

## Assignment 5.2

### Q1) Explain the concept of Tokenization.

**Ans)** Tokenization is the first step in any NLP pipeline. It has an important effect on the rest of your pipeline. A tokenizer breaks unstructured data and natural language text into chunks of information that can be considered discrete elements. The token occurrences in a document can be used directly as a vector representing that document. This immediately turns an unstructured string (text document) into a numerical data structure suitable for machine learning. They can also be used directly by a computer to trigger useful actions and responses. Or they might be used in a machine learning pipeline as features that trigger more complex decisions or behavior. Tokenization can separate sentences, words, characters, or subwords. When we split the text into sentences, we call it sentence tokenization. For words, we call it word tokenization.

### Example of sentence tokenization

```
sent_tokenize('Life is a matter of choices, and every choice you make makes you.')  
['Life is a matter of choices, and every choice you make makes you.']
```

### Example of word tokenization

```
word_tokenize("The sole meaning of life is to serve humanity")  
['The', 'sole', 'meaning', 'of', 'life', 'is', 'to', 'serve', 'humanity']
```

### Q2) How and when is Gram tokenization is used?

**Ans)** The Gram tokenization or N-gram tokenizer first breaks text down into words whenever it encounters one of a list of specified characters, then it emits N-grams of each word of the specified length. N-grams are like a sliding window that moves across the word - a continuous sequence of characters of the specified length. They are useful for querying languages that don't use spaces or that have long compound words, like German. Example: Let's understand n-grams practically with the help of the following sentence:

“I reside in Bengaluru”

SL.No.	Type of n-gram	Generated n-grams
1	Unigram	[“I”, “reside”, “in”, “Bengaluru”]
2	Bigram	[“I reside”, “reside in”, “in Bengaluru”]
3	Trigram	[“I reside in”, “reside in Bengaluru”]

From the table above, it's clear that unigram means taking only one word at a time, bigram means taking two words at a time and trigram means taking three words at a time

**Q3) What is meant by the TFIDF? Explain in detail.**

**Ans)** TFIDF works by proportionally increasing the number of times a word appears in the document but is counterbalanced by the number of documents in which it is present. Hence, words like ‘this’, ‘are’ etc., that are commonly present in all the documents are not given a very high rank. However, a word that is present too many times in a few of the documents will be given a higher rank as it might be indicative of the context of the document.

*Term Frequency:* Term frequency is defined as the number of times a word (i) appears in a document (j) divided by the total number of words in the document.

$$tf_{i,j} = \frac{n_{i,j}}{\sum_k n_{i,k}}$$

*Inverse Document Frequency:* Inverse document frequency refers to the log of the total number of documents divided by the number of documents that contain the word. The logarithm is added to dampen the importance of a very high value of IDF.

$$idf(w) = \log\left(\frac{N}{df_t}\right)$$

TFIDF is computed by multiplying the term frequency with the inverse document frequency.

$$w_{i,j} = tf_{i,j} \times \log\left(\frac{N}{df_i}\right)$$

Let us now see an illustration of TFIDF in the following sentences, which we refer to as documents.

**Document 1:** Text processing is necessary.

**Document 2:** Text processing is necessary and important.

Word	TF		IDF	TFIDF	
	Doc 1	Doc 2		Doc 1	Doc 2
Text	1/4	1/6	$\log(2/2) = 0$	0	0
Processing	1/4	1/6	$\log(2/2) = 0$	0	0
Is	1/4	1/6	$\log(2/2) = 0$	0	0
Necessary	1/4	1/6	$\log(2/2) = 0$	0	0
And	0/4	1/6	$\log(2/1) = 0.3$	0	0.05
Important	0/4	1/6	$\log(2/1) = 0.3$	0	0.05

The above table shows how the TFIDF of some words are zero and some words are non-zero depending on their frequency in the document and across all documents. The limitation of TFIDF is again that this vectorization doesn't help in bringing in the contextual meaning of the words as it is just based on the frequency.

**Q4) Perform NLP for the netflix\_titles.csv. Download netflix\_titles.csv. The dataset has a set of reviews for the shows/movies available on the platform.**

- Perform EDA of the dataset. Conclude on the relationship between the various features
- Conclude the most important features
- Identify similar content by matching text-based features
- Conclude whether Netflix is increasingly focusing on TV shows rather than movies

**Ans)** The answer is on the next page.

# 1. Perform EDA of the dataset. Conclude on the relationship between the various features

Netflix is a popular service that people across the world use for entertainment. In this EDA, I will explore the netflix-shows dataset through visualizations and graphs using matplotlib and seaborn.

## Package Import

First, we will import necessary packages.

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
import matplotlib
```

## Loading the Dataset

Now we are ready to load the dataset. We will do this using the standard `read_csv` command from Pandas. Let's take a glimpse at how the data looks like.

In [2]: `netflix_titles_df = pd.read_csv('netflix_titles.csv')  
netflix_titles_df.head()`

Out[2]:

	show_id	type	title	director	cast	country	date_added	release_year	rating
0	81145628	Movie	Norm of the North: King Sized Adventure	Richard Finn, Tim Maltby	Alan Marriott, Andrew Toth, Brian Dobson, Cole...	United States, India, South Korea, China	September 9, 2019	2019	TV-PG
1	80117401	Movie	Jandino: Whatever it Takes	NaN	Jandino Asporaat	United Kingdom	September 9, 2016	2016	TV-MA
2	70234439	TV Show	Transformers Prime	NaN	Peter Cullen, Sumalee Montano, Frank Welker, J...	United States	September 8, 2018	2013	TV-Y7-FV
3	80058654	TV Show	Transformers: Robots in Disguise	NaN	Will Friedle, Darren Criss, Constance Zimmer, ...	United States	September 8, 2018	2016	TV-Y7
4	80125979	Movie	#realityhigh	Fernando Lebrija	Nesta Cooper, Kate Walsh, John Michael Higgins...	United States	September 8, 2017	2017	TV-14

After a quick glimpse at the dataset, it looks like a typical movies/shows dataset without user ratings. We can also see that there are NaN values in some columns.

## Data Preparation and Cleaning

In [3]: `netflix_titles_df.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6234 entries, 0 to 6233
Data columns (total 12 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   show_id     6234 non-null    int64  
 1   type        6234 non-null    object  
 2   title       6234 non-null    object  
 3   director    4265 non-null    object  
 4   cast         5664 non-null    object  
 5   country     5758 non-null    object  
 6   date_added  6223 non-null    object  
 7   release_year 6234 non-null    int64  
 8   rating      6224 non-null    object  
 9   duration    6234 non-null    object  
 10  listed_in   6234 non-null    object  
 11  description 6234 non-null    object  
dtypes: int64(2), object(10)
memory usage: 584.6+ KB
```

There are 6,234 entries and 12 columns to work with for EDA. Right off the bat, there are a few columns that contain null values ('director', 'cast', 'country', 'date\_added', 'rating').

In [4]: `netflix_titles_df.nunique()`

```
Out[4]: show_id      6234
type          2
title         6172
director     3301
cast          5469
country       554
date_added   1524
release_year  72
rating        14
duration     201
listed_in    461
description   6226
dtype: int64
```

## Handling Null Values

We can see that for each of the columns, there are a lot of different unique values for some of them. It makes sense that `show_id` is large since it is a unique key used to identify a movie/show. Title, director, cast, country, date\_added, listed\_in, and description contain many unique values as well.

