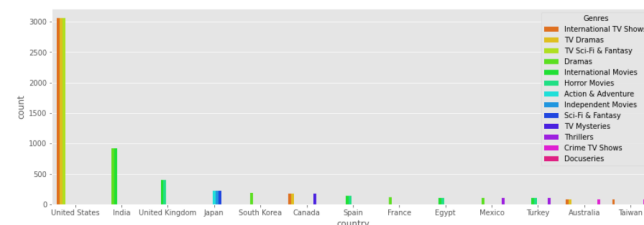


Most of the Movies/TV Shows have a duration of around 100 min.

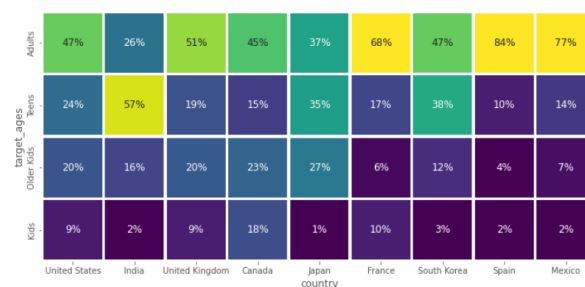


It is observed that 1608 TV Shows had only one season. The count of

longest running TV Shows is very less.



The United States produces maximum International TV Shows, TV Dramas, Sci-fi and Fantasy TV shows, International Movies. India, UK, Spain, Egypt, Mexico and Turkey are having most of the Content as Dramas and International Movies.



It is observed that content available for kids is less as compared to other categories. Content available for Adults is more in almost every country except India. In India, Most of the content is available for Teens. Netflix should focus on Teen and Adult Contents to generate maximum revenue. Spain and Mexico are producing the highest Adult Content on Netflix almost 84% and 77% respectively.

5. Text Processing:

1. Removed Punctuations
2. Removed Stopwords
3. Removed Short words

2013. The word2vec algorithm uses a neural network model to learn word associations from a large corpus of text. Once trained, such a model can detect synonymous words or suggest additional words for a partial sentence. Word2Vec creates vectors of the words that are distributed numerical representations of word features – these word features could comprise of words that represent the context of the individual words present in our vocabulary. Word embeddings eventually help in establishing the association of a word with another similar meaning word through the created vectors.

As seen in the image below where word embeddings are plotted, similar meaning words are closer in space, indicating their semantic similarity.

Word2Vec

Pre-trained vectors trained on a part of the Google News dataset (about 100 billion words). The model contains 300-dimensional vectors for 3 million words and phrases. The phrases were obtained using a simple data-driven approach described in 'Distributed Representations of Words and Phrases and their Compositionality'

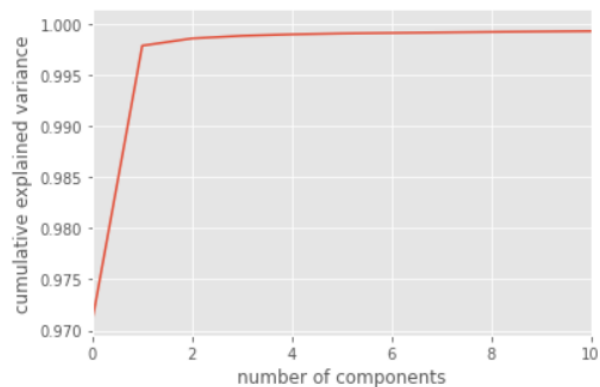
Encoding Categorical Variables

A one-hot encoding is a representation of categorical variables as binary vectors. This first requires that the categorical values be mapped to integer values. Then, each integer value is represented as a binary vector that is all zero values

except the index of the integer, which is marked with a 1.

Principal component analysis (PCA):

Principal component analysis (PCA) is a technique for **reducing the dimensionality of such datasets, increasing interpretability but at the same time minimizing information loss**. It does so by creating new uncorrelated variables that successively maximize variance.



Select 2 PCA Components

KMeans Clustering

K-Means Clustering is an Unsupervised Learning algorithm, which groups the unlabeled dataset into different clusters.

