**Netflix Business Case Study**

▾ Importing Libraries , Loading dataset

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

from google.colab import drive

drive.mount("/content/drive")
```

```
Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).
```

```
df = pd.read_csv('/content/drive/MyDrive/Colab Notebooks/netflix.csv')
```

```
df
```

|   | show_id | type | title | director | cast | country | date_added | release_ye |
|---|---------|------|-------|----------|------|---------|------------|------------|
| 0 | s1 | Movie | Dick Johnson Is Dead | Kirsten Johnson | NaN | United States | September 25, 2021 | 2C |
| 1 | s2 | TV Show | Blood & Water | NaN | Ama Qamata, Khosi Ngema, Gail Mabalane, Thaban... | South Africa | September 24, 2021 | 2C |
| 2 | s3 | TV Show | Ganglands | Julien Leclercq | Sami Bouajila, Tracy Gotoas, Samuel Jouy, Nabi... | NaN | September 24, 2021 | 2C |
| 3 | s4 | TV Show | Jailbirds New Orleans | NaN | NaN | NaN | September 24, 2021 | 2C |
| 4 | s5 | TV Show | Kota Factory | NaN | Mayur More, Jitendra Kumar, Ranjan Raj, Alam | India | September 24, 2021 | 2C |

```
df.shape
```

```
(8807, 12)
```

**The Dataset has 8807 Rows & 12 Columns**

```
df.nunique()
```

```
show_id        8807
type              2
title          8807
director       4528
cast           7692
country         748
date_added     1767
release_year     74
rating           17
duration        220
listed_in       514
description    8775
dtype: int64
```

**cheking the Unique Value in the Dataset**

Overall null values in each column of the dataset -

```
df.isna().sum()
```

```
show_id          0
type             0
title            0
director      2634
cast           825
country        831
date_added      10
release_year     0
rating           4
duration         3
listed_in        0
description      0
dtype: int64
```

**Form The Above analyse the Data set has many Missing Values in Director, Cast, Country**

```
df[df['duration'].isna()]
```

| | show_id | type | title | director | cast | country | date_added | release_year | rat |
|---|---|---|---|---|---|---|---|---|---|
| **5541** | s5542 | Movie | Louis C.K. 2017 | Louis C.K. | Louis C.K. | United States | April 4, 2017 | 2017 | 74 |

```
indices = df[df['duration'].isna()].index
```

```
df.loc[indices] = df.loc[indices].fillna(method = 'ffill' , axis = 1)
```

```
df.loc[indices ,'rating'] = 'Unknown'
```

```
df.loc[indices]
```

| | show_id | type | title | director | cast | country | date_added | release_year | ra |
|---|---|---|---|---|---|---|---|---|---|
| **5541** | s5542 | Movie | Louis C.K. 2017 | Louis C.K. | Louis C.K. | United States | April 4, 2017 | 2017 | Unk |

**Filling the null values**

```
df.fillna("Unknown", inplace = True )
```

```
df
```

| | show_id | type | title | director | cast | country | date_added | release_y |
|---|---|---|---|---|---|---|---|---|
| **0** | s1 | Movie | Dick Johnson Is Dead | Kirsten Johnson | Unknown | United States | September 25, 2021 | 2 |

Ama

```
df.groupby(["rating"])["title"].count()
```

```
rating
G             41
NC-17          3
NR            80
PG           287
PG-13        490
R            799
TV-14       2160
TV-G         220
TV-MA       3207
TV-PG        863
TV-Y         307
TV-Y7        334
TV-Y7-FV       6
UR             3
Unknown        7
Name: title, dtype: int64
```

Mayur

**The rating column has NR which is same as UR.**

Show      Factory                      24, 2021

Region

```
df.loc[df['rating'] == 'UR' , 'rating'] = 'NR'
df.rating.value_counts()
```

```
TV-MA       3207
TV-14       2160
TV-PG        863
R            799
PG-13        490
TV-Y7        334
TV-Y         307
PG           287
TV-G         220
NR            83
G             41
Unknown        7
TV-Y7-FV       6
NC-17          3
Name: rating, dtype: int64
```

- dropped the null from date_added column

```
df.drop(df.loc[df['date_added'].isna()].index , axis = 0 , inplace = True)
```

```
df['date_added'].isna().sum()
```

```
0
```

```
df['date_added']
```

```
0        September 25, 2021
1        September 24, 2021
2        September 24, 2021
3        September 24, 2021
4        September 24, 2021
                ...
8802     November 20, 2019
8803          July 1, 2019
8804      November 1, 2019
8805      January 11, 2020
8806         March 2, 2019
Name: date_added, Length: 8807, dtype: object
```

We can add the new column 'year_added' by extracting the year from 'date_added' column

**Non Graphical Analysis**

```
# 2 types of content present in dataset - either Movie or TV Show
df['type'].unique()

     array(['Movie', 'TV Show'], dtype=object)


movies  = df.loc[df['type'] == 'Movie']
tv_shows = df.loc[df['type'] == 'TV Show']


movies.duration.value_counts()

     90 min     152
     94 min     146
     97 min     146
     93 min     146
     91 min     144
                ...
     208 min      1
     5 min        1
     16 min       1
     186 min      1
     191 min      1
     Name: duration, Length: 205, dtype: int64


tv_shows.duration.value_counts()

     1 Season     1793
     2 Seasons     425
     3 Seasons     199
     4 Seasons      95
     5 Seasons      65
     6 Seasons      33
     7 Seasons      23
     8 Seasons      17
     9 Seasons       9
     10 Seasons      7
     13 Seasons      3
     15 Seasons      2
     12 Seasons      2
     11 Seasons      2
     17 Seasons      1
     Name: duration, dtype: int64
```

Since movie and TV shows both have different format for duration, we can change duration for movies as minutes & TV shows as seasons

```
movies['duration'] = movies['duration'].str[:-3]
movies['duration'] = movies['duration'].astype('float')


tv_shows['duration'] = tv_shows.duration.str[:-7].apply(lambda x : x.strip())
tv_shows['duration'] = tv_shows['duration'].astype('float')


tv_shows.rename({'duration': 'duration_in_seasons'} ,axis = 1 , inplace = True)
movies.rename({'duration': 'duration_in_minutes'} ,axis = 1 , inplace = True)


tv_shows.duration_in_seasons

     1       2.0
     2       1.0
     3       1.0
     4       2.0
     5       1.0
             ...
     8795    2.0
     8796    2.0
     8797    3.0
     8800    1.0
     8803    2.0
     Name: duration_in_seasons, Length: 2666, dtype: float64


movies.duration_in_minutes

     0        90.0
     6        91.0
     7       125.0
     9       104.0
     12      127.0
             ...
     8801     96.0
     8802    158.0
     8804     88.0
```

```
8805      88.0
8806     111.0
Name: duration_in_minutes, Length: 6131, dtype: float64
```

when was first movie added on netflix and when is the most recent movie added on netflix as per data i.e. dataset duration

```
timeperiod = pd.Series((df['date_added'].min().strftime('%B %Y') , df['date_added'].max().strftime('%B %Y')))
timeperiod.index = ['first' , 'Most Recent']
timeperiod
```

```
first          January 2008
Most Recent    September 2021
dtype: object
```

The oldest and the most recent movie/TV show released on the Netflix in which year?

```
df.release_year.min() , df.release_year.max()
```

```
(1925, 2021)
```

```
df.loc[(df.release_year == df.release_year.min()) | (df.release_year == df.release_year.max())].sort_values('release_year')
```

| | show_id | type | title | director | cast | country | date_added | release_v |
|---|---|---|---|---|---|---|---|---|
| **4250** | s4251 | TV Show | Pioneers: First Women Filmmakers* | NaN | NaN | NaN | 2018-12-30 | |
| **966** | s967 | Movie | Get the Grift | Pedro Antonio | Marcus Majella, Samantha Schmütz, Caito Mainie... | Brazil | 2021-04-28 | |
| **967** | s968 | TV Show | Headspace Guide to Sleep | NaN | Evelyn Lewis Prieto | NaN | 2021-04-28 | |
| **968** | s969 | TV Show | Sexify | NaN | Aleksandra Skraba, Maria Sobocińska, Sandra Dr... | Poland | 2021-04-28 | |
| **972** | s973 | TV Show | Fatma | NaN | Burcu Biricik, Uğur Yücel, Mehmet Yılmaz Ak, H... | Turkey | 2021-04-27 | |

Which are different ratings available on Netflix in each type of content? Check the number of content released in each type.