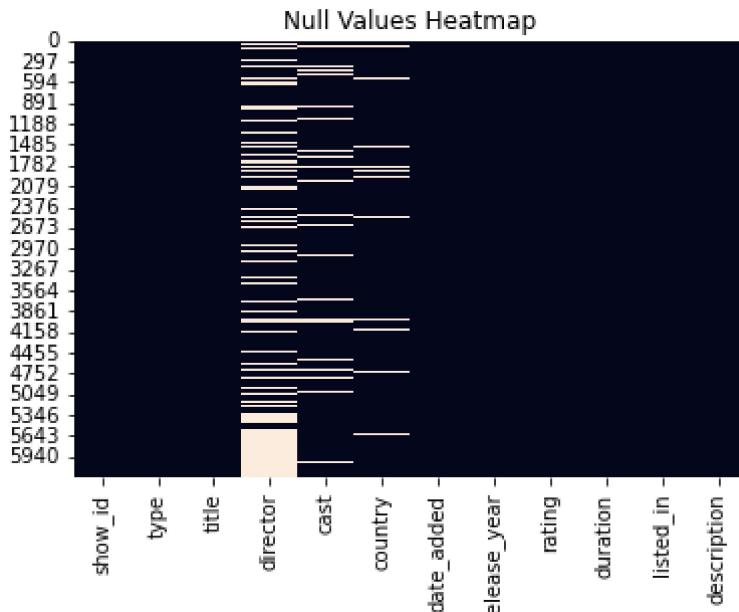
In [5]: `netflix_titles_df.isnull().values.any()`

Out[5]: True

In [6]: `netflix_titles_df.isnull().sum().sum()`

Out[6]: 3036

In [7]: `sns.heatmap(netflix_titles_df.isnull(), cbar=False)  
plt.title('Null Values Heatmap')  
plt.show()`



In [8]: `netflix_titles_df.isnull().sum()`

Out[8]:

show_id	0
type	0
title	0
director	1969
cast	570
country	476
date_added	11
release_year	0
rating	10
duration	0
listed_in	0
description	0

dtype: int64

Above in the heatmap and table, we can see that there are quite a few null values in the dataset. There are a total of 3,036 null values across the entire dataset with 1,969 missing points under 'director', 570 under 'cast', 476 under 'country', 11 under 'date\_added', and 10 under 'rating'. We will have to handle all null data points before we can dive into EDA and modeling.

```
In [9]: netflix_titles_df['director'].fillna('No Director', inplace=True)
netflix_titles_df['cast'].fillna('No Cast', inplace=True)
netflix_titles_df['country'].fillna('Country Unavailable', inplace=True)
netflix_titles_df.dropna(subset=['date_added', 'rating'], inplace=True)
```

```
In [10]: netflix_titles_df.isnull().any()
```

```
Out[10]: show_id      False
type          False
title         False
director      False
cast          False
country       False
date_added    False
release_year  False
rating        False
duration      False
listed_in     False
description   False
dtype: bool
```

For null values, the easiest way to get rid of them would be to delete the rows with the missing data. However, this wouldn't be beneficial to our EDA since there is loss of information. Since 'director', 'cast', and 'country' contain the majority of null values, I will choose to treat each missing value as unavailable. The other two labels 'date\_added' and 'rating' contains an insignificant portion of the data so I will drop them from the dataset. After, we can see that there are no more null values in the dataset.

## Splitting the Dataset

Since the dataset can either contain movies or shows, it'd be nice to have datasets for both so we can take a deep dive into just Netflix movies or Netflix TV shows so we will create two new datasets. One for movies and the other one for shows.

In [11]: `netflix_movies_df = netflix_titles_df[netflix_titles_df['type']=='Movie'].copy()  
netflix_movies_df.head()`

Out[11]:

	show_id	type	title	director	cast	country	date_added	release_year	rating
0	81145628	Movie	Norm of the North: King Sized Adventure	Richard Finn, Tim Maltby	Alan Marriott, Andrew Toth, Brian Dobson, Cole...	United States, India, South Korea, China	September 9, 2019	2019	TV-PG
1	80117401	Movie	Jandino: Whatever it Takes	No Director	Jandino Asporaat	United Kingdom	September 9, 2016	2016	TV-MA
4	80125979	Movie	#realityhigh	Fernando Lebrija	Nesta Cooper, Kate Walsh, John Michael Higgins...	United States	September 8, 2017	2017	TV-14
6	70304989	Movie	Automata	Gabe Ibáñez	Antonio Banderas, Dylan McDermott, Melanie Gri...	Bulgaria, United States, Spain, Canada	September 8, 2017	2014	R
7	80164077	Movie	Fabrizio Copano: Solo pienso en mi	Rodrigo Toro, Francisco Schultz	Fabrizio Copano	Chile	September 8, 2017	2017	TV-MA

In [12]: `netflix_shows_df = netflix_titles_df[netflix_titles_df['type']=='TV Show'].copy()  
netflix_shows_df.head()`

Out[12]:

	show_id	type	title	director	cast	country	date_added	release_year	r
2	70234439	TV Show	Transformers Prime	No Director	Peter Cullen, Sumalee Montano, Frank Welker, J...	United States	September 8, 2018	2013	Y
3	80058654	TV Show	Transformers: Robots in Disguise	No Director	Will Friedle, Darren Criss, Constance Zimmer, ...	United States	September 8, 2018	2016	T
5	80163890	TV Show	Apaches	No Director	Alberto Ammann, Eloy Azorin, Verónica Echegui,...	Spain	September 8, 2017	2016	
8	80117902	TV Show	Fire Chasers	No Director	No Cast	United States	September 8, 2017	2017	
26	80244601	TV Show	Castle of Stars	No Director	Chaiyapol Pupart, Jintanutda Lummakanon, Worra...	Country Unavailable	September 7, 2018	2015	T



## Data Preparation

In the duration column, there appears to be a discrepancy between movies and shows. Movies are based on the duration of the movie and shows are based on the number of seasons. To make EDA easier, I will convert the values in these columns into integers for both the movies and shows datasets.

In [13]: `netflix_movies_df.duration = netflix_movies_df.duration.str.replace(' min', '').as  
netflix_shows_df.rename(columns={'duration':'seasons'}, inplace=True)  
netflix_shows_df.replace({'seasons':{'1 Season':'1 Seasons'}}, inplace=True)  
netflix_shows_df.seasons = netflix_shows_df.seasons.str.replace(' Seasons', '').as`

## Exploratory Analysis and Visualization

First we will begin analysis on the entire Netflix dataset consisting of both movies and shows.

Revisiting the data, let us see how it looked like again.

In [14]: `netflix_titles_df.head()`

Out[14]:

	show_id	type	title	director	cast	country	date_added	release_year	rating
0	81145628	Movie	Norm of the North: King Sized Adventure	Richard Finn, Tim Maltby	Alan Marriott, Andrew Toth, Brian Dobson, Cole...	United States, India, South Korea, China	September 9, 2019	2019	TV-PG
1	80117401	Movie	Jandino: Whatever it Takes	No Director	Jandino Asporaat	United Kingdom	September 9, 2016	2016	TV-MA
2	70234439	TV Show	Transformers Prime	No Director	Peter Cullen, Sumalee Montano, Frank Welker, J...	United States	September 8, 2018	2013	TV-Y7-FV
3	80058654	TV Show	Transformers: Robots in Disguise	No Director	Will Friedle, Darren Criss, Constance Zimmer, ...	United States	September 8, 2018	2016	TV-Y7
4	80125979	Movie	#realityhigh	Fernando Lebrija	Nesta Cooper, Kate Walsh, John Michael Higgins...	United States	September 8, 2017	2017	TV-14

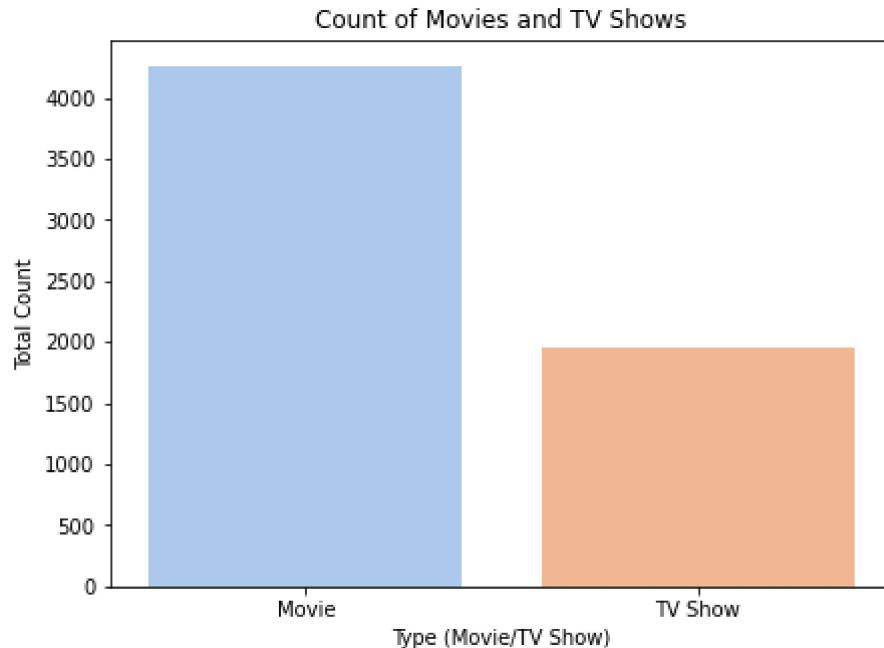
## Netflix Film Types: Movie or TV Show

It'd be interesting to see the comparison between the total number of movies and shows in this dataset just to get an idea of which one is the majority.