It is an iterative algorithm that divides the unlabeled dataset into k different clusters in such a way that each dataset belongs to only one group that has similar properties.

It is a centroid-based algorithm, where each cluster is associated with a centroid. The main aim of this algorithm is to minimize the sum of distances between the data

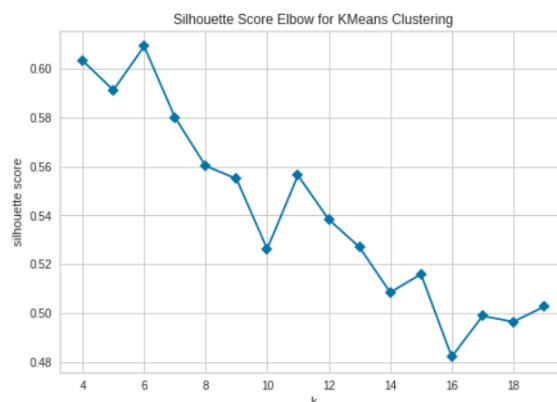
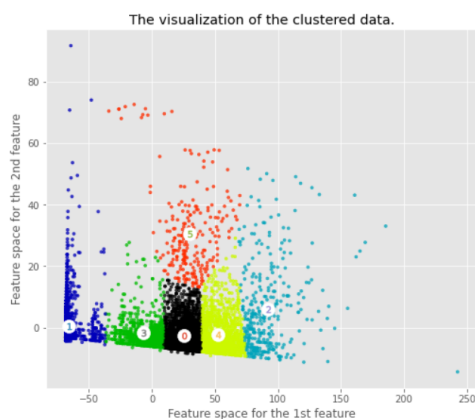
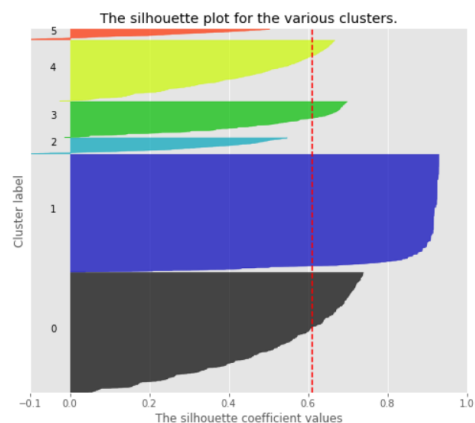
point and their corresponding clusters.

Calculation of 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. The Silhouette score is calculated for each sample of different clusters. To calculate the Silhouette score for each observation/data point, the following distances need to be found out for each observation belonging to all the clusters: Mean distance between the observation and all other data points in the same cluster. This distance can also be called a mean intra-cluster distance. The mean distance is denoted by a . Mean distance between the observation and all other data points of the next nearest cluster. This distance can also be called a mean nearest-cluster distance. The mean distance is denoted by b . The Silhouette Coefficient for a sample is $S = (b - a) / \max(a, b)$

The value of the silhouette coefficient is between [-1, 1]. A score of 1 denotes the best meaning that the data point i is very compact within the cluster to which it belongs and far away from the other clusters. The worst value is -1. Values near 0 denote overlapping clusters.

```
For n_clusters = 4 The average silhouette_score is : 0.6042559
For n_clusters = 5 The average silhouette_score is : 0.5925321
For n_clusters = 6 The average silhouette_score is : 0.6094308
For n_clusters = 7 The average silhouette_score is : 0.5814071
For n_clusters = 8 The average silhouette_score is : 0.5542456
For n_clusters = 9 The average silhouette_score is : 0.5299943
For n_clusters = 10 The average silhouette_score is : 0.550438
For n_clusters = 11 The average silhouette_score is : 0.546725
For n_clusters = 12 The average silhouette_score is : 0.534712
For n_clusters = 13 The average silhouette_score is : 0.523382
For n_clusters = 14 The average silhouette_score is : 0.483397
For n_clusters = 15 The average silhouette_score is : 0.502128
For n_clusters = 16 The average silhouette_score is : 0.500163
For n_clusters = 17 The average silhouette_score is : 0.487304
For n_clusters = 18 The average silhouette_score is : 0.478807
For n_clusters = 19 The average silhouette_score is : 0.511422
```



Maximum silhouette score is 0.6 for cluster number 6

Elbow Method to get number of clusters

The K-Elbow Visualizer implements the “elbow” method of selecting the optimal number of clusters for K-means clustering.