```
df.groupby(['type' , 'rating'])['show_id'].count()
```

```
type      rating
Movie     G                41
          NC-17             3
          NR               78
          Not Available     5
          PG              287
          PG-13           490
          R               797
          TV-14          1427
          TV-G            126
          TV-MA          2062
          TV-PG           540
          TV-Y            131
```

```
        TV-Y7                139
        TV-Y7-FV               5
TV Show  NR                    4
        Not Available          2
        R                      2
        TV-14                730
        TV-G                  94
        TV-MA               1143
        TV-PG                321
        TV-Y                 175
        TV-Y7                194
        TV-Y7-FV               1
Name: show_id, dtype: int64
```

Working on the columns having maximum null values and the columns having comma separated multiple values for each record

- Country column

```
df['country'].value_counts()
```

```
United States                              2812
India                                       972
United Kingdom                              418
Japan                                       244
South Korea                                 199
                                            ...
Romania, Bulgaria, Hungary                    1
Uruguay, Guatemala                            1
France, Senegal, Belgium                      1
Mexico, United States, Spain, Colombia        1
United Arab Emirates, Jordan                  1
Name: country, Length: 748, dtype: int64
```

We see that many movies are produced in more than 1 country. Hence, the country column has comma separated values of countries.

This makes it difficult to analyse how many movies were produced in each country. We can use explode function in pandas to split the country column into different rows.

we are Creating a separate table for country , to avoid the duplicasy of records in our origional table after exploding.

```
country_tb = df[['show_id' , 'type' , 'country']]
country_tb.dropna(inplace = True)
country_tb['country'] = country_tb['country'].apply(lambda x : x.split(','))
country_tb = country_tb.explode('country')
country_tb
```

|      | show_id | type    | country       |
|------|---------|---------|---------------|
| 0    | s1      | Movie   | United States |
| 1    | s2      | TV Show | South Africa  |
| 4    | s5      | TV Show | India         |
| 7    | s8      | Movie   | United States |
| 7    | s8      | Movie   | Ghana         |
| ...  | ...     | ...     | ...           |
| 8801 | s8802   | Movie   | Jordan        |
| 8802 | s8803   | Movie   | United States |
| 8804 | s8805   | Movie   | United States |
| 8805 | s8806   | Movie   | United States |
| 8806 | s8807   | Movie   | India         |

10010 rows × 3 columns

```
# some duplicate values are found, which have unnecessary spaces. some empty strings found
country_tb['country'] = country_tb['country'].str.strip()
```

```
country_tb.loc[country_tb['country'] == '']
```

| | show_id | type | country |
|---|---|---|---|
| 193 | s194 | TV Show | |
| 365 | s366 | Movie | |
| 1192 | s1193 | Movie | |
| 2224 | s2225 | Movie | |
| 4653 | s4654 | Movie | |

```
country_tb = country_tb.loc[country_tb['country'] != '']
```

```
country_tb['country'].nunique()
```

```
    122
```

Netflix has movies from the total 122 countries.

Total movies and tv shows in each country

```
x = country_tb.groupby(['country' , 'type'])['show_id'].count().reset_index()
x.pivot(index = ['country'] , columns = 'type' , values = 'show_id').sort_values('Movie',ascending = False)
```

| type | Movie | TV Show |
|---|---|---|
| **country** | | |
| **United States** | 2752.0 | 932.0 |
| **India** | 962.0 | 84.0 |
| **United Kingdom** | 534.0 | 271.0 |
| **Canada** | 319.0 | 126.0 |
| **France** | 303.0 | 90.0 |
| ... | ... | ... |
| **Azerbaijan** | NaN | 1.0 |
| **Belarus** | NaN | 1.0 |
| **Cuba** | NaN | 1.0 |
| **Cyprus** | NaN | 1.0 |
| **Puerto Rico** | NaN | 1.0 |

122 rows × 2 columns

- Director column

```
df['director'].value_counts()
```

```
    Rajiv Chilaka                      19
    Raúl Campos, Jan Suter             18
    Marcus Raboy                       16
    Suhas Kadav                        16
    Jay Karas                          14
                                       ..
    Raymie Muzquiz, Stu Livingston      1
    Joe Menendez                        1
    Eric Bross                          1
    Will Eisenberg                      1
    Mozez Singh                         1
    Name: director, Length: 4528, dtype: int64
```

There are some movies which are directed by multiple directors. Hence multiple names of directors are given in comma separated format. We will explode the director column as well. It will create many duplicate records in originaltable hence we created separate table for directors.

```
dir_tb = df[['show_id' , 'type' , 'director']]
dir_tb.dropna(inplace = True)
dir_tb['director'] = dir_tb['director'].apply(lambda x : x.split(','))
dir_tb
```

| | show_id | type | director |
|---|---|---|---|
| **0** | s1 | Movie | [Kirsten Johnson] |
| **2** | s3 | TV Show | [Julien Leclercq] |
| **5** | s6 | TV Show | [Mike Flanagan] |
| **6** | s7 | Movie | [Robert Cullen, José Luis Ucha] |
| **7** | s8 | Movie | [Haile Gerima] |
| **...** | ... | ... | ... |
| **8801** | s8802 | Movie | [Majid Al Ansari] |
| **8802** | s8803 | Movie | [David Fincher] |
| **8804** | s8805 | Movie | [Ruben Fleischer] |
| **8805** | s8806 | Movie | [Peter Hewitt] |
| **8806** | s8807 | Movie | [Mozez Singh] |

6173 rows × 3 columns

```python
dir_tb = dir_tb.explode('director')
```

```python
dir_tb['director'] = dir_tb['director'].str.strip()
```

```python
# checking if empty stirngs are there in director column
dir_tb.director.apply(lambda x : True if len(x) == 0 else False).value_counts()
```

```
    False    6978
    Name: director, dtype: int64
```

```python
dir_tb
```

| | show_id | type | director |
|---|---|---|---|
| **0** | s1 | Movie | Kirsten Johnson |
| **2** | s3 | TV Show | Julien Leclercq |
| **5** | s6 | TV Show | Mike Flanagan |
| **6** | s7 | Movie | Robert Cullen |
| **6** | s7 | Movie | José Luis Ucha |
| **...** | ... | ... | ... |
| **8801** | s8802 | Movie | Majid Al Ansari |
| **8802** | s8803 | Movie | David Fincher |
| **8804** | s8805 | Movie | Ruben Fleischer |
| **8805** | s8806 | Movie | Peter Hewitt |
| **8806** | s8807 | Movie | Mozez Singh |

6978 rows × 3 columns

```python
dir_tb['director'].nunique()
```

```
    4993
```

There are total 4993 unique directors in the dataset.

Total movies and tv shows directed by each director

```python
x = dir_tb.groupby(['director' , 'type'])['show_id'].count().reset_index()
x.pivot(index= ['director'] , columns = 'type' , values = 'show_id').sort_values('Movie' ,ascending = False)
```
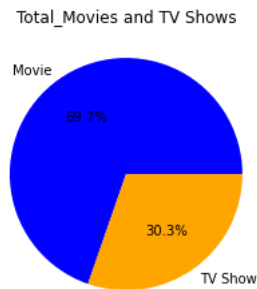
| type | Movie | TV Show |
|---|---|---|
| **director** | | |
| **Rajiv Chilaka** | 22.0 | NaN |
| **Jan Suter** | 21.0 | NaN |
| **Raúl Campos** | 19.0 | NaN |
| **Suhas Kadav** | 16.0 | NaN |
| **Marcus Raboy** | 15.0 | 1.0 |
| ... | ... | ... |
| **Vijay S. Bhanushali** | NaN | 1.0 |
| **YC Tom Lee** | NaN | 1.0 |

- 'listed_in' column to understand more about genres