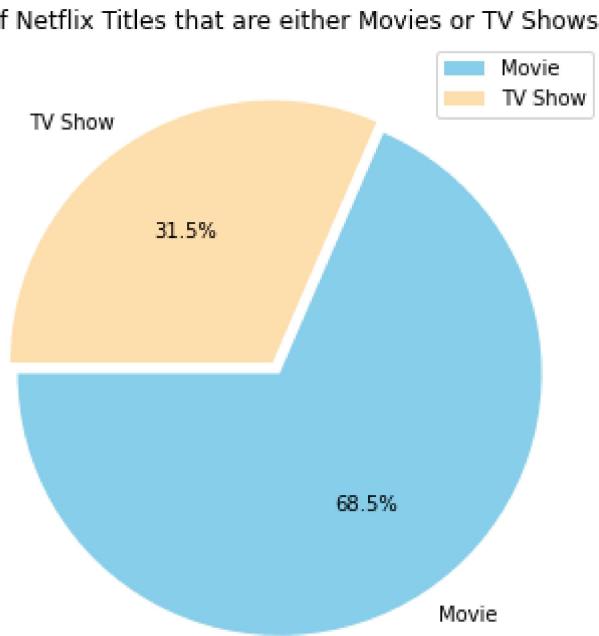
```
In [15]: plt.figure(figsize=(7,5))
g = sns.countplot(netflix_titles_df.type, palette="pastel");
plt.title("Count of Movies and TV Shows")
plt.xlabel("Type (Movie/TV Show)")
plt.ylabel("Total Count")
plt.show()
```

C:\anaconda3\lib\site-packages\seaborn\\_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.

```
warnings.warn(
```



```
In [16]: plt.figure(figsize=(12,6))
plt.title("% of Netflix Titles that are either Movies or TV Shows")
g = plt.pie(netflix_titles_df.type.value_counts(), explode=(0.025,0.025), labels=
plt.legend()
plt.show()
```



So there are roughly 4,000+ movies and almost 2,000 shows with movies being the majority. This makes sense since shows are always an ongoing thing and have episodes. If we were to do a headcount of TV show episodes vs. movies, I am sure that TV shows would come out as the majority. However, in terms of title, there are far more movie titles (68.5%) than TV show titles (31.5%).

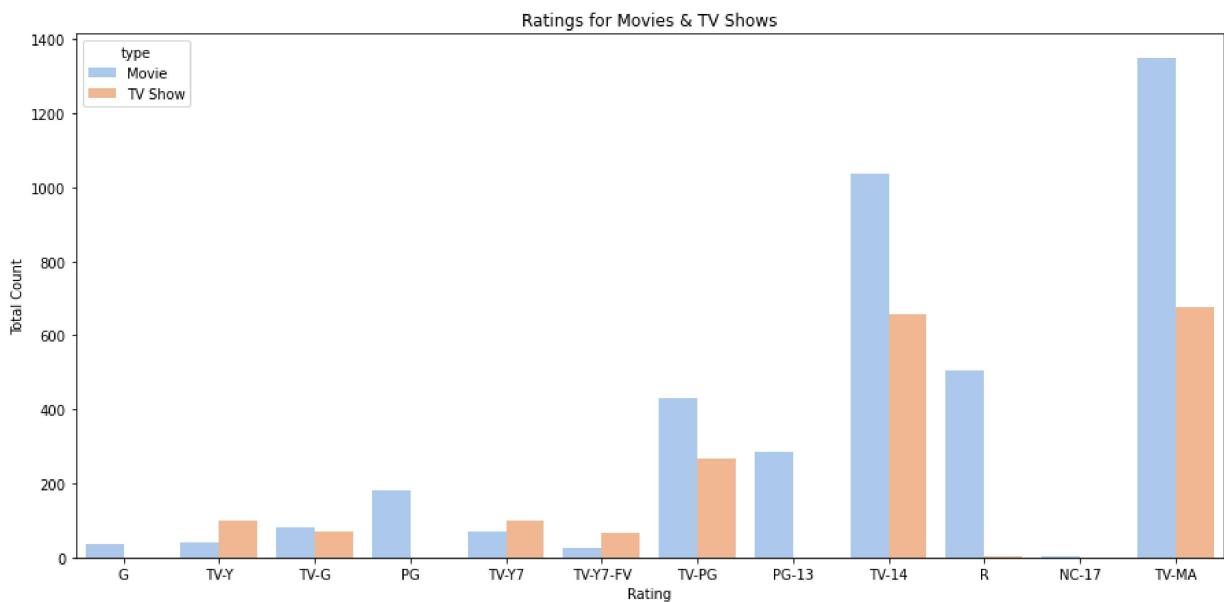
## Netflix Film Ratings

Now, we will explore the ratings which are based on the film rating system. The ordering of the ratings will be based on the age of the respective audience from youngest to oldest. We will not include the ratings 'NR' and 'UR' in the visuals since they stand for unrated and non-rated content.

```
In [17]: order = ['G', 'TV-Y', 'TV-G', 'PG', 'TV-Y7', 'TV-Y7-FV', 'TV-PG', 'PG-13', 'TV-14', 'R', 'NC-17', 'TV-MA']
plt.figure(figsize=(15,7))
g = sns.countplot(netflix_titles_df.rating, hue=netflix_titles_df.type, order=order)
plt.title("Ratings for Movies & TV Shows")
plt.xlabel("Rating")
plt.ylabel("Total Count")
plt.show()
```

C:\anaconda3\lib\site-packages\seaborn\\_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.

```
warnings.warn(
```



In [18]:

```
fig, ax = plt.subplots(1,2, figsize=(19, 5))
g1 = sns.countplot(netflix_movies_df.rating, order=order, palette="Set2", ax=ax[0])
g1.set_title("Ratings for Movies")
g1.set_xlabel("Rating")
g1.set_ylabel("Total Count")
g2 = sns.countplot(netflix_shows_df.rating, order=order, palette="Set2", ax=ax[1])
g2.set(yticks=np.arange(0,1600,200))
g2.set_title("Ratings for TV Shows")
g2.set_xlabel("Rating")
g2.set_ylabel("Total Count")
fig.show()
```

C:\anaconda3\lib\site-packages\seaborn\\_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.

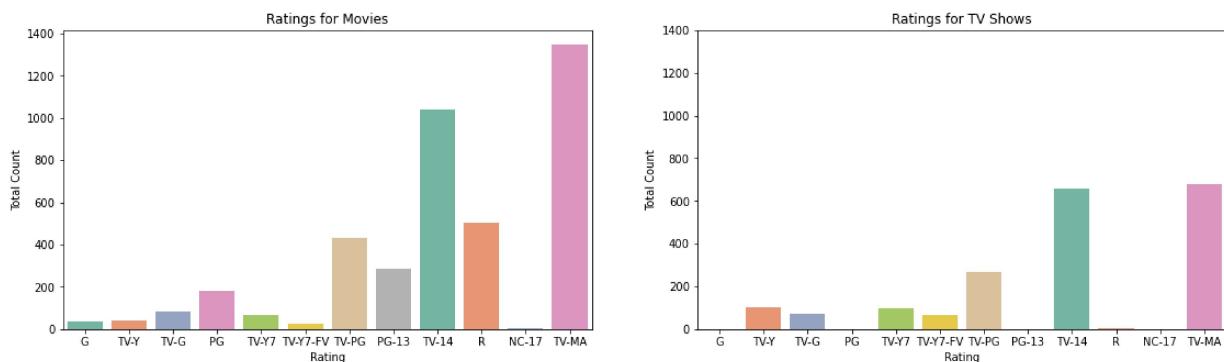
warnings.warn(

C:\anaconda3\lib\site-packages\seaborn\\_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.

warnings.warn(

C:\Users\Shaina Mehta\AppData\Local\Temp\ipykernel\_15356\1913264221.py:11: User Warning: Matplotlib is currently using module://matplotlib\_inline.backend\_inline, which is a non-GUI backend, so cannot show the figure.

fig.show()



Overall, there is much more content for a more mature audience. For the mature audience, there is much more movie content than there are TV shows. However, for the younger audience (under the age of 17), it is the opposite, there are slightly more TV shows than there are movies.