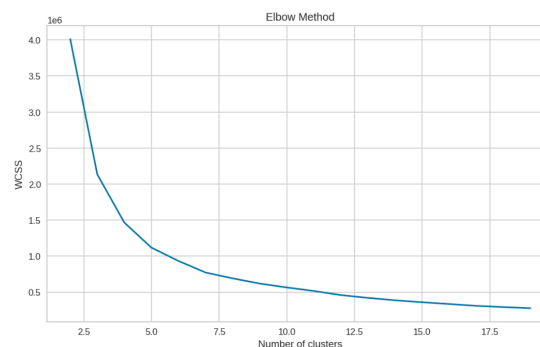
The elbow method runs k-means clustering on the dataset for a range of values for k (say from 1-10) and then for each value of k computes an average score for all clusters. By default, the distortion score is computed, the sum of square distances from each point to its assigned center.

When these overall metrics for each model are plotted, it is possible to visually determine the best value for k. If the line chart looks like an arm, then the “elbow” (the point of inflection on the curve) is the best value of k. The “arm” can be either up or down, but if there is a strong inflection point, it is a good indication that the underlying model fits best at that point.

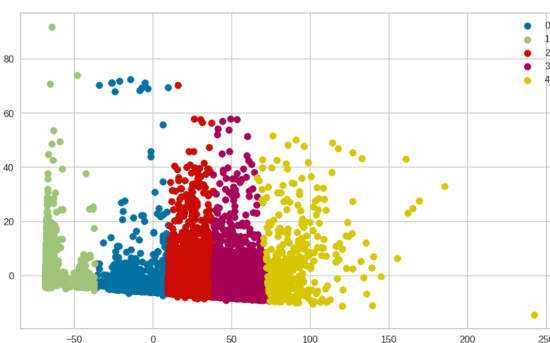
Elbow Curve Method

Perform K-means clustering with all these different values of K. For each of the K values, we calculate average distances to the centroid across all data points.

Plot these points and find the point where the average distance from the centroid falls suddenly (“Elbow”).



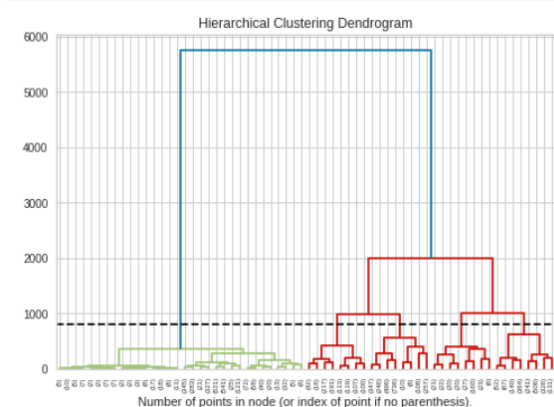
Perform Clustering considering k=5.



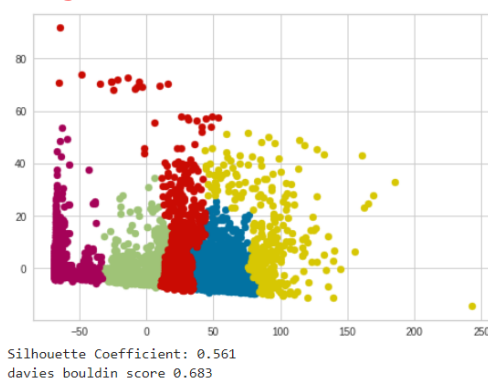
Select number of clusters for Agglomerative clustering using Dendrogram

A dendrogram is a diagram that shows the hierarchical relationship between objects. It is most commonly created as an output from hierarchical clustering. The key to interpreting a hierarchical cluster analysis is to look at the point at which any given pair of cards “join together” in the tree diagram.