```
genre_tb = df[['show_id' , 'type', 'listed_in']]


genre_tb['listed_in'] = genre_tb['listed_in'].apply(lambda x : x.split(','))
genre_tb = genre_tb.explode('listed_in')
genre_tb['listed_in'] = genre_tb['listed_in'].str.strip()


genre_tb
```

| | show_id | type | listed_in |
|---|---|---|---|
| **0** | s1 | Movie | Documentaries |
| **1** | s2 | TV Show | International TV Shows |
| **1** | s2 | TV Show | TV Dramas |
| **1** | s2 | TV Show | TV Mysteries |
| **2** | s3 | TV Show | Crime TV Shows |
| **...** | ... | ... | ... |
| **8805** | s8806 | Movie | Children & Family Movies |
| **8805** | s8806 | Movie | Comedies |
| **8806** | s8807 | Movie | Dramas |
| **8806** | s8807 | Movie | International Movies |
| **8806** | s8807 | Movie | Music & Musicals |

19303 rows × 3 columns

```
genre_tb.listed_in.unique()
```

```
array(['Documentaries', 'International TV Shows', 'TV Dramas',
       'TV Mysteries', 'Crime TV Shows', 'TV Action & Adventure',
       'Docuseries', 'Reality TV', 'Romantic TV Shows', 'TV Comedies',
       'TV Horror', 'Children & Family Movies', 'Dramas',
       'Independent Movies', 'International Movies', 'British TV Shows',
       'Comedies', 'Spanish-Language TV Shows', 'Thrillers',
       'Romantic Movies', 'Music & Musicals', 'Horror Movies',
       'Sci-Fi & Fantasy', 'TV Thrillers', "Kids' TV",
       'Action & Adventure', 'TV Sci-Fi & Fantasy', 'Classic Movies',
       'Anime Features', 'Sports Movies', 'Anime Series',
       'Korean TV Shows', 'Science & Nature TV', 'Teen TV Shows',
       'Cult Movies', 'TV Shows', 'Faith & Spirituality', 'LGBTQ Movies',
       'Stand-Up Comedy', 'Movies', 'Stand-Up Comedy & Talk Shows',
       'Classic & Cult TV'], dtype=object)
```

```
genre_tb.listed_in.nunique()
```

```
42
```

Total 42 genres present in dataset

```
df.merge(genre_tb , on = 'show_id' ).groupby(['type_y'])['listed_in_y'].nunique()
```

```
type_y
Movie      20
```

```
TV Show   22
Name: listed_in_y, dtype: int64
```

Movies have 20 genres and TV shows have 22 genres.

```python
# total movies/TV shows in each genre
x = genre_tb.groupby(['listed_in' , 'type'])['show_id'].count().reset_index()
x.pivot(index = 'listed_in' , columns = 'type' , values = 'show_id').sort_index()
```

| | Documentaries | 869.0 | NaN |
|---|---|---|---|
| | Dramas | 1.2743 | NaN |

```
# Exploring cast column
```

| | Horror Movies | 357.0 | NaN |

```
cast_tb = df[['show_id' , 'type' ,'cast']]
cast_tb.dropna(inplace = True)
cast_tb['cast'] = cast_tb['cast'].apply(lambda x : x.split(','))
cast_tb = cast_tb.explode('cast')
cast_tb
```

| | show_id | type | cast |
|---|---|---|---|
| 1 | s2 | TV Show | Ama Qamata |
| 1 | s2 | TV Show | Khosi Ngema |
| 1 | s2 | TV Show | Gail Mabalane |
| 1 | s2 | TV Show | Thabang Molaba |
| 1 | s2 | TV Show | Dillon Windvogel |
| ... | ... | ... | ... |
| 8806 | s8807 | Movie | Manish Chaudhary |
| 8806 | s8807 | Movie | Meghna Malik |
| 8806 | s8807 | Movie | Malkeet Rauni |
| 8806 | s8807 | Movie | Anita Shabdish |
| 8806 | s8807 | Movie | Chittaranjan Tripathy |

64057 rows × 3 columns

```
cast_tb['cast'] = cast_tb['cast'].str.strip()
```

| | TV Comedies | NaN | 574.0 |

```
# checking empty strings
cast_tb[cast_tb['cast'] == '']
```

| | show_id | type | cast |
|---|---|---|---|
| | TV Sci-Fi & Fantasy | NaN | 83.0 |

```
# Total actors on the Netflix
cast_tb.cast.nunique()
```

    36403

```
# Total movies/TV shows by each actor
x = cast_tb.groupby(['cast' , 'type'])['show_id'].count().reset_index()
x.pivot(index = 'cast' , columns = 'type' , values = 'show_id').sort_values('TV Show' , ascending = False)
```

| type | Movie | TV Show |
|---|---|---|
| cast | | |
| Takahiro Sakurai | 7.0 | 25.0 |
| Yuki Kaji | 10.0 | 19.0 |
| Junichi Suwabe | 4.0 | 17.0 |
| Daisuke Ono | 5.0 | 17.0 |
| Ai Kayano | 2.0 | 17.0 |
| ... | ... | ... |
| Şerif Sezer | 1.0 | NaN |
| Şevket Çoruh | 1.0 | NaN |
| Şinasi Yurtsever | 3.0 | NaN |
| Şükran Ovalı | 1.0 | NaN |
| Şọpẹ́ Dìrísù | 1.0 | NaN |

36403 rows × 2 columns

## ▾ 4. Visual Analysis - Univariate & Bivariate

- 4.1. Distribution of content across the different types

```
types = df.type.value_counts()
plt.pie(types,  labels=types.index, autopct='%1.1f%%' , colors = ['blue' , 'orange'])
plt.title('Total_Movies and TV Shows')
plt.show()
```

Total_Movies and TV Shows



It is observed that , around 70% content is Movies and around 30% content is TV shows.

- 4.2 Distribution of 'date_added' column

How has the number of movies/TV shows added on Netflix per year changed over the time?

```
d = df.groupby(['year_added' ,'type' ])['show_id'].count().reset_index()
d.rename({'show_id' : 'total movies/TV shows'}, axis = 1 , inplace = True)
```

```
plt.figure(figsize = (12,6))
sns.lineplot(data = d , x = 'year_added' , y = 'total movies/TV shows' , hue = 'type', marker = 'o'  , ms = 6)
plt.xlabel('year_added' , fontsize = 12)
plt.ylabel('total movies/TV shows' , fontsize = 12)
plt.title('total movies and TV shows by the year_added' , fontsize = 12)
plt.show()
```



Observation:

- The content added on the Netflix surged drastically after 2015.
- 2019 marks the highest number of movies and TV shows added on the Netflix.
- Year 2020 and 2021 has seen the drop in content added on Netflix, possibly because of Pandemic. But still , TV shows content have not dropped as drastic as movies. In recent years TV shows are focussed more than Movies.

- 4.3 Distribution of 'Release_year' column

How has the number of movies released per year changed over the last 20-30 years?