In [19]:

```
netflix_titles_df['year_added'] = pd.DatetimeIndex(netflix_titles_df['date_added'])
netflix_movies_df['year_added'] = pd.DatetimeIndex(netflix_movies_df['date_added'])
netflix_shows_df['year_added'] = pd.DatetimeIndex(netflix_shows_df['date_added'])
netflix_titles_df['month_added'] = pd.DatetimeIndex(netflix_titles_df['date_added'])
netflix_movies_df['month_added'] = pd.DatetimeIndex(netflix_movies_df['date_added'])
netflix_shows_df['month_added'] = pd.DatetimeIndex(netflix_shows_df['date_added'])
```

## Content added each year

Now we will take a look at the amount content Netflix has added throughout the previous years.

Since we are interested in when Netflix added the title onto their platform, we will add a

'year\_added' column shows the year of the date from the 'date\_added' column as shown above.

```
In [20]: netflix_year = netflix_titles_df['year_added'].value_counts().to_frame().reset_index()
netflix_year = netflix_year[netflix_year.year != 2020]
netflix_year
```

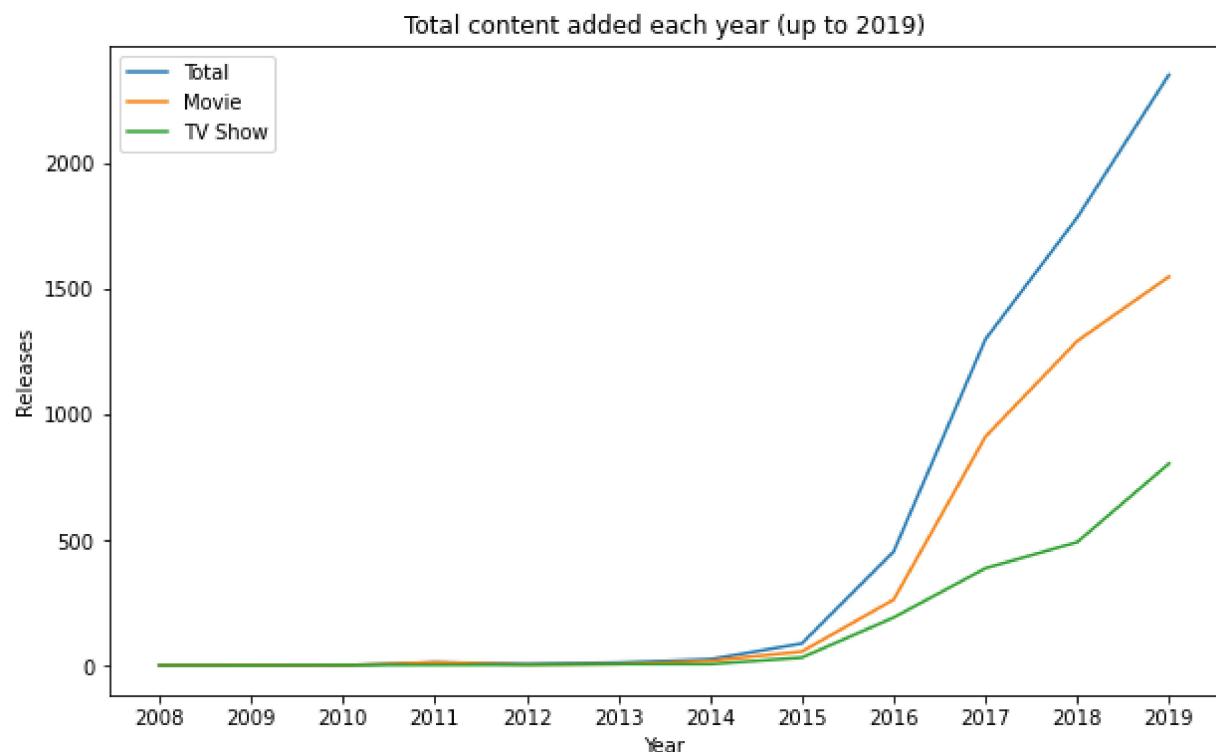
Out[20]:

	year	count
0	2019	2349
1	2018	1781
2	2017	1297
3	2016	453
5	2015	88
6	2014	25
7	2011	13
8	2013	12
9	2012	7
10	2009	2
11	2008	2
12	2010	1

```
In [21]: netflix_year2 = netflix_titles_df[['type', 'year_added']]
movie_year = netflix_year2[netflix_year2['type']=='Movie'].year_added.value_counts()
movie_year = movie_year[movie_year.year != 2020]
show_year = netflix_year2[netflix_year2['type']=='TV Show'].year_added.value_counts()
show_year = show_year[show_year.year != 2020]
```

In [22]:

```
fig, ax = plt.subplots(figsize=(10, 6))
sns.lineplot(data=netflix_year, x='year', y='count')
sns.lineplot(data=movie_year, x='year', y='count')
sns.lineplot(data=show_year, x='year', y='count')
ax.set_xticks(np.arange(2008, 2020, 1))
plt.title("Total content added each year (up to 2019)")
plt.legend(['Total', 'Movie', 'TV Show'])
plt.ylabel("Releases")
plt.xlabel("Year")
plt.show()
```

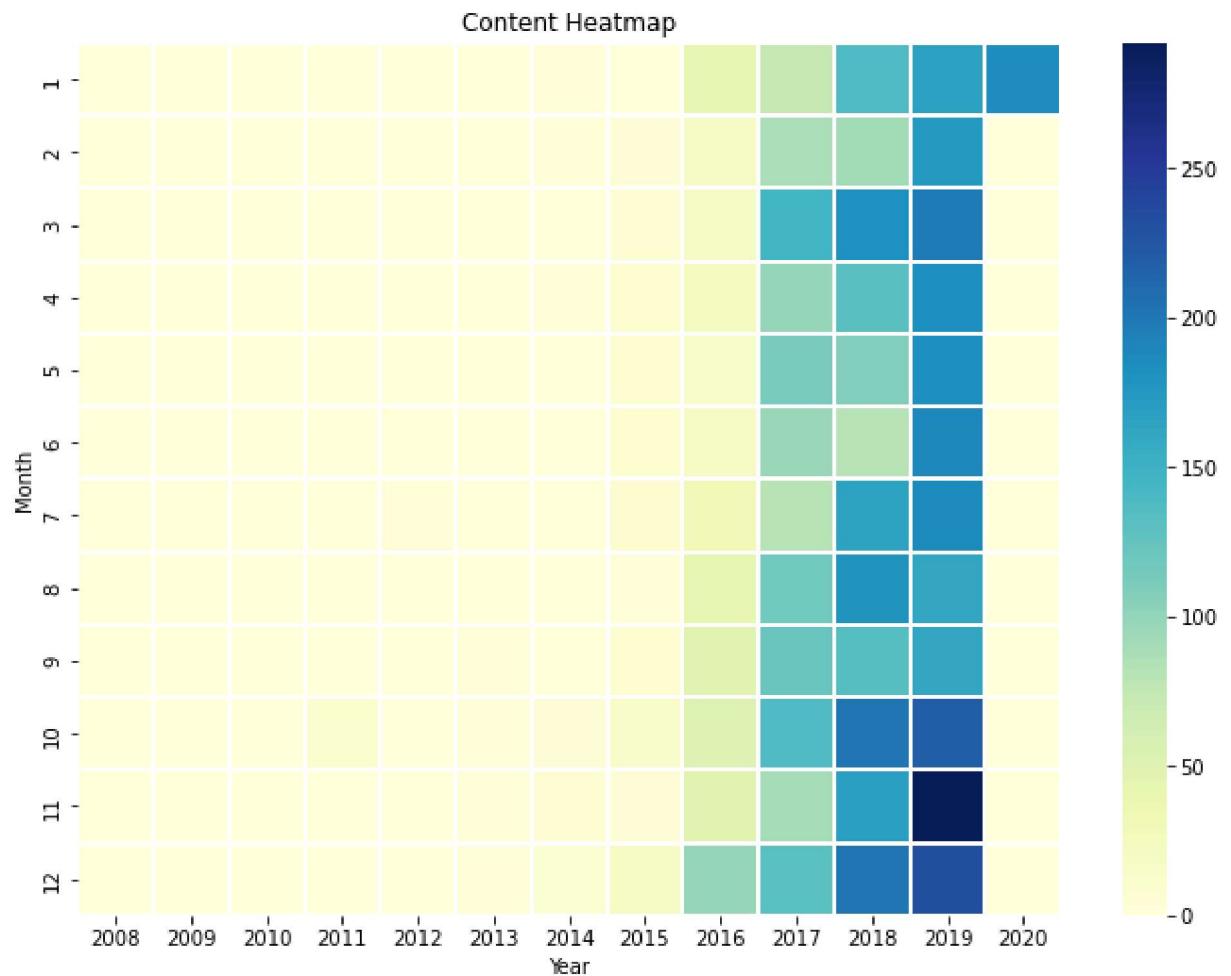


Based on the above timeline, we can see that the popular streaming platform started gaining traction after 2014. Since then, the amount of content added has been tremendous. I decided to exclude content added during 2020 since the data does not include a full years worth of data. We