Cards that join together sooner are more similar to each other than those that join together later.

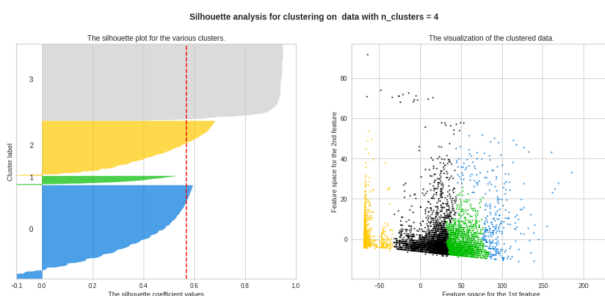


Number of clusters from Dendrogram are 5



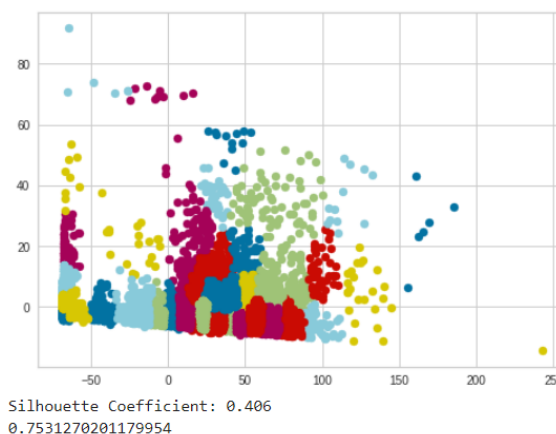
Select number of clusters for Agglomerative clustering using silhouette_score

```
For n_clusters = 4 The average silhouette_score is : 0.5689844746234257
For n_clusters = 5 The average silhouette_score is : 0.5608615188910879
For n_clusters = 6 The average silhouette_score is : 0.5374665903101778
For n_clusters = 7 The average silhouette_score is : 0.5560047393055707
For n_clusters = 8 The average silhouette_score is : 0.535151933511449
For n_clusters = 9 The average silhouette_score is : 0.5394598325700358
```



Affinity Propagation Clustering

Affinity Propagation, instead, takes as input measures of similarity between pairs of data points, and simultaneously considers all data points as potential examples. Real-valued messages are exchanged between data points until a high-quality set of examples and corresponding clusters gradually emerges.

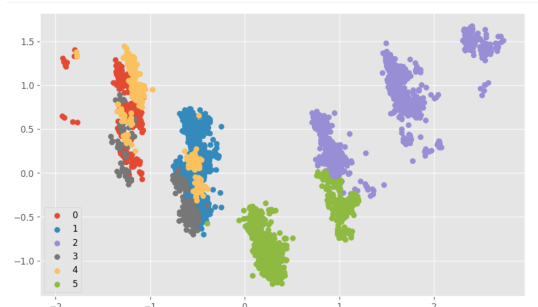
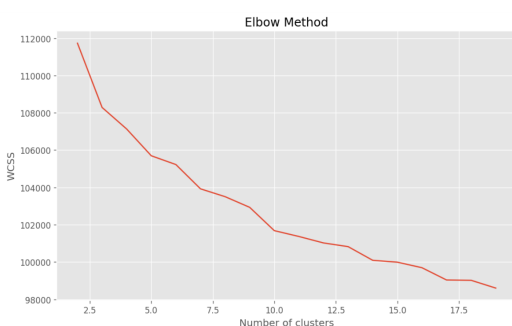


PART B: Modeling with CountVectorizer and TfidfVectorizer

K Mean with CountVectorizer and TfidfVectorizer

TF-IDF stands for term frequency-inverse document frequency and it is a measure, used in the fields of information retrieval (IR) and machine learning, that can quantify the importance or relevance of string representations (words, phrases, lemmas, etc) in a

document amongst a collection of documents. CountVectorizer means breaking down a sentence or any text into words by performing preprocessing tasks like converting all words to lowercase, thus removing special characters. In NLP models can't understand textual data they only accept numbers, so this textual data needs to be vectorized.