```
d = df.groupby(['type' , 'release_year'])['show_id'].count().reset_index()
d.rename({'show_id' : 'total movies/TV shows'}, axis = 1 , inplace = True)
d
```

| | type | release_year | total movies/TV shows |
|---|---|---|---|
| 0 | Movie | 1942 | 2 |
| 1 | Movie | 1943 | 3 |
| 2 | Movie | 1944 | 3 |
| 3 | Movie | 1945 | 3 |
| 4 | Movie | 1946 | 1 |
| ... | ... | ... | ... |
| 114 | TV Show | 2017 | 265 |
| 115 | TV Show | 2018 | 379 |
| 116 | TV Show | 2019 | 397 |
| 117 | TV Show | 2020 | 436 |
| 118 | TV Show | 2021 | 315 |

119 rows × 3 columns

```
plt.figure(figsize = (12,6))
sns.lineplot(data = d , x = 'release_year' , y = 'total movies/TV shows' , hue = 'type' , marker = 'o'  , ms = 6 )
plt.xlabel('release_year' , fontsize = 12)
plt.ylabel('total movies/TV shows' , fontsize = 12)
plt.title('total movies and TV shows by the release_year' , fontsize = 12)
plt.xlim( left = 2000 , right = 2021)
plt.xticks(np.arange(2000 , 2021 , 2))
plt.show()
```



Observation:

- 2018 marks the highest number of movie and TV show releases.
- Since 2018, A drop in movies is seen and rise in TV shows is observed clearly, and TV shows surpasses the movies count in mid 2020.
- In recent years TV shows are focussed more than Movies.
- The yearly number of releases has surged drastically from 2015.

- 4.4 Total movies/TV shows by each director

```
# total Movies directed by top 10 directors
top_10_dir = dir_tb.director.value_counts().head(10).index
df_new = dir_tb.loc[dir_tb['director'].isin(top_10_dir)]
```

```
plt.figure(figsize= (8 , 6))
sns.countplot(data = df_new , y = 'director' , order = top_10_dir , orient = 'v')
plt.xlabel('total_movies/TV shows' , fontsize = 12)
plt.xlabel('Movies/TV shows count')
plt.ylabel('Directors' , fontsize = 12)
plt.title('Total_movies/TVshows_by_director')
plt.show()
```

Total_movies/TVshows_by_director

Observation:

- The top 3 directors on Netflix in terms of count of movies directed by them are - Rajiv Chilaka, Jan Suter, Raúl Campos

- 4.4 Checking Outliers for number of movies directed by each director

```
x = dir_tb.director.value_counts()
x

        Rajiv Chilaka    22
        Jan Suter        21
        Raúl Campos      19
        Suhas Kadav      16
        Marcus Raboy     16
                         ..
        Raymie Muzquiz    1
        Stu Livingston    1
        Joe Menendez      1
        Eric Bross        1
        Mozez Singh       1
        Name: director, Length: 4993, dtype: int64
```

```
def calculate_outliers(data):
    # Calculate the first quartile (Q1)
    q1 = np.percentile(data, 25)

    # Calculate the third quartile (Q3)
    q3 = np.percentile(data, 75)

    # Calculate the interquartile range (IQR)
    iqr = q3 - q1

    # Determine the lower and upper bounds for outliers
    lower_bound = q1 - 1.5 * iqr
    upper_bound = q3 + 1.5 * iqr

    # Identify outliers in the dataset
    outliers = [value for value in data if value < lower_bound or value > upper_bound]

    return outliers
```

```
def calculate_max_occurred_value(data):
    # Calculate the unique values and their counts in the dataset
    unique_values, value_counts = np.unique(data, return_counts=True)

    # Find the index of the maximum count
    max_count_index = np.argmax(value_counts)

    # Retrieve the corresponding unique value with the maximum count
    max_occurred_value = unique_values[max_count_index]

    return max_occurred_value
```

```
outliers = calculate_outliers(x)  # Implement your outlier calculation method
max_occurred_value = calculate_max_occurred_value(x)  # Implement your method to find the maximum-occurred value
set(outliers)
```

```
{2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 15, 16, 19, 21, 22}
```
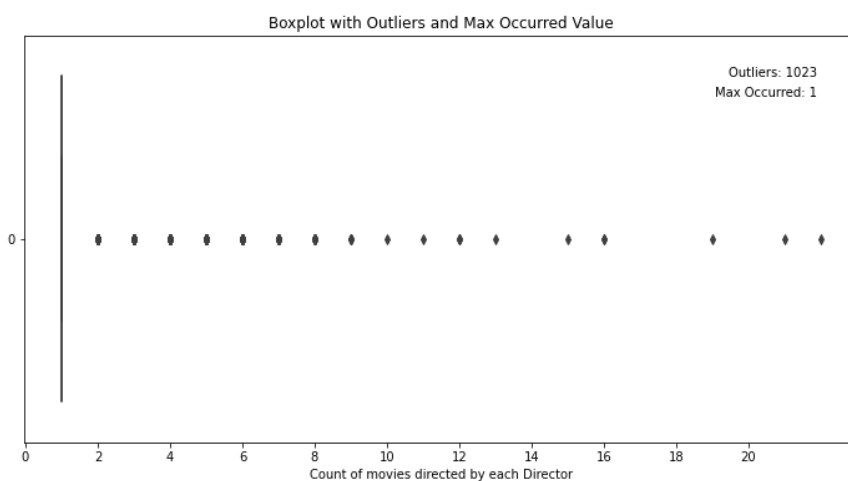
```
max_occurred_value
```

```
    1
```

```python
plt.figure(figsize = (12,6))
sns.boxplot(data=x, showfliers=True, whis=1.5 , orient = 'h')

# Calculate the outliers and maximum-occurred value
outliers = calculate_outliers(x)  # Implement your outlier calculation method
max_occurred_value = calculate_max_occurred_value(x)  # Implement your method to find the maximum-occurred value

# Annotate the plot
plt.text(0.95, 0.9, f"Outliers: {len(outliers)}", transform=plt.gca().transAxes, ha='right')
plt.text(0.95, 0.85, f"Max Occurred: {max_occurred_value}", transform=plt.gca().transAxes, ha='right')


plt.xlabel("Count of movies directed by each Director")
plt.xticks(np.arange(0,22,2))
plt.title("Boxplot with Outliers and Max Occurred Value")

# Show the plot
plt.show()
```



It is Observed that maximum occured value is 1, which means maximum directors on the Netflix have directed 1 movie/Tv show. There are few directors who have directed more than 1 movies/tv shows and they are outliers.

- 4.5 Total movies/TV shows by each country

```python
# Lets check for top 10 countries
top_10_country = country_tb.country.value_counts().head(10).index
df_new = country_tb.loc[country_tb['country'].isin(top_10_country)]


x = df_new.groupby(['country' , 'type'])['show_id'].count().reset_index()
x.pivot(index = 'country' , columns = 'type' , values = 'show_id').sort_values('Movie',ascending = False)
```

|  | type | Movie | TV Show |
|---|---|---|---|
| country |  |  |  |

```
plt.figure(figsize= (8,5))
sns.countplot(data = df_new , x = 'country' , order = top_10_country , hue = 'type')
plt.xticks(rotation = 90 , fontsize = 12)
plt.ylabel('total_movies/TV shows' , fontsize = 12)
plt.xlabel('')
plt.title('Total_movies/TVshows_by_country')
plt.show()
```



Total_movies/TVshows_by_country

```
top_10_country = country_tb.country.value_counts().head(10).index
country_tb['cat'] = country_tb['country'].apply(lambda x : x if x in top_10_country else 'Other Countries' )


x = country_tb.cat.value_counts()

plt.figure(figsize = (8,8))
plt.pie(x , labels = x.index, autopct='%1.1f%%')
plt.title('Total Content produced in each country' , fontsize = 15)
plt.show()
```



Total Content produced in each country

- Observation:
  - United States is the HIGHEST contributor country on Netflix, followed by India and United Kingdom.
  - Maximum content of Netflix which is around 75% , is coming from these top 10 countries. Rest of the world only contributes 25% of the content.

- 4.6 Total content distribution by release year of the content

```
plt.figure(figsize= (12,6))
sns.boxplot(data = df , x = 'release_year')
plt.xlabel('release_year' , fontsize = 12)
plt.title('Total_movies/TVshows_by_release_year')
plt.xticks(np.arange(1940 , 2021 , 5))
plt.xlim((1940 , 2022))
plt.show()
```



- Netflix have major content which is released in the year range 2000-2021
- It seems that the content older than year 2000 is almost missing from the Netflix.

- 4.7 Total movies/TV shows distribution by rating of the content

```
m = movies.loc[~movies.rating.isin(['Not Available' , 'NC-17' , 'TV-Y7-FV'])]
m = m.rating.value_counts()
t = tv_shows.loc[~tv_shows.rating.isin(['Not Available' , 'R' , 'NR', 'TV-Y7-FV'])]
t = t.rating.value_counts()
```