can see that there has been a consistent growth in the number of movies on Netflix compared to shows.

```
In [23]: month_year_df = netflix_titles_df.groupby('year_added')['month_added'].value_count()

plt.figure(figsize=(11,8))
sns.heatmap(month_year_df, linewidths=0.025, cmap="YlGnBu")
plt.title("Content Heatmap")
plt.ylabel("Month")
plt.xlabel("Year")
plt.show()
```



In the above heatmap, we can see that around 2014 is when Netflix began to increase their content count. We can see over the years and months, Netflix continues to slowly increase the amount of content that is being added into their platform. We can see in 2020, the data stops at January since that is the latest month available in the dataset.

## Netflix Film Duration

In [24]:

```
fig, ax = plt.subplots(1,2, figsize=(19, 5))
g1 = sns.distplot(netflix_movies_df.duration, color='skyblue', ax=ax[0]);
g1.set_xticks(np.arange(0,360,30))
g1.set_title("Duration Distribution for Netflix Movies")
g1.set_ylabel("% of All Netflix Movies")
g1.set_xlabel("Duration (minutes)")
g2 = sns.countplot(netflix_shows_df.seasons, color='skyblue', ax=ax[1]);
g2.set_title("Netflix TV Shows Seasons")
g2.set_ylabel("Count")
g2.set_xlabel("Season(s)")
fig.show()
```

C:\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

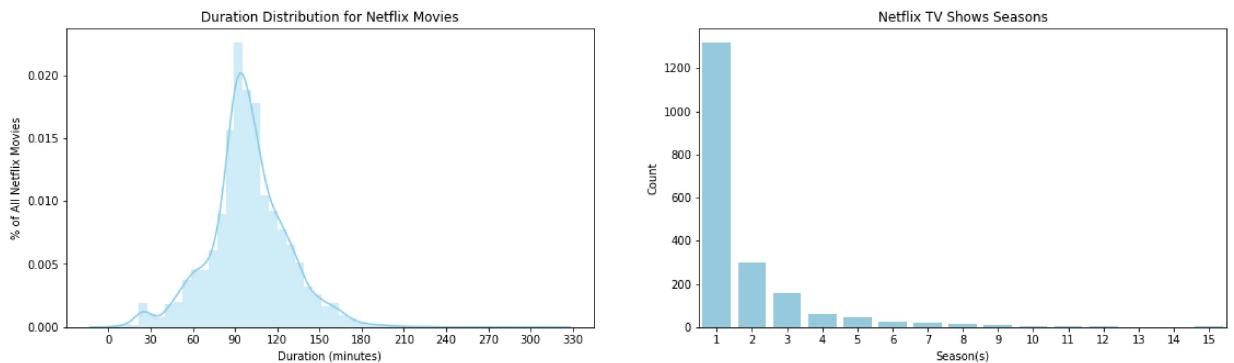
```
warnings.warn(msg, FutureWarning)
```

C:\anaconda3\lib\site-packages\seaborn\\_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.

```
warnings.warn(
```

C:\Users\Shaina Mehta\AppData\Local\Temp\ipykernel\_15356\1788411507.py:11: UserWarning: Matplotlib is currently using module://matplotlib\_inline.backend\_inline, which is a non-GUI backend, so cannot show the figure.

```
fig.show()
```

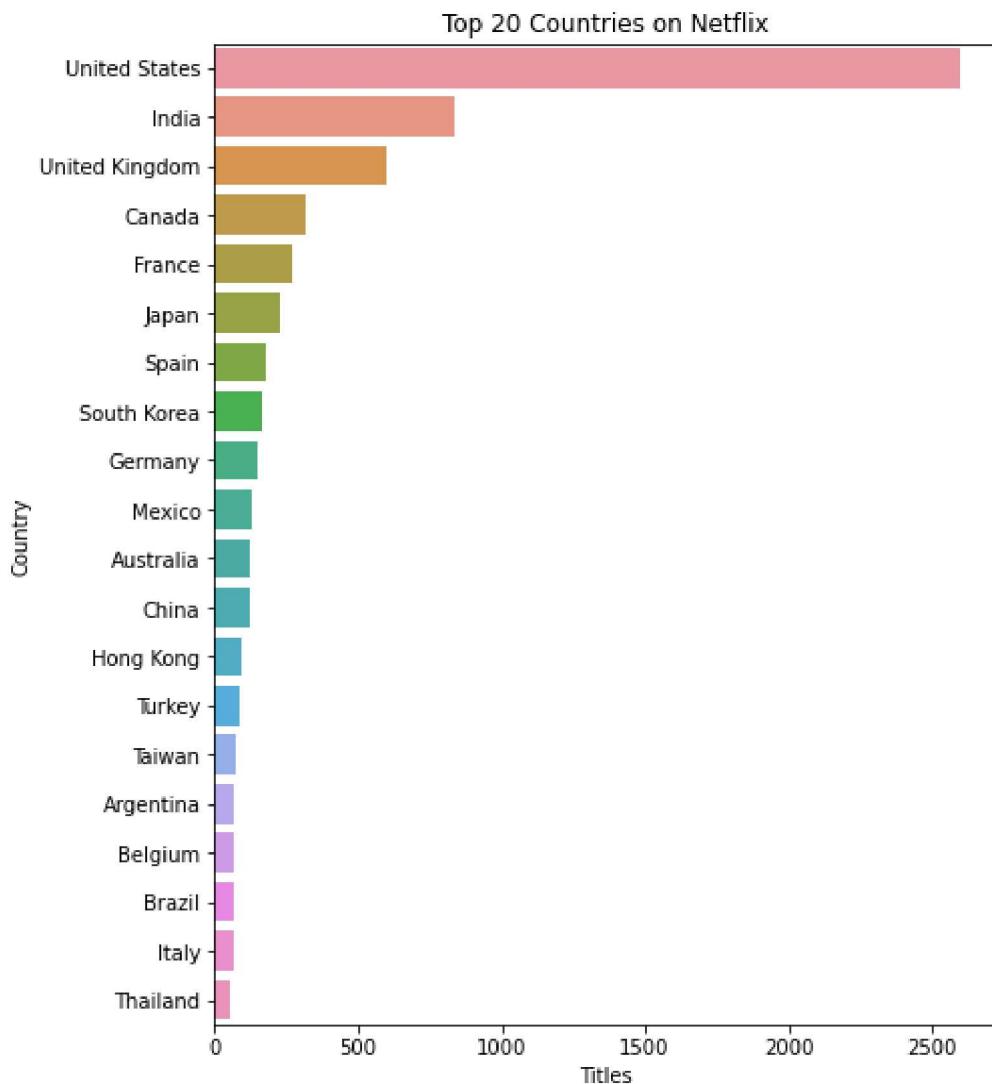


Now we will look into the duration of Netflix films. Since movies are measured in time and shows are measured by seasons, we need to split the dataset between movies and TV shows. Above on the left, we can see that the duration for Netflix movies closely resembles a normal distribution with the average viewing time spanning about 90 minutes which seems to make sense. Netflix TV shows on the other hand seems to be heavily skewed to the right where the majority of shows only have 1 season.