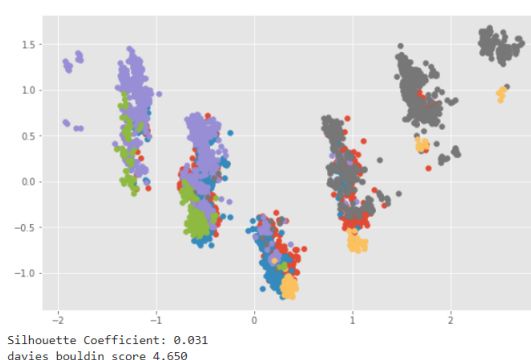
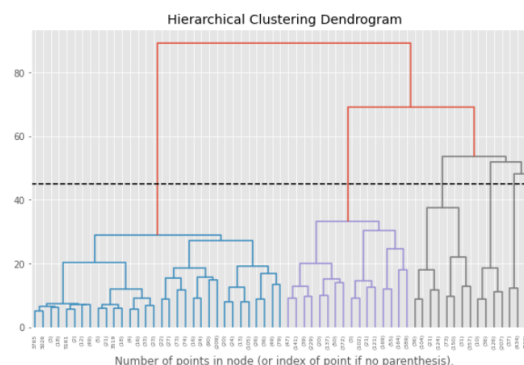


Agglomerative Clustering with CountVectorizer and TfidfVectorizer

For n_clusters = 4 The average silhouette_score is : 0.028335436335219656
 For n_clusters = 5 The average silhouette_score is : 0.02866348585584631
 For n_clusters = 6 The average silhouette_score is : 0.0308242009559484
 For n_clusters = 7 The average silhouette_score is : 0.029358281176396637
 For n_clusters = 8 The average silhouette_score is : 0.03139167971810257
 For n_clusters = 9 The average silhouette_score is : 0.02934520186084311



Agglomerative Clustering with Dendrogram and CountVectorizer



7. Recommendation System

Cosine similarity measures the similarity between two vectors of an inner product space. It is measured by the cosine of the angle between two vectors and determines whether two vectors are pointing in roughly the same direction. It is often used to measure document similarity in text analysis

Recommendations	
0	Charlie's Angels: Full Throttle
1	Malibu Rescue: The Series
2	Sex, Explained
3	Dynasty
4	The Dukes of Hazzard
5	The Who Was? Show
6	The Legend of 420
7	The Seventies
8	DreamWorks How to Train Your Dragon Legends
9	Hellboy

Summary

Clustering Models with word2vec

SL No.	Model Name	Optimal Number of cluster
1	K-Means with silhouette_score with word2vec	6
2	K-Means with Elbow method with word2vec	5
3	Agglomerative Clustering with dendrogram with word2vec	5
4	Agglomerative Clustering with silhouette_score with word2vec	4
5	Affinity propagation clustering with word2vec	5

Clustering Models with CountVectorizer

SL No.	Model Name	Optimal Number of cluster
1	K-Means with Elbow method with countvectorizer	6
3	Agglomerative Clustering with dendrogram with countvectorizer	6
4	Agglomerative Clustering with silhouette_score with countvectorizer	6

8. Conclusion

Exploratory Data Analysis

1. The attribute 'director', 'cast', 'country', 'date_added', 'rating' consists of missing values. To tackle missing values, we will replace 'country' and 'rating' missing values by the most frequent entity that is 'United States' and 'TV-MA' respectively. missing values in 'cast' by 'unknown'. There are around 30.68 % values missing in 'director', hence we decide to drop it. 69% of the content available on Netflix are movies; the remaining 31% are TV Shows.
2. Netflix has 5377 movies, which is more than double

the quantity of TV shows. In recent years more TV Shows are released as compared to Movies on Netflix. Less number of TV shows and Movies were released in 2020-2021 due to the coronavirus pandemic. Most of the Movies/TV Shows were added in the month of December and January.