```
fig, ax = plt.subplots(1,2, figsize=(14,8))
ax[0].pie(m , labels = m.index, autopct='%1.1f%%')
ax[0].set_title('Total_movies_by_rating')

ax[1].pie(t , labels = t.index, autopct='%1.1f%%')
ax[1].set_title('Total_TV_shows_by_rating')

plt.tight_layout()
plt.show()
```

Total_movies_by_rating

Total_TV_shows_by_rating

TV-MA

Highest number of movies and TV shows are rated TV-MA (for mature audiences), followed by TV-14 & R/TV-PG
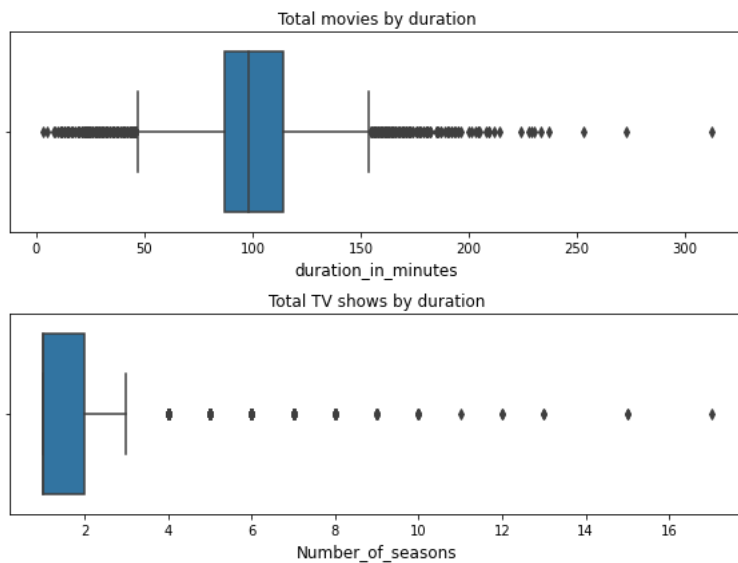
TV-14

- 4.8 Total movies/TV shows distribution by duration of the content

2.1%

TV-G

TV-G

```
fig, ax = plt.subplots(2,1, figsize=(8,6))

sns.boxplot (data = movies , x = 'duration_in_minutes' ,ax =ax[0])
ax[0].set_xlabel('duration_in_minutes' ,  fontsize = 12)
ax[0].set_title('Total movies by duration')

sns.boxplot (data = tv_shows , x = 'duration_in_seasons' , ax = ax[1])
ax[1].set_xlabel('Number_of_seasons' ,  fontsize = 12)
ax[1].set_title('Total TV shows by duration')

plt.tight_layout()
plt.show()
```

Total movies by duration

duration_in_minutes

Total TV shows by duration

Number_of_seasons

- Movie Duration: 50 mins - 150 mins is the range excluding potential outliers (values lying outside the whiskers of boxplot)
- TV Show Duration: 1-3 seasons is the range for TV shows excluding potential outliers
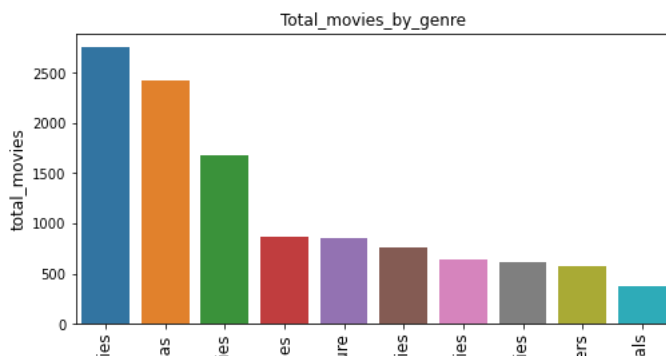
- 4.9 Total movies/TV shows in each Genre

```
# Lets check the count for top 10 genres in Movies and TV_shows


top_10_movie_genres = genre_tb[genre_tb['type'] == 'Movie'].listed_in.value_counts().head(10).index
df_movie = genre_tb.loc[genre_tb['listed_in'].isin(top_10_movie_genres)]


top_10_TV_genres = genre_tb[genre_tb['type'] == 'TV Show'].listed_in.value_counts().head(10).index
df_tv = genre_tb.loc[genre_tb['listed_in'].isin(top_10_TV_genres)]


plt.figure(figsize= (8,4))
sns.countplot(data = df_movie , x = 'listed_in' , order = top_10_movie_genres)
plt.xticks(rotation = 90 , fontsize = 12)
plt.ylabel('total_movies' , fontsize = 12)
plt.xlabel('Genres' , fontsize = 12)
plt.title('Total_movies_by_genre')
plt.show()
```
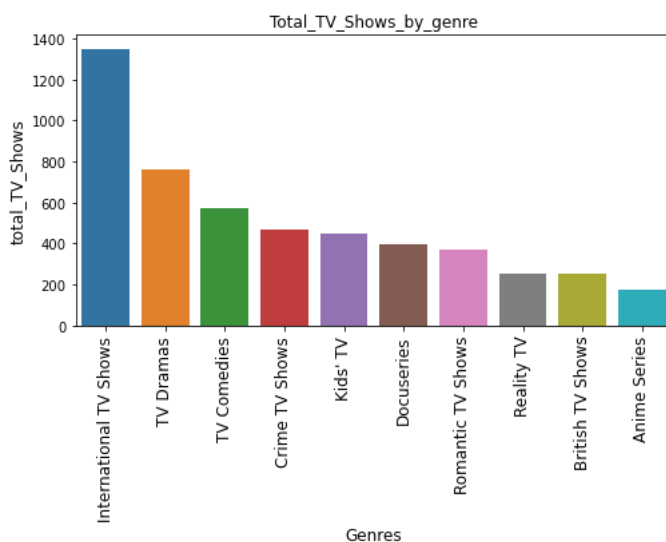
Total_movies_by_genre

```
plt.figure(figsize= (8,4))
sns.countplot(data = df_tv , x = 'listed_in' , order = top_10_TV_genres)
plt.xticks(rotation = 90 , fontsize = 12)
plt.ylabel('total_TV_Shows' , fontsize = 12)
plt.xlabel('Genres' , fontsize = 12)
plt.title('Total_TV_Shows_by_genre')
plt.show()
```



Total_TV_Shows_by_genre

- International Movies and TV Shows , Dramas , and Comedies are the top 3 genres on Netflix for both Movies and TV shows.
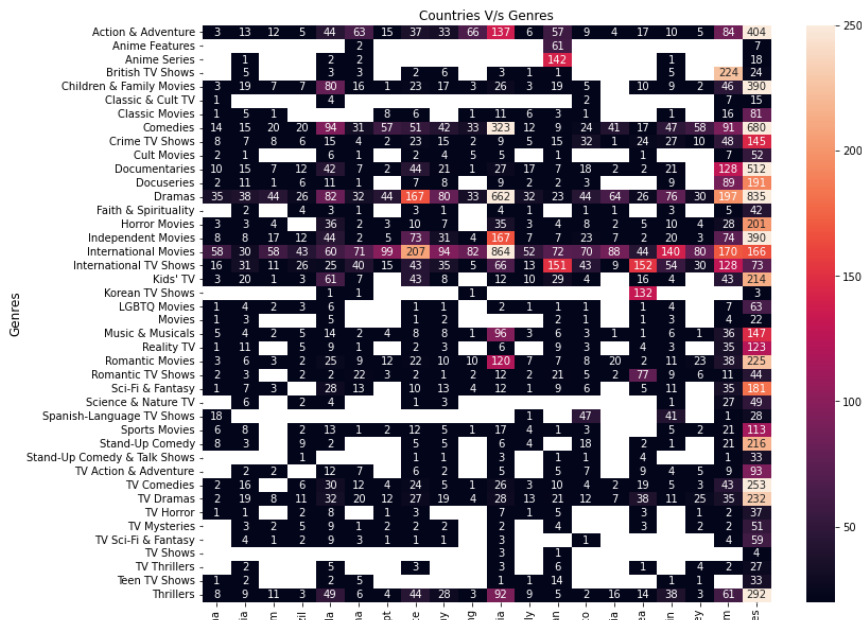
# 5. Bivariate Analysis

- 5.1 Lets check popular genres in top 20 countries