## Countries with the most content available

```
In [25]: filtered_countries = netflix_titles_df.set_index('title').country.str.split(',')
filtered_countries = filtered_countries[filtered_countries != 'Country Unavailable']

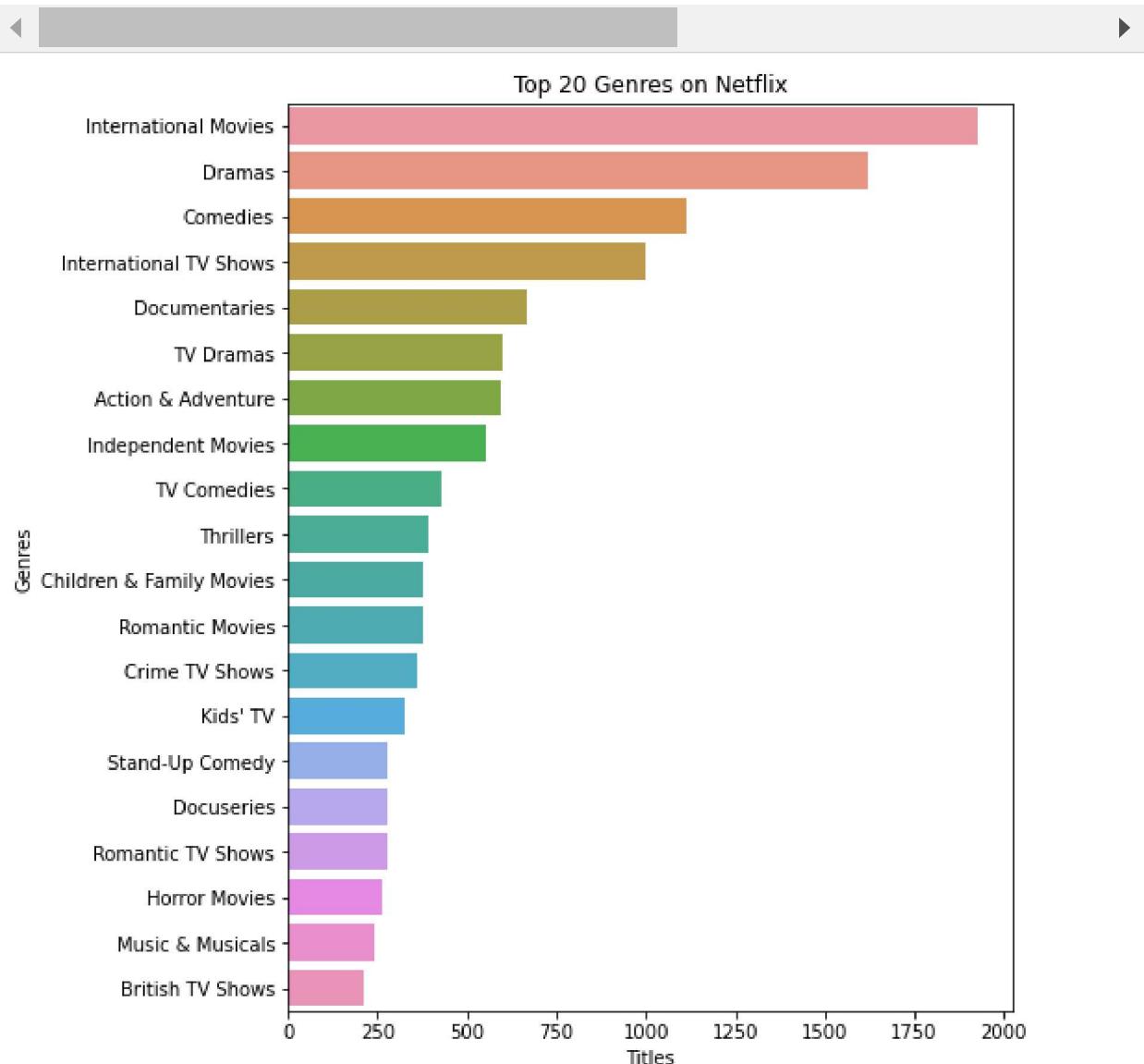
plt.figure(figsize=(7,9))
g = sns.countplot(y = filtered_countries, order=filtered_countries.value_counts())
plt.title('Top 20 Countries on Netflix')
plt.xlabel('Titles')
plt.ylabel('Country')
plt.show()
```



Now we will explore the countries with the most content on Netflix. Films typically are available in multiple countries as shown in the original dataset. Therefore, we need to separate all countries within a film before we can analyze the data. After separating countries and removing titles with no countries available, we can plot a Top 20 list to see which countries have the highest availability of films on Netflix. Unsurprisingly, the United States stands out on top since Netflix is an American company. India surprisingly comes in second followed by the UK and Canada. China interestingly is not even close to the top even though it has about 18% of the world's population. Reasons for this could be for political reasons and the banning of certain applications which isn't uncommon between the United States and China.

## Popular Genres

```
In [26]: filtered_genres = netflix_titles_df.set_index('title').listed_in.str.split(', ',  
  
plt.figure(figsize=(7,9))  
g = sns.countplot(y = filtered_genres, order=filtered_genres.value_counts().index)  
plt.title('Top 20 Genres on Netflix')  
plt.xlabel('Titles')  
plt.ylabel('Genres')  
plt.show()
```



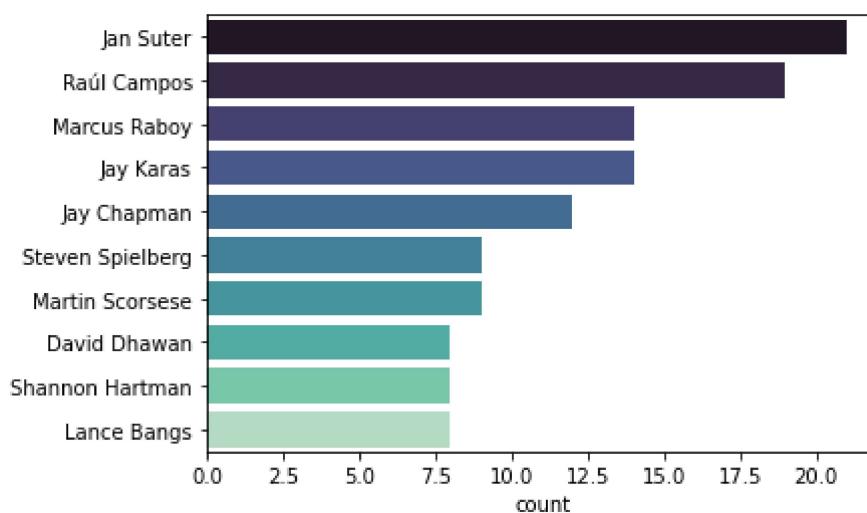
In terms of genres, international movies takes the cake surprisingly followed by dramas and comedies. Even though the United States has the most content available, it looks like Netflix has decided to release a ton of international movies. The reason for this could be that most Netflix subscribers aren't actually in the United States, but rather the majority of viewers are actually international subscribers.

# Asking and Answering Questions

Shaina Mehta  
7CSE 4Y  
A2305219268

**Who are the top 10 directors on Netflix with the most releases?**

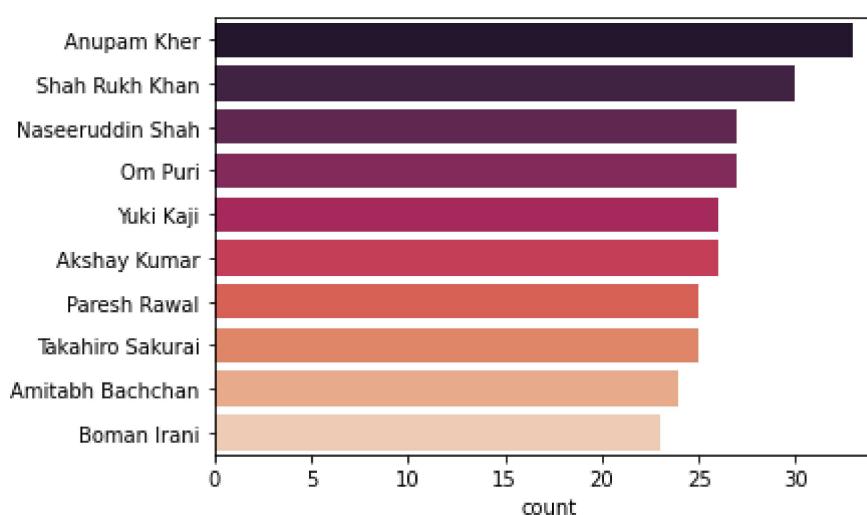
```
In [27]: filtered_directors = netflix_titles_df[netflix_titles_df.director != 'No Director']
sns.countplot(y = filtered_directors, order=filtered_directors.value_counts().index)
plt.show()
```



As stated previously regarding the top genres, it's no surprise that the most popular directors on Netflix with the most titles are mainly international as well.