3. Number of Movies added on Netflix are more as compared to TV Shows throughout the year. In recent few year more number of TV Shows were added on NetFlix as compared to Movies, We can say Netflix is more focusing on TV Shows than Movies.
4. The United States, India, United Kingdom, Japan, South Korea, Canada, Spain, France, Egypt and Turkey are the Top 10 countries which produce most of the content on Netflix. The United States produced most of the content on Netflix. Also, the number of movies released are more than TV Shows in the United States. In India, Canada, Spain, France, Egypt and Turkey, Most of the content on Netflix is Movies. The United Kingdom has almost equal production of Movies and TV Shows. In Japan and South Korea, Number of TV Shows are available on Netflix.
5. It is observed that, in each category, Quantity of Movies is more than the Quantity of

TV Shows. The Availability of the Adult Content is more on Netflix and Least for the Kids.

6. Popular Movies ratings are TV-MA, TV-14, R, TV-PG, PG-14 and PG. It is observed that Adults and Teens are mostly active on Netflix. Popular TV Shows ratings are TV-MA, TV-14, R, TV-PG, PG-14 and PG. Top 5 Genres in 'TV Shows' are Kid's TV, TV Dramas, TV Crime Shows, TV Comedies, TV Romantic. Top 5 Genres in 'Movies' are Documentaries, Stand up Comedy, Dramas and International Movies, Comedies and Independent Movies. It is observed that 1608 TV Shows has only one season. The count of longest running TV Shows is very less.
7. Famous Actors on Netflix based on the Frequency of their occurrence on screen are Anupam Kher, Takahiro Sakurai, Shah Rukh Khan, Om Puri and Boman Irani and so on. Most of the Movies/TV Shows have a duration of around 100 min. The United States produces maximum International TV Shows, TV Dramas, Sci-fi and Fantasy TV shows, International Movies. India, UK, Spain, Egypt, Mexico and Turkey are having most of the Content as Dramas and International Movies.
8. It is observed that content available for kids is less as compared to other categories. Content available for Adults is more in almost

every country except India. In India, Most of the content is available for Teens. Netflix should focus on Teen and Adult Contents to generate maximum revenue. Spain and Mexico are producing the highest Adult Content on Netflix almost 84% and 77% respectively.

Clustering with Word2vec

1. K-Means with 0.6092 silhouette_score with word2vec has an optimum number of clusters as 6.
2. K-Means with Elbow method with word2vec has 5 optimum clusters.
3. Agglomerative Clustering with dendrogram with word2vec has 5 optimum clusters.
4. Agglomerative Clustering with 0.53 silhouette_score with word2vec gives 4 clusters.
5. Affinity propagation clustering with word2vec has 5 optimum clusters.

Clustering with CountVectorizer

1. It is observed that , after using CountVectorizer and tfidfVectorizer, we get the less silhouette_score as 0.032
2. Hence we can say word2vec word embedding method is more suitable for our model.

Winner Model

I have used Principal component Analysis for feature reduction,

Recommended System is also designed for getting recommendation of movies and TV Shows. Hyperparameter Tuning is done in every model to get optimum results.

K-Means with word2vec with 6 optimum clusters with 0.6092 silhouette_score

9. Future Work

We can use different Word2vec vectors for further analysis. Some other methods of clustering like Density-based, Distribution-based, Centroid-based, DBSCAN clustering algorithm, Gaussian Mixture Model algorithm, BIRCH algorithm.

References:

1. https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.CountVectorizer.html
2. https://scikit-learn.org/stable/modules/generated/sklearn.metrics.silhouette_score.html
3. https://scikit-learn.org/stable/auto_examples/cluster/plot_agglomerative_dendrogram.html
4. <https://help.netflix.com/en/node/2064>
5. <https://scikit-learn.org/stable/modules/generated/sklearn.cluster.AgglomerativeClustering.html>
6. <https://scikit-learn.org/stable/modules/generated/sklearn.cluster.AffinityPropagation.html>

Thank You.....!