```
top_20_country = country_tb.country.value_counts().head(20).index
top_20_country = country_tb.loc[country_tb['country'].isin(top_20_country)]


x = top_20_country.merge(genre_tb , on = 'show_id').drop_duplicates()
country_genre = x.groupby([ 'country' , 'listed_in'])['show_id'].count().sort_values(ascending = False).reset_index()
country_genre = country_genre.pivot(index = 'listed_in' , columns = 'country' , values = 'show_id')


plt.figure(figsize = (12,10))
sns.heatmap(data = country_genre , annot = True , fmt=".0f" , vmin = 20 , vmax = 250 )
plt.xlabel('Countries' , fontsize = 12)
plt.ylabel('Genres' , fontsize = 12)
plt.title('Countries V/s Genres' , fontsize = 12)
```

```
Text(0.5, 1.0, 'Countries V/s Genres')
```



Countries V/s Genres

Popular genres across countries: Action & Adventure, Children & Family Movies, Comedies, Dramas, International Movies & TV Shows, TV Dramas, Thrillers

Country-specific genres: Korean TV shows (Korea), British TV Shows (UK), Anime features and Anime series (Japan), Spanish TV Shows (Argentina, Mexico and Spain)

United States and UK have a good mix of almost all genres.

Maximum International movies are produced in India.

### 5.2 Country-wise Rating of Content

```
x = top_20_country.merge(df , on = 'show_id').groupby(['country_x' , 'rating'])['show_id'].count().reset_index()


country_rating = x.pivot(index = ['country_x'] , columns = 'rating' , values = 'show_id')


plt.figure(figsize = (10,8))
sns.heatmap(data = country_rating , annot = True , fmt=".0f"  , vmin = 10 , vmax=200)
plt.ylabel('Countries' , fontsize = 12)
plt.xlabel('Rating' , fontsize = 12)
plt.title('Countries V/s Rating' , fontsize = 12)
```

```
Text(0.5, 1.0, 'Countries V/s Rating')
```

Countries V/s Rating

- Overall, Netflix has an large amount of adult content across all countries (TV-MA & TV-14).
- India also has many titles rated TV-PG, other than TV-MA & TV-14.
- Only US, Canada, UK, France and Japan have content for young audiences (TV-Y & TV-Y7).
- There is scarce content for general audience (TV-G & G) across all countries except US.

- 5.3 The top actors by country

Germany

```
x = cast_tb.merge(country_tb , on = 'show_id').drop_duplicates()
x = x.groupby(['country' , 'cast'])['show_id'].count().reset_index()
x.loc[x['country'].isin(['United States'])].sort_values('show_id' , ascending = False).head(5)
```

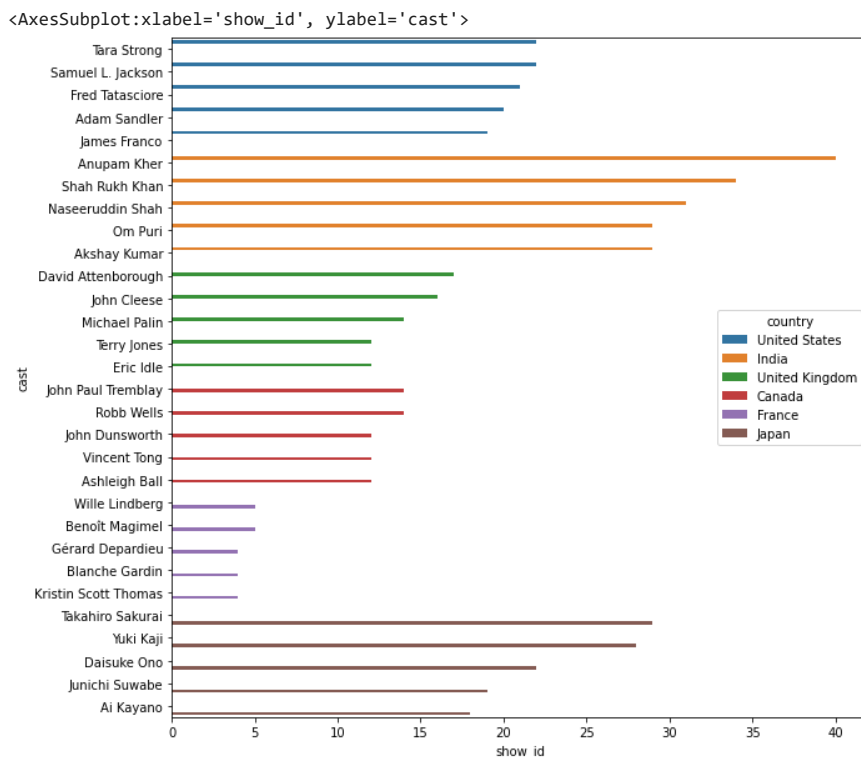|       | country       | cast              | show_id |
|-------|---------------|-------------------|---------|
| 49405 | United States | Tara Strong       | 22      |
| 48330 | United States | Samuel L. Jackson | 22      |
| 40463 | United States | Fred Tatasciore   | 21      |
| 35733 | United States | Adam Sandler      | 20      |
| 41672 | United States | James Franco      | 19      |

```
country_list = ['India' , 'United Kingdom' , 'Canada' , 'France' , 'Japan']
top_5_actors = x.loc[x['country'].isin(['United States'])].sort_values('show_id' , ascending = False).head(5)


for i in country_list:
    new = x.loc[x['country'].isin([i])].sort_values('show_id' , ascending = False).head(5)
    top_5_actors = pd.concat( [top_5_actors , new] , ignore_index = True)



# top 5 actors in top countries and their movies/tv shows count
top_5_actors
```

|   | country | cast | show_id |
|---|---|---|---|
| 0 | United States | Tara Strong | 22 |
| 1 | United States | Samuel L. Jackson | 22 |
| 2 | United States | Fred Tatasciore | 21 |
| 3 | United States | Adam Sandler | 20 |
| 4 | United States | James Franco | 19 |

```
plt.figure(figsize = (10,10))
sns.barplot(data = top_5_actors , y = 'cast' , x = 'show_id' , hue = 'country')
```

```
<AxesSubplot:xlabel='show_id', ylabel='cast'>
```



- 5.4 Top 5 directors by Genre

```
genre_list = [ 'Children & Family Movies', 'Comedies','Dramas', 'International Movies', 'Documentaries' ,
              'International TV Shows', 'Sci-Fi & Fantasy', 'Thrillers', 'Horror Movies']

x = dir_tb.merge(genre_tb , on = 'show_id').groupby([ 'listed_in' , 'director',])['show_id'].count().reset_index()

top_5_dir = x.loc[x['listed_in'] == 'Action & Adventure'].sort_values('show_id' , ascending = False).head()

for i in genre_list:
    new = x.loc[x['listed_in'] == i].sort_values('show_id' , ascending = False).head()
    top_5_dir = pd.concat([top_5_dir , new])