**Who are the top 10 actors on Netflix based on number of titles?**

```
In [28]: filtered_cast = netflix_titles_df[netflix_titles_df.cast != 'No Cast'].set_index('cast')
sns.countplot(y = filtered_cast, order=filtered_cast.value_counts().index[:10], palette='magma')
plt.show()
```



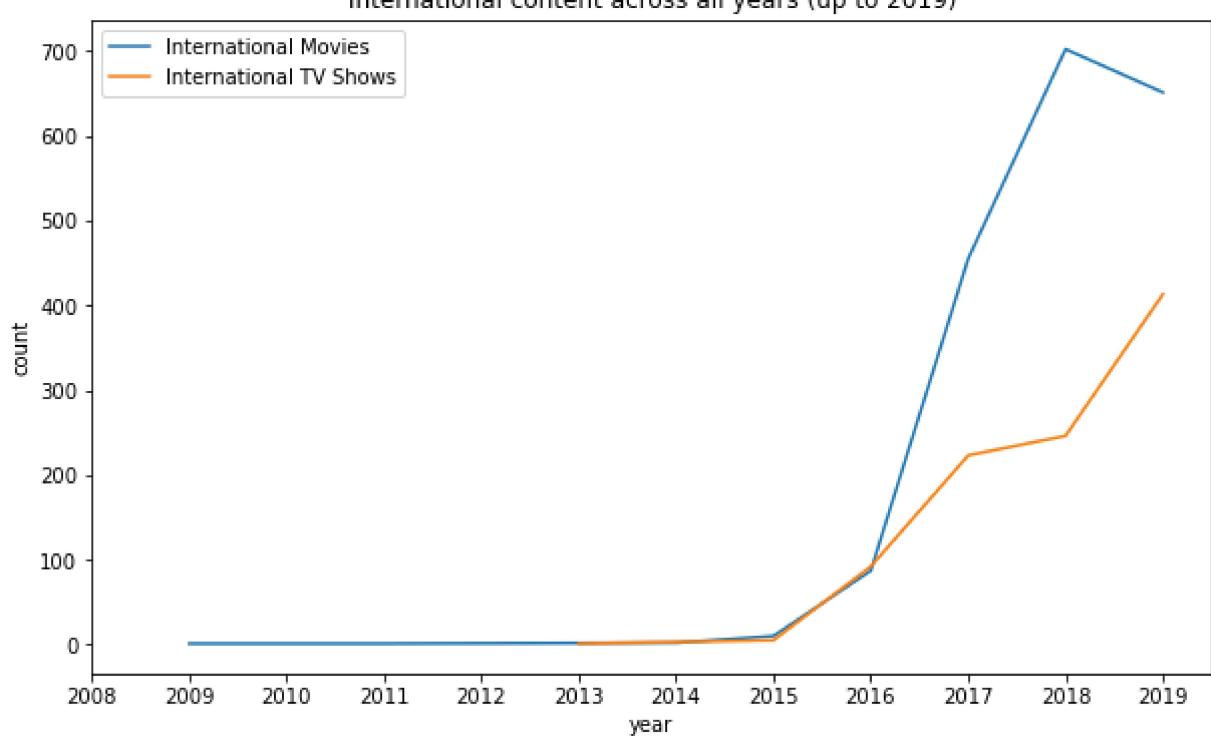
In this list, we can see that the most popular actors on Netflix based on the number of titles are all international as well. This reinforces the sentiment that the majority of Netflix subscribers are international.

## How does the timeline look like for the addition of International Movies compared to International TV Shows?

```
In [29]: international_movies = netflix_titles_df[netflix_titles_df['listed_in'].str.contains('Movie')].value_counts().to_frame().reset_index()
intmov_year = international_movies[['year_added']].value_counts().to_frame().reset_index()
intmov_year = intmov_year[intmov_year.year != 2020]

international_shows = netflix_titles_df[netflix_titles_df['listed_in'].str.contains('TV Show')].value_counts().to_frame().reset_index()
intshow_year = international_shows[['year_added']].value_counts().to_frame().reset_index()
intshow_year = intshow_year[intshow_year.year != 2020]

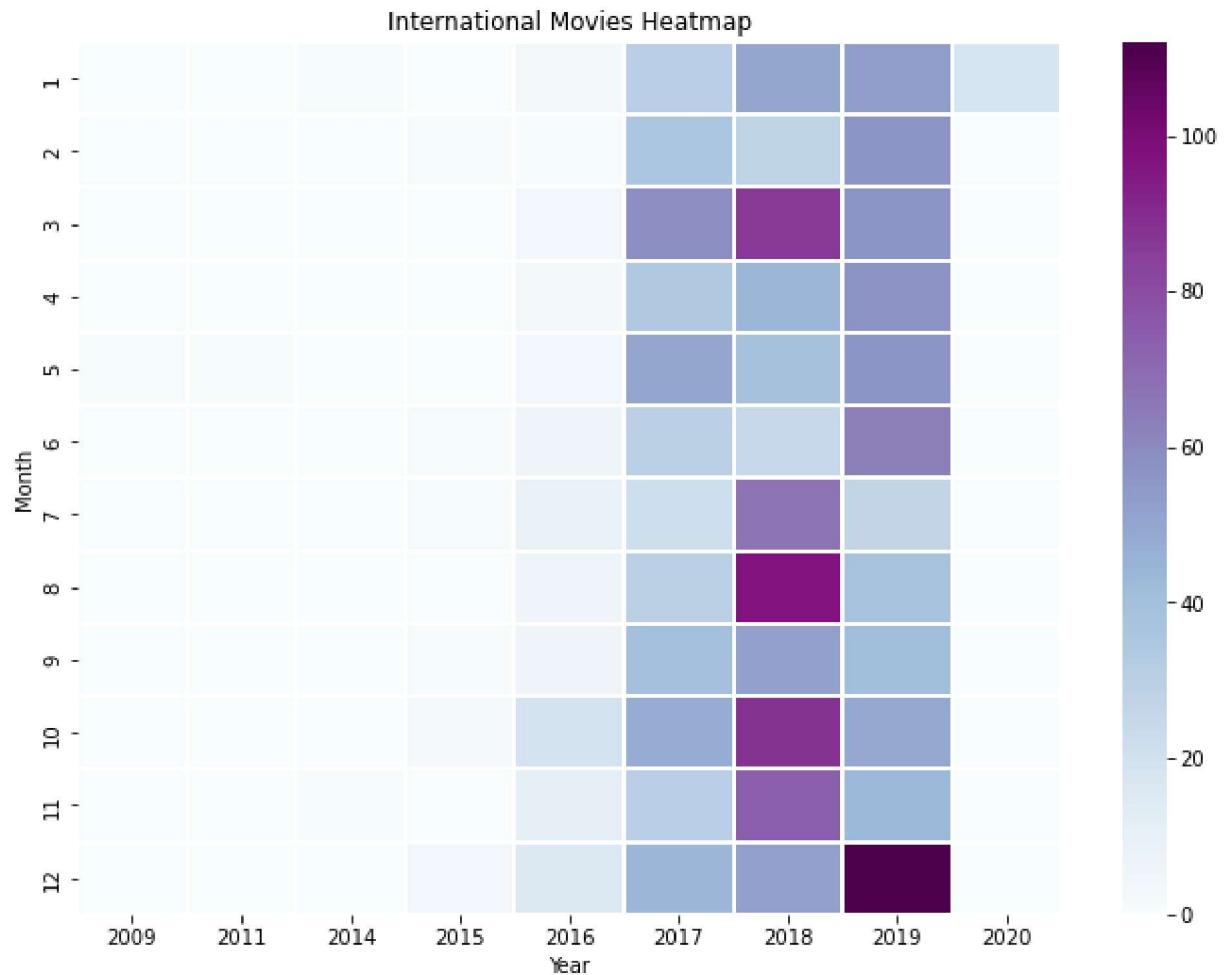
fig, ax = plt.subplots(figsize=(10, 6))
sns.lineplot(data=intmov_year, x='year', y='count')
sns.lineplot(data=intshow_year, x='year', y='count')
ax.set(xticks=np.arange(2008, 2020, 1))
plt.title("International content across all years (up to 2019)")
plt.legend(['International Movies', 'International TV Shows'])
plt.show()
```