top_5_dir
```

| 4590 | Dramas | Hanung Bramantyo | 8 |
|---|---|---|---|
| 5544 | Dramas | S.S. Rajamouli | 7 |
| 7509 | International Movies | Cathy Garcia-Molina | 13 |
| 9330 | International Movies | Youssef Chahine | 10 |
| 9340 | International Movies | Yılmaz Erdoğan | 9 |
| 7620 | International Movies | David Dhawan | 8 |
| 8208 | International Movies | Kunle Afolayan | 8 |
| 3834 | Documentaries | Vlad Yudin | 6 |
| 3799 | Documentaries | Thierry Donard | 5 |
| 3217 | Documentaries | Edward Cotterill | 4 |
| 3262 | Documentaries | Frank Capra | 4 |
| 3075 | Documentaries | Barry Avrich | 4 |
| 9373 | International TV Shows | Alastair Fothergill | 3 |
| 9419 | International TV Shows | Hsu Fu-chun | 2 |
| 9436 | International TV Shows | Jung-ah Im | 2 |
| 9501 | International TV Shows | Shin Won-ho | 2 |
| 9478 | International TV Shows | Pali Yahya | 1 |
| 10752 | Sci-Fi & Fantasy | Lilly Wachowski | 4 |
| 10744 | Sci-Fi & Fantasy | Lana Wachowski | 4 |
| 10684 | Sci-Fi & Fantasy | Guillermo del Toro | 3 |
| 10790 | Sci-Fi & Fantasy | Paul W.S. Anderson | 3 |
| 10635 | Sci-Fi & Fantasy | Barry Sonnenfeld | 3 |
| 11974 | Thrillers | Rathindran R Prasad | 4 |
| 11698 | Thrillers | David Fincher | 4 |
| 11612 | Thrillers | Anurag Kashyap | 3 |
| 11636 | Thrillers | Brad Anderson | 3 |
| 11754 | Thrillers | Gregory Hoblit | 3 |
| 6280 | Horror Movies | Rocky Soraya | 6 |
| 6260 | Horror Movies | Poj Arnon | 5 |
| 6267 | Horror Movies | Rathindran R Prasad | 4 |
| 6191 | Horror Movies | Leigh Janiak | 3 |
| 6052 | Horror Movies | Banjong Pisanthanakun | 3 |

- 5.5 Top 5 genres in each country

```
x = genre_tb.merge(country_tb , on = 'show_id').drop_duplicates()
x = x.groupby(['country' , 'listed_in'])['show_id'].count().reset_index()
x.loc[x['country'] == 'United States'].sort_values('show_id' , ascending = False).head(5)

country_list = ['India'  , 'United Kingdom' , 'Canada' , 'France' , 'Japan']
top_5_genre = x.loc[x['country'].isin(['United States'])].sort_values('show_id' , ascending = False).head(5)

for i in country_list:
    new = x.loc[x['country'] == i].sort_values('show_id' , ascending = False).head(5)
    top_5_genre = pd.concat( [top_5_genre , new] , ignore_index = True)
```
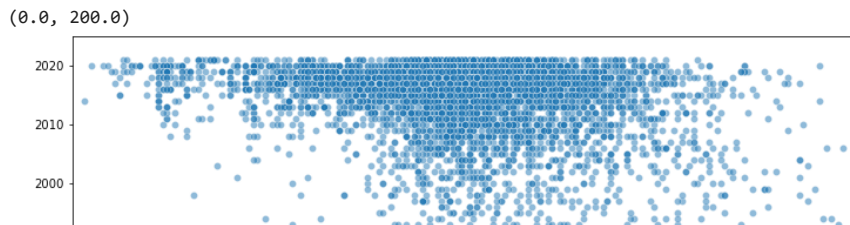
```
top_5_genre
```

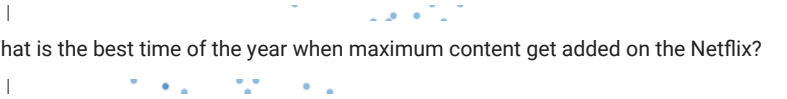|     | country        | listed_in               | show_id |
|-----|----------------|-------------------------|---------|
| 0   | United States  | Dramas                  | 835     |
| 1   | United States  | Comedies                | 680     |
| 2   | United States  | Documentaries           | 512     |
| 3   | United States  | Action & Adventure      | 404     |
| 4   | United States  | Independent Movies      | 390     |
| 5   | India          | International Movies     | 864     |
| 6   | India          | Dramas                  | 662     |
| 7   | India          | Comedies                | 323     |
| 8   | India          | Independent Movies      | 167     |
| 9   | India          | Action & Adventure      | 137     |
| 10  | United Kingdom | British TV Shows        | 224     |
| 11  | United Kingdom | Dramas                  | 197     |
| 12  | United Kingdom | International Movies     | 170     |
| 13  | United Kingdom | International TV Shows   | 128     |
| 14  | United Kingdom | Documentaries           | 128     |
| 15  | Canada         | Comedies                | 94      |
| 16  | Canada         | Dramas                  | 82      |
| 17  | Canada         | Children & Family Movies | 80     |
| 18  | Canada         | Kids' TV                | 61      |
| 19  | Canada         | International Movies     | 60      |
| 20  | France         | International Movies     | 207     |
| 21  | France         | Dramas                  | 167     |
| 22  | France         | Independent Movies      | 73      |
| 23  | France         | Comedies                | 51      |
| 24  | France         | Thrillers               | 44      |
| 25  | Japan          | International TV Shows   | 151     |
| 26  | Japan          | Anime Series            | 142     |
| 27  | Japan          | International Movies     | 72      |
| 28  | Japan          | Anime Features          | 61      |
| 29  | Japan          | Action & Adventure      | 57      |

- 5.6 Variation in duration of movies by Release year

```
plt.figure(figsize = (12,8))
sns.scatterplot(movies['duration_in_minutes'], movies['release_year'],  alpha=0.5)
plt.xlim((0,200))
```

(0.0, 200.0)



- Observation
  - The movies shorter than 150 minutes duration have increased drastically after 2000 while movies longer than 150 minutes are not much popular.
  - There is a huge surge in the number of shorter duration movies (less than 75 mins) post 2010. Overall, Short movies have been popular in last 10 years.
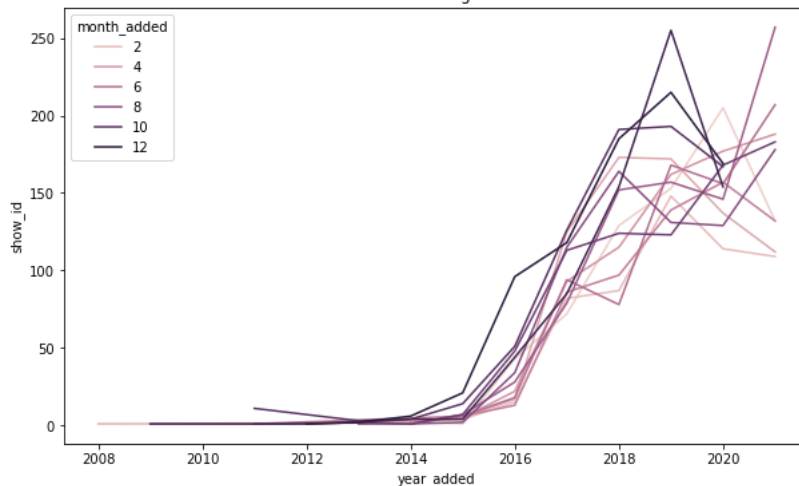
- 5.7 What is the best time of the year when maximum content get added on the Netflix?

```
month_year = df.groupby(['year_added' , 'month_added'])['show_id'].count().reset_index()
```
                                    duration_in_minutes

```
plt.figure(figsize = (10,6))
sns.lineplot(data=month_year, x = 'year_added', y = 'show_id', hue='month_added')
plt.title('Year and Month of Adding Shows on Netflix')
```

Text(0.5, 1.0, 'Year and Month of Adding Shows on Netflix')



- The number of shows getting added is increasing with each year until 2020.
- Also, months in the last quarter of the year (Oct-Dec) have more shows being added than the other months of the year. This could be because US has its festive season in Dec and India also has Diwali in Oct-Nov.

- 5.8 Which countries are adding more number of content over the time?