Based on the timeline, we can see that there are far more international movie releases than there are international tv show releases. However, near 2018, the growth of international movies started to decline while international tv shows constantly showed significant growth in the past few years.

```
In [30]: intmov_month_year_df = international_movies.groupby('year_added')['month_added']

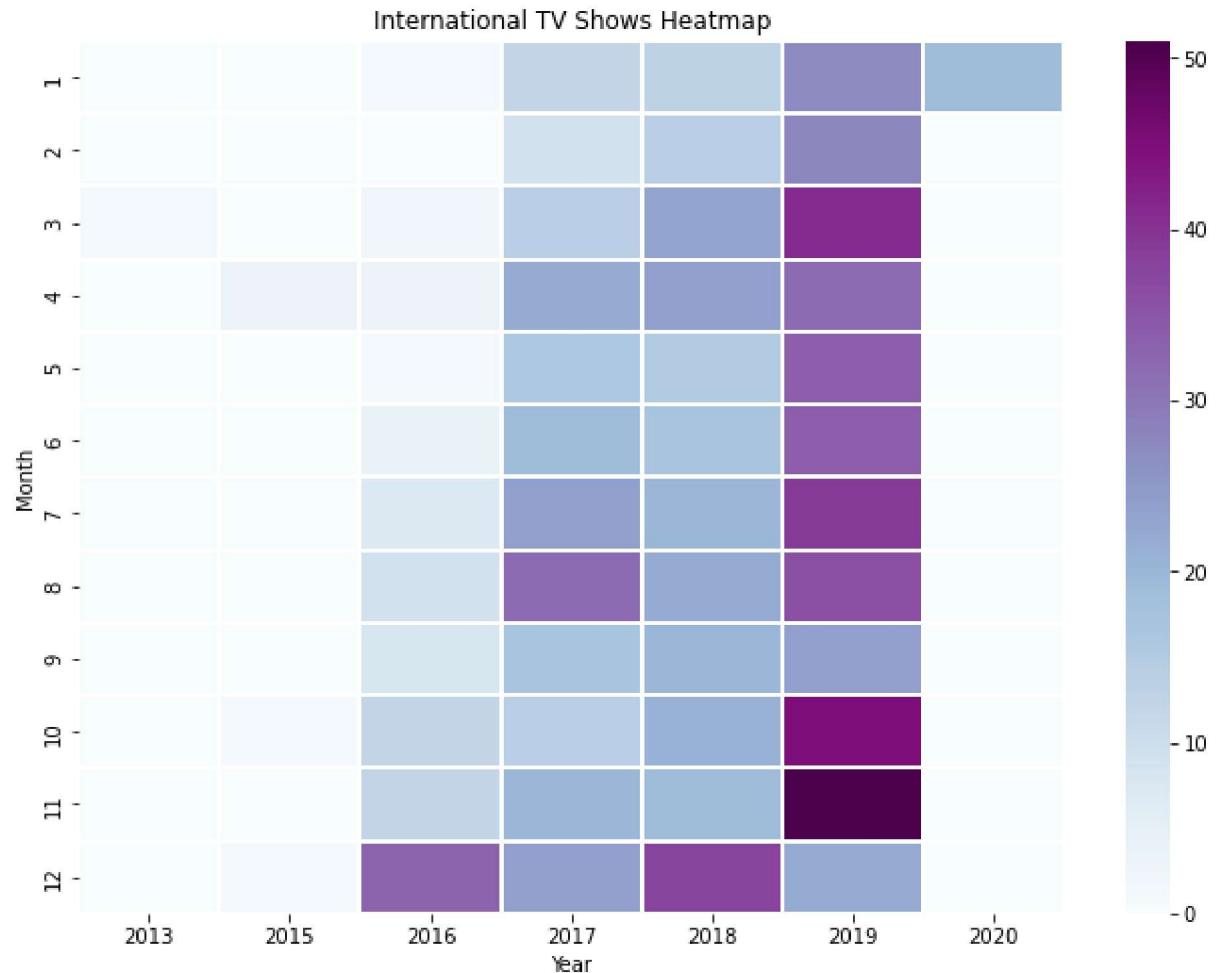
fig, ax = plt.subplots(figsize=(11, 8))
sns.heatmap(intmov_month_year_df, linewidths=0.025, cmap="BuPu")
plt.title("International Movies Heatmap")
plt.ylabel("Month")
plt.xlabel("Year")
plt.show()
```



In the heatmap above, we can see that a majority of international movies were added throughout the year in 2018. Then in December 2019, Netflix added the most international movie content.

```
In [31]: intsho_month_year_df = international_shows.groupby('year_added')['month_added'].value_counts().unstack().fillna(0).reset_index()

fig, ax = plt.subplots(figsize=(11, 8))
sns.heatmap(intsho_month_year_df, linewidths=0.025, cmap="BuPu")
plt.title("International TV Shows Heatmap")
plt.ylabel("Month")
plt.xlabel("Year")
plt.show()
```



In the above heatmap, we can see that the majority of international TV shows were added throughout the year 2019.

## 2. Conclude the most important features

On seeing the dataset, one can say that most important features present in the dataset are:

1. type
2. title
3. director
4. cast
5. country
6. rating
7. listed\_in
8. description

## 3. Identify similar content by matching text-based features

Here, we are going to create the NLP pipeline and then we will convert it into TD-IDF vectors, reduce it using singular value decomposition and then we will find the similar content using cosine similarity.

```
In [32]: netflix_titles_df['complete_description'] = netflix_titles_df['listed_in'] + " "
```

```
In [33]: import nltk
from nltk.corpus import stopwords
from nltk.stem import WordNetLemmatizer
import re
import string
```

```
In [34]: def wordopt(text):
    text = text.lower()
    text = re.sub('\[.*?\]', '', text)
    text = re.sub("\\"W", " ", text)
    text = re.sub('https?://\S+|www\.\S+', '', text)
    text = re.sub('<.*?>+', '', text)
    text = re.sub('[%s]' % re.escape(string.punctuation), ' ', text)
    text = re.sub('\n', ' ', text)
    text = re.sub('\w*\d\w*', ' ', text)
    return text
```

```
In [35]: def nlpSteps(text):
    tokens = nltk.word_tokenize(text)
    stop_words = stopwords.words('english')
    txt_no_stopwords = [word for word in tokens if word not in stop_words]
    lemmatizer = WordNetLemmatizer()
    lemmatized_words = []
    for w in txt_no_stopwords:
        lemmatized_words.append(lemmatizer.lemmatize(w))
    text = " ".join(lemmatized_words)
    return text
```

```
In [36]: netflix_titles_df['complete description'] = netflix_titles_df['complete descripti
```

```
In [37]: from sklearn.feature_extraction.text import TfidfVectorizer
vectors = TfidfVectorizer()
x = vectors.fit_transform(netflix_titles_df['complete description'])
```

```
In [38]: x.shape
```

```
Out[38]: (6214, 14865)
```

```
In [57]: from sklearn.decomposition import TruncatedSVD
svd = TruncatedSVD(n_components=2000)
matrix = svd.fit_transform(x)
```

```
In [58]: matrix.shape
```

```
Out[58]: (6214, 2000)
```

```
In [75]: from sklearn.metrics.pairwise import cosine_similarity
value = netflix_titles_df.loc[netflix_titles_df['title'] == 'Apaches'].index[0]
a = np.array(matrix[value]).reshape(1, -1)
score = cosine_similarity(matrix, a).reshape(-1)
```

```
In [80]: similar = pd.DataFrame(score, index = netflix_titles_df.title, columns=['Score'])
similar.sort_values('Score', ascending=False, inplace=True)
similar.head(11)
```

Out[80]:

	Score
	title
<b>Apaches</b>	1.000000
<b>El Cartel</b>	0.497653
<b>Tijuana</b>	0.478061
<b>Historia de un clan</b>	0.443355
<b>Las muñecas de la mafia</b>	0.438642
<b>Jack Taylor</b>	0.414341
<b>The Good Bandit</b>	0.412130
<b>Green Frontier</b>	0.407304
<b>Enemigo íntimo</b>	0.395472
<b>Four Seasons in Havana</b>	0.395314
<b>Miss Dynamite</b>	0.389976

You can see that we have find the similar movies based on their genre and their synosis.

## 4. Conclude whether Netflix is increasingly focusing on TV shows rather than movies

From the above EDA, we can say that Netflix is more focussing on Movies rather than TV shows.