```
country_list = country_tb.country.value_counts().head(12).index
top_12_country = country_tb.loc[country_tb['country'].isin(country_list)]
country_year = top_12_country.merge(df , on = 'show_id')[['show_id','country_x' ,'type_x' , 'year_added' ]]
country_year.columns = ['show_id', 'country', 'type', 'year_added']
```

```
country_year = country_year.groupby(['country' , 'year_added'])['show_id'].count().reset_index()
```

```
plt.figure(figsize = (10,6))
sns.lineplot(data = country_year , x = 'year_added' , y = 'show_id' , hue = 'country' , palette ='rainbow' )
```

```
<AxesSubplot:xlabel='year_added', ylabel='show_id'>
```
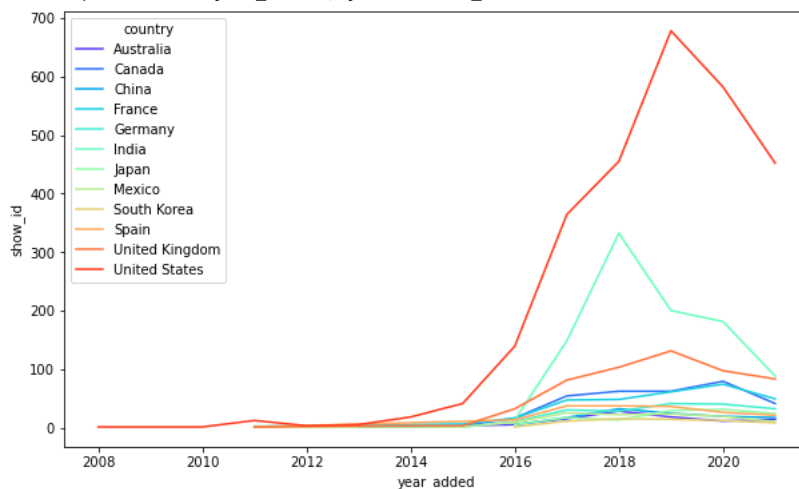


Observation : United Stated have always added highset number of movies/TV shows over the time. Since 2016, India has seen spike in popularity of content and added more number of content, followed by United Kingdom at 3rd position.



```
movie_type = country_year.loc[country_year.type == 'Movie'].groupby(['country' , 'year_added'])['show_id'].count().reset_index()
tv_type = country_year.loc[country_year.type == 'TV Show'].groupby(['country' , 'year_added'])['show_id'].count().reset_index()


plt.figure(figsize = (10,6))
sns.lineplot(data = movie_type , x = 'year_added' , y = 'show_id' , hue = 'country' , palette ='rainbow' )
```
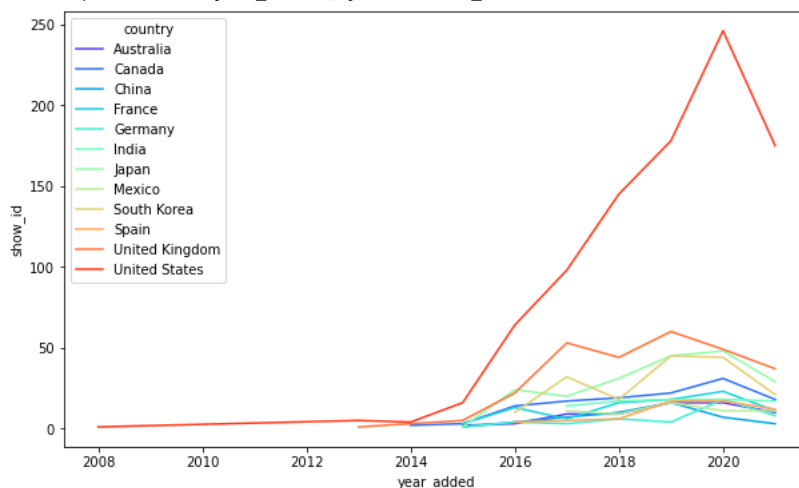
```
<AxesSubplot:xlabel='year_added', ylabel='show_id'>
```



```
plt.figure(figsize = (10,6))
sns.lineplot(data = tv_type , x = 'year_added' , y = 'show_id' , hue = 'country' , palette ='rainbow' )
```

```
<AxesSubplot:xlabel='year_added', ylabel='show_id'>
```



Observation: It is observed that United States tops in both movies and TV Shows. India is at 2nd positon in movies but In TV shows United Kingdom is at 2nd position, followed by India ,South Korea , Australia. It shows in countries like United Kingdom , South Korea , Australia TV Shows popularity is rising more than movies

## ▼ Insights based on Non-Graphical and Visual Analysis

- Around 70% content on Netflix is Movies and around 30% content is TV shows.
- The movies and TV shows uploading on the Netflix started from the year 2008, It had very lesser content till 2014.
- Year 2015 marks the drastic surge in the content getting uploaded on Netflix. It continues the uptrend since then and 2019 marks the highest number of movies and TV shows added on the Netflix. Year 2020 and 2021 has seen the drop in content added on Netflix, possibly because of Pandemic. But still , TV shows content have not dropped as drastic as movies.
- Since 2018, A drop in the movies is seen , but rise in TV shows is observed clearly. Being in continuous uptrend , TV shows surpassed the movies count in mid 2020. It shows the rise in popularity of tv shows in recent years.
- Netflix has movies from variety of directors. Around 4993 directors have their movies or tv shows on Netflix.
- Netflix has movies from total 122 countries, United States being the highset contributor with almost 37% of all the content.
- The release year for shows is concentrated in the range 2005-2021.
- 50 mins - 150 mins is the range of movie durations, excluding potential outliers.
- 1-3 seasons is the range for TV shows seasons, excluding potential outliers.
- various ratings of content is avaialble on netfilx, for the various viewers categories like kids, adults , families. Highest number of movies and TV shows are rated TV-MA (for mature audiences).
- Content in most of the ratings is available in lesser quanitity except in US. Ratings like TV-Y7 , TV-Y7 FV , PG ,TV-G , G , TV-Y , TV-PG are very less avaialble in all countries except US.
- International Movies and TV Shows , Dramas , and Comedies are the top 3 genres on Netflix for both Movies and TV shows.
- Mostly country specific popular genres are observed in each country. Only United States have a good mix of almost all genres. Eg. Korean TV shows (Korea), British TV Shows (UK), Anime features and Anime series (Japan) and so on.
- Indian Actors have been acted in maximum movies on netflix. Top 5 actors are in India based on quantity of movies.
- Shorter duration movies have been popular in last 10 years.

## ▼ Business Insights

- Netflix have majority of content which is released after the year 2000. It is observed that the content older than year 2000 is very scarce on Netflix. Senior Citizen could be the target audience for such content, which is almost missing currently.
- Maximum content (more than 80%) is
  - TV-MA - Content intended for mature audiences aged 17 and above.
  - TV-14 - Content suitable for viewers aged 14 and above.
  - TV-PG - Parental guidance suggested (similar ratings - PG-13 , PG)
  - R - Restricted Content, that may not be suitable for viewers under age 17.

These ratings' movies target Matured and Adult audience. Rest 20 % of the content is for kids aged below 13. It shows that Netflix is currently serving mostly Mature audiences or Children with parental guidance.

- Most popular genres on Netflix are International Movies and TV Shows , Dramas , Comedies, Action & Adventure, Children & Family Movies, Thrillers.
- Maximum content of Netflix which is around 75% , is coming from the top 10 countries. Rest of the world only contributes 25% of the content. More countries can be focussed in future to grow the business.
- Liking towards the shorter duration content is on the rise. (duration 75 to 150 minutes and seasons 1 to 3) This can be considered while production of new content on Netflix.
- drop in content is seen across all the countries and type of content in year 2020 and 2021, possibly because of Pandemic.

## ▼ Recommendations

- Very limited genres are focussed in most of the countries except US. It seems the current available genres suits best for US and few countries but maximum countries need some more genres which are highly popular in the region. eg. Indian Mythological content is highly popular. We can create such more country specific genres and It might also be liked acorss the world just like Japanese Anime.

- Country specific insights - The content need to be targetting the demographic of any country. Netflix can produce higher number of content in the perticular rating as per demographic of the country. Eg.
  - The country like India , which is highly populous , has maximum content available only in three rating TV-MA, TV-14 , TV-PG. It is unlikely to serve below 14 age and above 35 year age group .