

In [1]:

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

In [2]:

```
df = pd.read_csv("https://d2bei9khq929f0.cloudfront.net/public_assets/assets/000/000/940/or
```

In [3]:

```
df.head()
```

Out[3]:

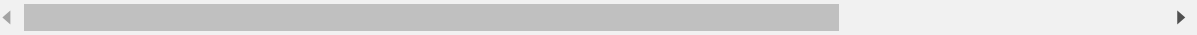
	show_id	type	title	director	cast	country	date_added	release_year	rating
0	s1	Movie	Dick Johnson Is Dead	Kirsten Johnson	NaN	United States	September 25, 2021	2020	PG-13
1	s2	TV Show	Blood & Water	NaN	Ama Qamata, Khosi Ngema, Gail Mabalane, Thaban...	South Africa	September 24, 2021	2021	TV-MA
2	s3	TV Show	Ganglands	Julien Leclercq	Sami Bouajila, Tracy Gotoas, Samuel Jouy, Nabi...	NaN	September 24, 2021	2021	TV-MA
3	s4	TV Show	Jailbirds New Orleans	NaN	NaN	NaN	September 24, 2021	2021	TV-MA
4	s5	TV Show	Kota Factory	NaN	Mayur More, Jitendra Kumar, Ranjan Raj, Alam K...	India	September 24, 2021	2021	TV-MA

In [4]:

```
df.tail()
```

Out[4]:

	show_id	type	title	director	cast	country	date_added	release_year	rat
8802	s8803	Movie	Zodiac	David Fincher	Mark Ruffalo, Jake Gyllenhaal, Robert Downey J...	United States	November 20, 2019	2007	
8803	s8804	TV Show	Zombie Dumb	NaN	NaN	NaN	July 1, 2019	2018	TV
8804	s8805	Movie	Zombieland	Ruben Fleischer	Jesse Eisenberg, Woody Harrelson, Emma Stone, ...	United States	November 1, 2019	2009	
8805	s8806	Movie	Zoom	Peter Hewitt	Tim Allen, Courteney Cox, Chevy Chase, Kate Ma...	United States	January 11, 2020	2006	
8806	s8807	Movie	Zubaan	Mozez Singh	Vicky Kaushal, Sarah-Jane Dias, Raaghav Chanan...	India	March 2, 2019	2015	TV



In [5]:

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 8807 entries, 0 to 8806
Data columns (total 12 columns):
 #   Column          Non-Null Count  Dtype
---  -
 0   show_id         8807 non-null   object
 1   type            8807 non-null   object
 2   title           8807 non-null   object
 3   director        6173 non-null   object
 4   cast            7982 non-null   object
 5   country         7976 non-null   object
 6   date_added      8797 non-null   object
 7   release_year    8807 non-null   int64
 8   rating          8803 non-null   object
 9   duration        8804 non-null   object
10   listed_in       8807 non-null   object
11   description      8807 non-null   object
dtypes: int64(1), object(11)
memory usage: 825.8+ KB
```

In [6]:

```
df.shape
```

Out[6]:

```
(8807, 12)
```

In [7]:

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

Out[7]:

```
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
```

In [8]:

```
df.describe()
```

Out[8]:

	release_year
count	8807.000000
mean	2014.180198
std	8.819312
min	1925.000000
25%	2013.000000
50%	2017.000000
75%	2019.000000
max	2021.000000

In [9]:

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

Out[9]:

```
Movie      6131
TV Show    2676
Name: type, dtype: int64
```

In [92]:

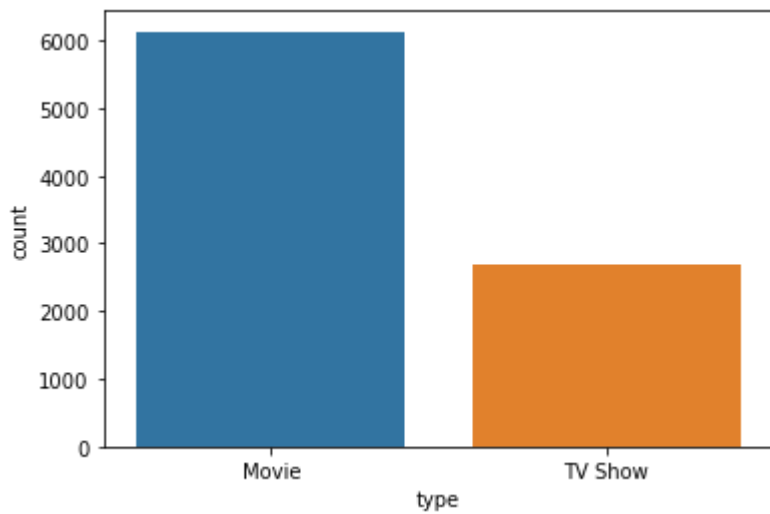
```
df['type'].value_counts(normalize=True)*100
```

Out[92]:

```
Movie      69.615079
TV Show    30.384921
Name: type, dtype: float64
```

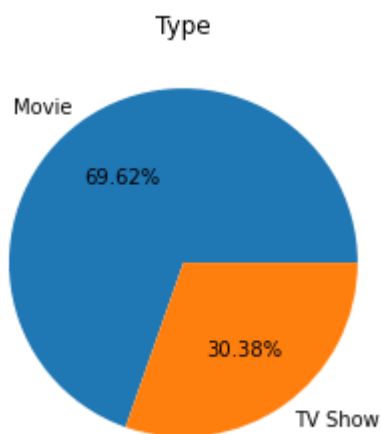
In [10]:

```
sns.countplot(data= df, x= 'type')  
plt.show()
```



In [11]:

```
plt.pie(df['type'].value_counts(), autopct= '%.2f%', labels= ['Movie', 'TV Show'])  
plt.title("Type")  
plt.show()
```



In [97]:

```
df['rating'].value_counts(normalize=True)*100
```

Out[97]:

TV-MA	36.430762
TV-14	24.537090
TV-PG	9.803476
R	9.076451
PG-13	5.566284
TV-Y7	3.794161
TV-Y	3.487447
PG	3.260252
TV-G	2.499148
NR	0.908781
G	0.465750
TV-Y7-FV	0.068159
NC-17	0.034079
UR	0.034079
74 min	0.011360
84 min	0.011360
66 min	0.011360

Name: rating, dtype: float64

In [99]:

```
df[df['rating']=='R']['type'].value_counts()
```

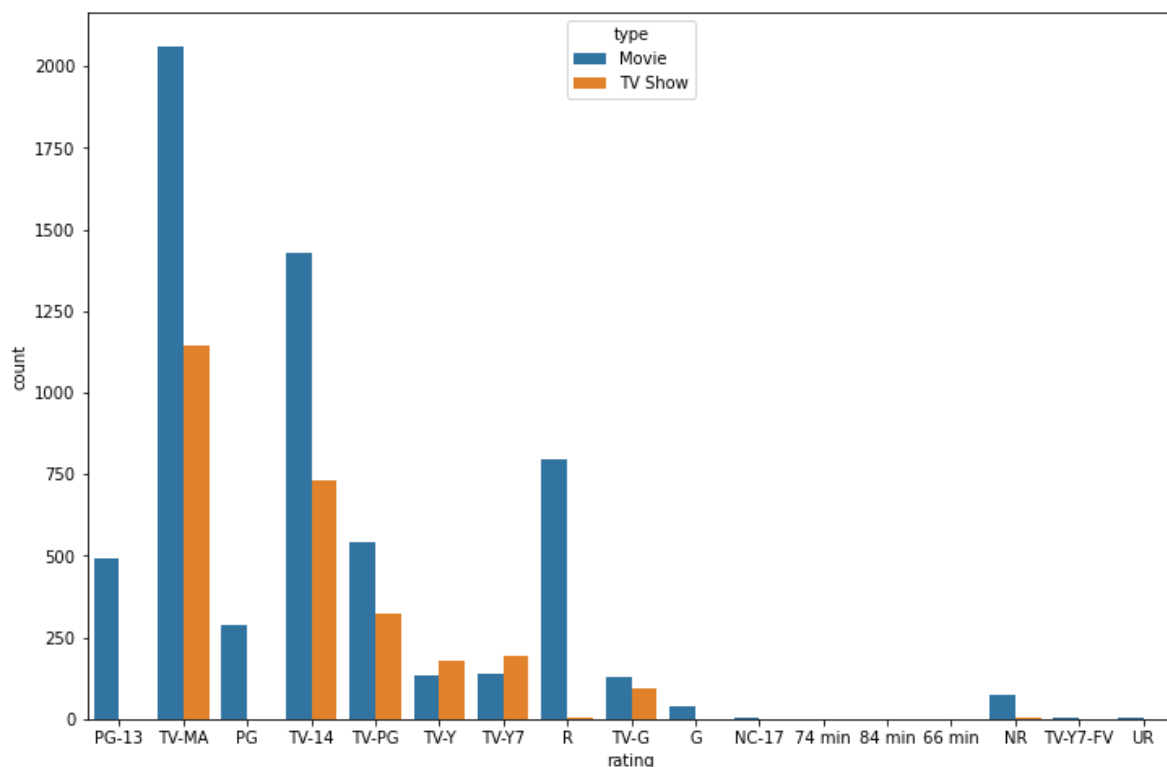
Out[99]:

Movie	797
TV Show	2

Name: type, dtype: int64

In [12]:

```
plt.figure(figsize=(12,8))
sns.countplot(data= df, x= 'rating' , hue= 'type')
plt.show()
```



In [101]:

```
df['release_year'].value_counts(normalize=True)*100
```

Out[101]:

```
2018    13.023731
2017    11.717952
2019    11.695242
2020    10.820938
2016    10.241853
```

```
...
1959    0.011355
1925    0.011355
1961    0.011355
1947    0.011355
1966    0.011355
```

```
Name: release_year, Length: 74, dtype: float64
```

In []:

In []:

In [15]:

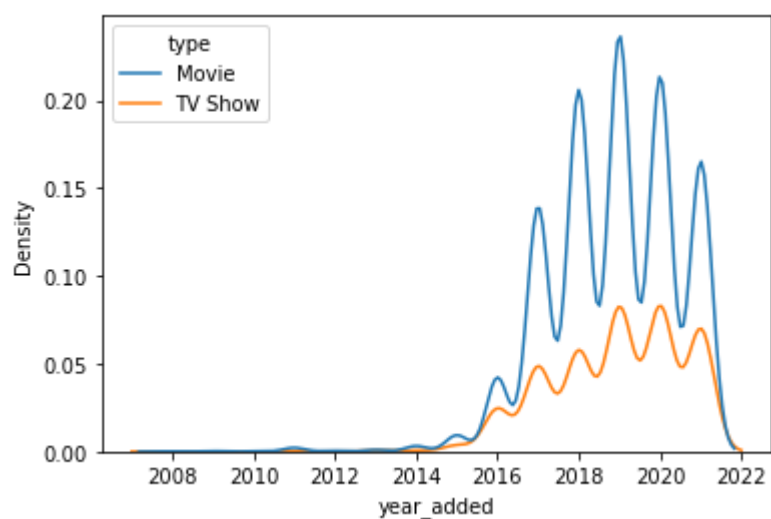
```
df['year_added'] = pd.to_datetime(df['date_added']).dt.year
# df['year_added'] = df['year_added'][df['year_added'].isna()==False].astype(int)
df.head()
```

Out[15]:

	show_id	type	title	director	cast	country	date_added	release_year	rating
0	s1	Movie	Dick Johnson Is Dead	Kirsten Johnson	NaN	United States	September 25, 2021	2020	PG-13
1	s2	TV Show	Blood & Water	NaN	Ama Qamata, Khosi Ngema, Gail Mabalane, Thaban...	South Africa	September 24, 2021	2021	TV-MA
2	s3	TV Show	Ganglands	Julien Leclercq	Sami Bouajila, Tracy Gotoas, Samuel Jouy, Nabi...	NaN	September 24, 2021	2021	TV-MA
3	s4	TV Show	Jailbirds New Orleans	NaN	NaN	NaN	September 24, 2021	2021	TV-MA
4	s5	TV Show	Kota Factory	NaN	Mayur More, Jitendra Kumar, Ranjan Raj, Alam K...	India	September 24, 2021	2021	TV-MA

In [16]:

```
sns.kdeplot(data=df, x='year_added' , hue='type')  
plt.show()
```



In [17]:

```
df[(df['release_year'] >=1990) & (df['release_year'] <= 2022)]['release_year'].value_counts
```

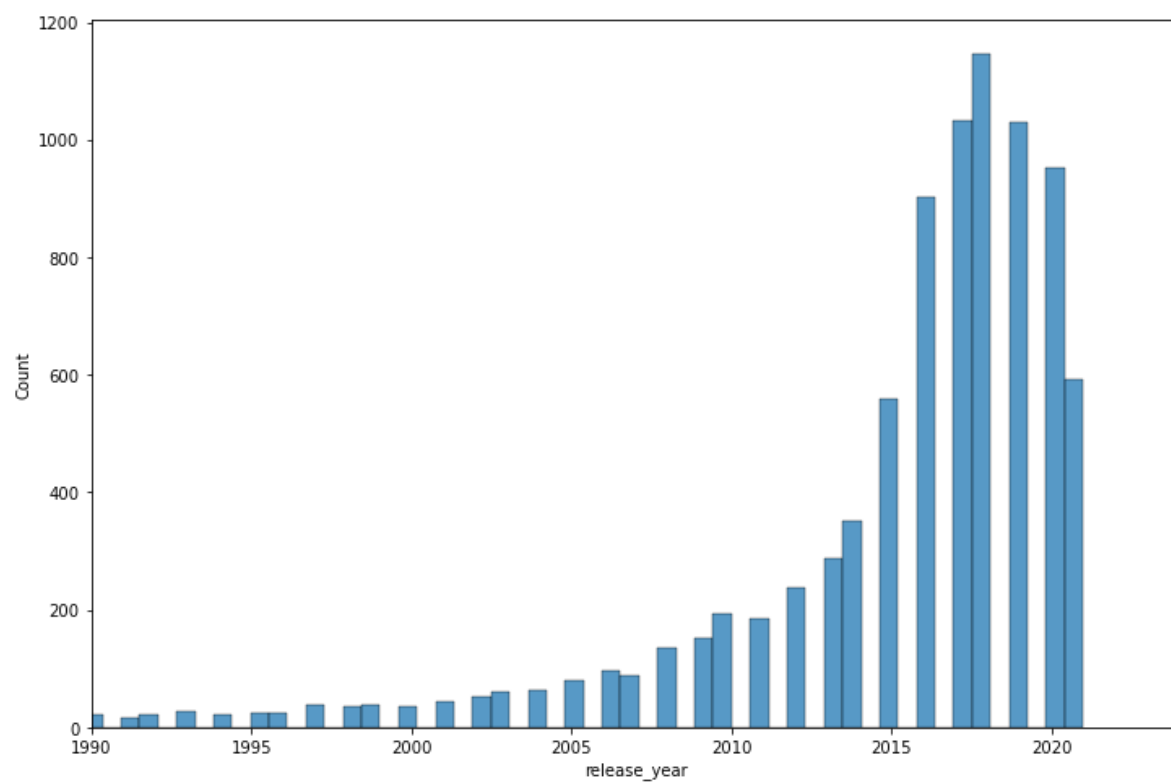
Out[17]:

2018	1147
2017	1032
2019	1030
2020	953
2016	902
2021	592
2015	560
2014	352
2013	288
2012	237
2010	194
2011	185
2009	152
2008	136
2006	96
2007	88
2005	80
2004	64
2003	61
2002	51
2001	45
1999	39
1997	38
2000	37
1998	36
1993	28
1995	25
1996	24
1992	23
1990	22
1994	22
1991	17

Name: release_year, dtype: int64

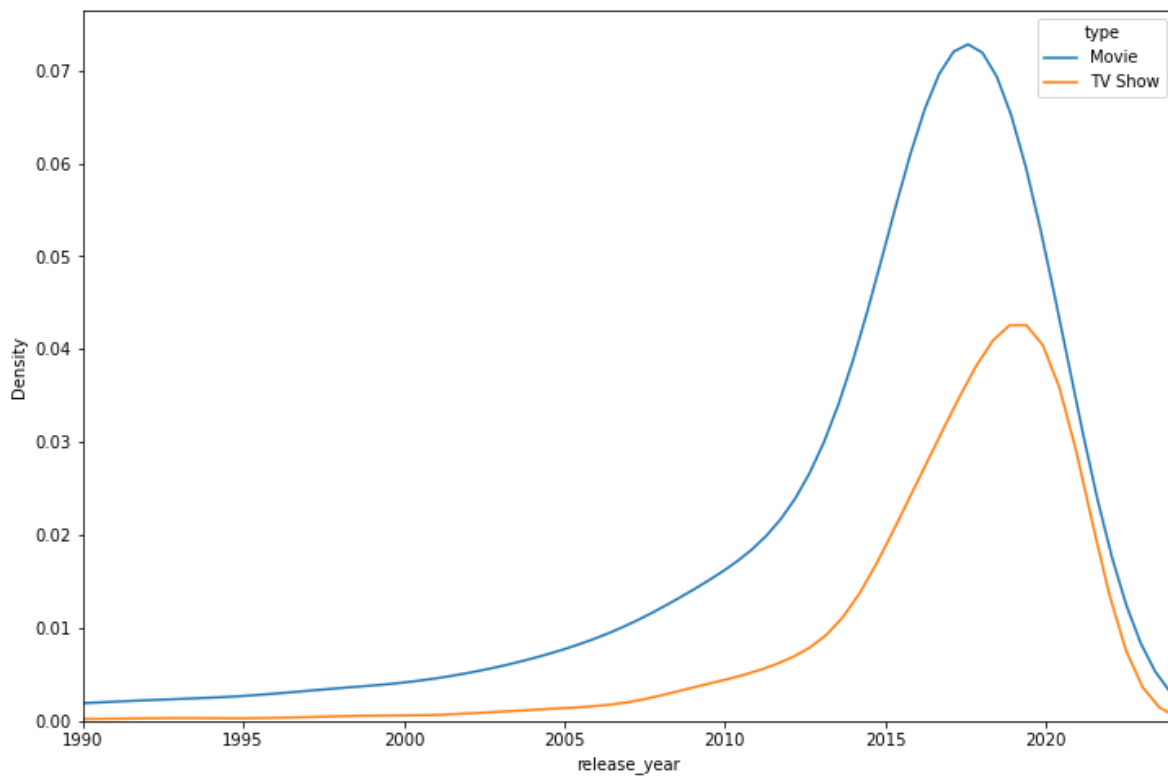
In [18]:

```
plt.figure(figsize=(12,8))  
sns.histplot(data= df, x= 'release_year')  
plt.xlim(1990,2024)  
plt.show()
```



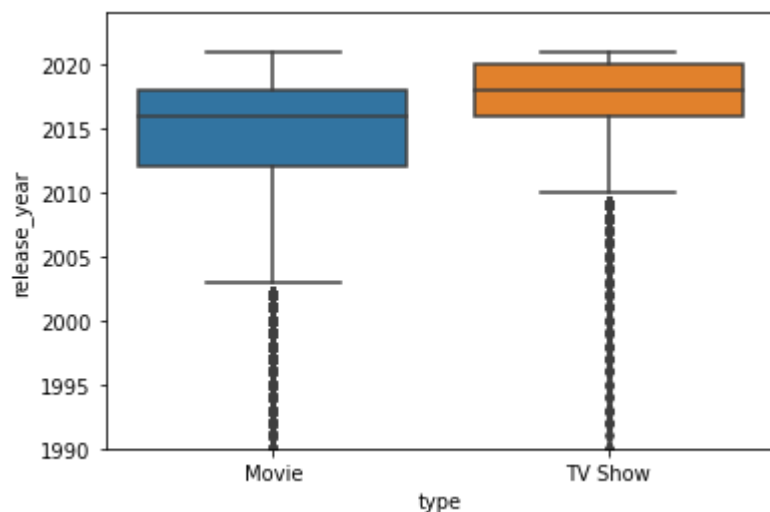
In [19]:

```
plt.figure(figsize=(12,8))
sns.kdeplot(data= df, x= 'release_year', hue = 'type')
plt.xlim(1990,2024)
plt.show()
```



In [20]:

```
sns.boxplot(data= df, y= 'release_year' , x= 'type')
plt.ylim(1990,2024)
plt.show()
```



In []:

Working on cast

In [21]:

```
# Seperating the cast members into columns for each title
constraint=df['cast'].apply(lambda x: str(x).split(',')).tolist()
cast_expand=pd.DataFrame(constraint,index=df['title'])
```

In [22]:

```
#stacking all cast columns into rows for each title
cast_expand = cast_expand.stack()
cast_expand = pd.DataFrame(cast_expand, columns=['cast'])
```

In [23]:

```
cast_expand
```

Out[23]:

title		cast
Dick Johnson Is Dead	0	nan
Blood & Water	0	Ama Qamata
	1	Khosi Ngema
	2	Gail Mabalane
	3	Thabang Molaba
...
Zubaan	3	Manish Chaudhary
	4	Meghna Malik
	5	Malkeet Rauni
	6	Anita Shabdish
	7	Chittaranjan Tripathy

64951 rows × 1 columns

In [24]:

```
# top 10 cast members with most content  
cast_expand.value_counts().head(11)
```

Out[24]:

```
cast  
nan                825  
Anupam Kher        43  
Shah Rukh Khan     35  
Julie Teiwani      33  
Naseeruddin Shah   32  
Takahiro Sakurai   32  
Rupa Bhimani       31  
Om Puri            30  
Akshay Kumar       30  
Yuki Kaji          29  
Amitabh Bachchan   28  
dtype: int64
```

In [25]:

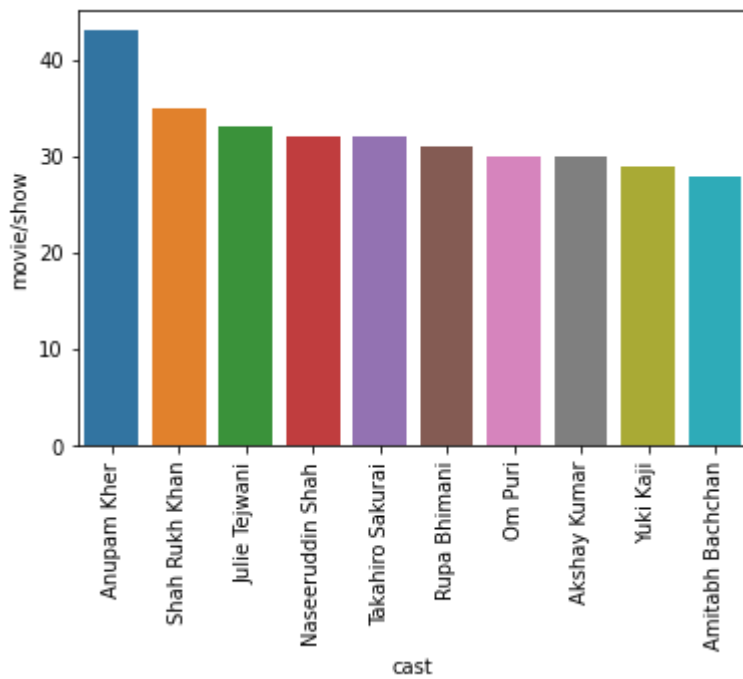
```
top10_cast = cast_expand.value_counts()[1:11]  
top10_cast = pd.DataFrame(top10_cast, columns = ['movie/show']).reset_index()  
top10_cast
```

Out[25]:

	cast	movie/show
0	Anupam Kher	43
1	Shah Rukh Khan	35
2	Julie Teiwani	33
3	Naseeruddin Shah	32
4	Takahiro Sakurai	32
5	Rupa Bhimani	31
6	Om Puri	30
7	Akshay Kumar	30
8	Yuki Kaji	29
9	Amitabh Bachchan	28

In [26]:

```
sns.barplot(data= top10_cast, x= 'cast' , y='movie/show')  
plt.xticks(rotation = 90)  
plt.show()
```



In [27]:

```
temp = df[['title','type']].set_index('title')
temp
```

Out[27]:

type	
title	
Dick Johnson Is Dead	Movie
Blood & Water	TV Show
Ganglands	TV Show
Jailbirds New Orleans	TV Show
Kota Factory	TV Show
...	...
Zodiac	Movie
Zombie Dumb	TV Show
Zombieland	Movie
Zoom	Movie
Zubaan	Movie

8807 rows × 1 columns

In [28]:

```
categorised_cast = pd.merge(cast_expand, temp, left_index=True, right_index=True)
categorised_cast
```

Out[28]:

		cast	type
title			
Dick Johnson Is Dead	0	nan	Movie
Blood & Water	0	Ama Qamata	TV Show
	1	Khosi Ngema	TV Show
	2	Gail Mabalane	TV Show
	3	Thabang Molaba	TV Show
...
Zubaan	3	Manish Chaudhary	Movie
	4	Meghna Malik	Movie
	5	Malkeet Rauni	Movie
	6	Anita Shabdish	Movie
	7	Chittaranjan Tripathy	Movie

64951 rows × 2 columns

In [29]:

```
top_movie_cast = pd.DataFrame(categorised_cast[categorised_cast['type'] == 'Movie'].value_counts())
top_movie_cast
```

Out[29]:

		No. of Movies
cast	type	
Anupam Kher	Movie	42
Shah Rukh Khan	Movie	35
Naseeruddin Shah	Movie	32
Om Puri	Movie	30
Akshay Kumar	Movie	30
Amitabh Bachchan	Movie	28
Julie Tejjwani	Movie	28
Paresh Rawal	Movie	28
Rupa Bhimani	Movie	27
Boman Irani	Movie	27

In [30]:

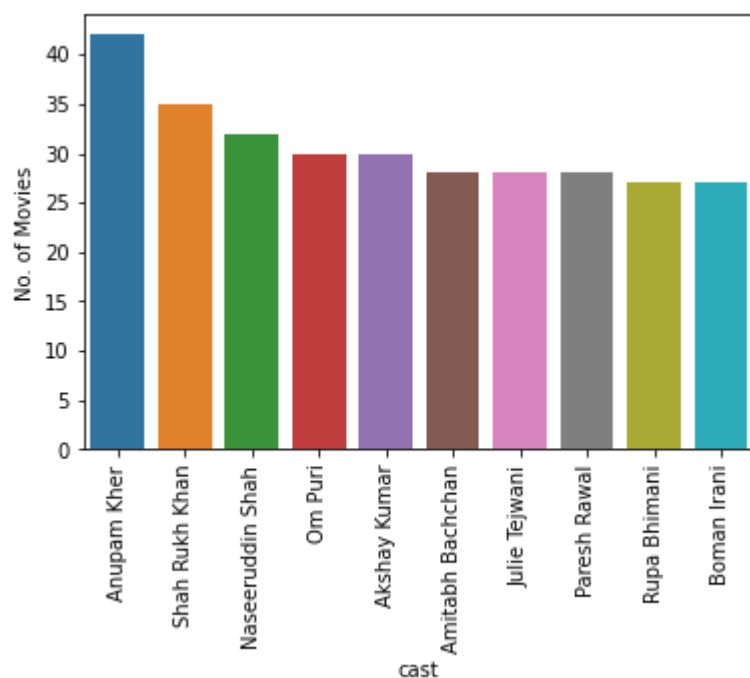
```
top_movie_cast.reset_index(1,drop=True,inplace=True)
top_movie_cast.reset_index(inplace = True)
top_movie_cast
```

Out[30]:

	cast	No. of Movies
0	Anupam Kher	42
1	Shah Rukh Khan	35
2	Naseeruddin Shah	32
3	Om Puri	30
4	Akshay Kumar	30
5	Amitabh Bachchan	28
6	Julie Teiwani	28
7	Paresh Rawal	28
8	Rupa Bhimani	27
9	Boman Irani	27

In [31]:

```
sns.barplot(data=top_movie_cast, x= 'cast', y='No. of Movies')
plt.xticks(rotation = 90)
plt.show()
```



In [32]:

```
top_show_cast = pd.DataFrame(categorised_cast[categorised_cast['type'] == 'TV Show'].value_  
top_show_cast
```

Out[32]:

No. of TV Shows		
cast	type	
Takahiro Sakurai	TV Show	25
Yuki Kaji	TV Show	19
Daisuke Ono	TV Show	17
Ai Kayano	TV Show	17
Junichi Suwabe	TV Show	17
Yuichi Nakamura	TV Show	16
Yoshimasa Hosoya	TV Show	15
Jun Fukuyama	TV Show	15
David Attenborough	TV Show	14
Vincent Tong	TV Show	13

In [33]:

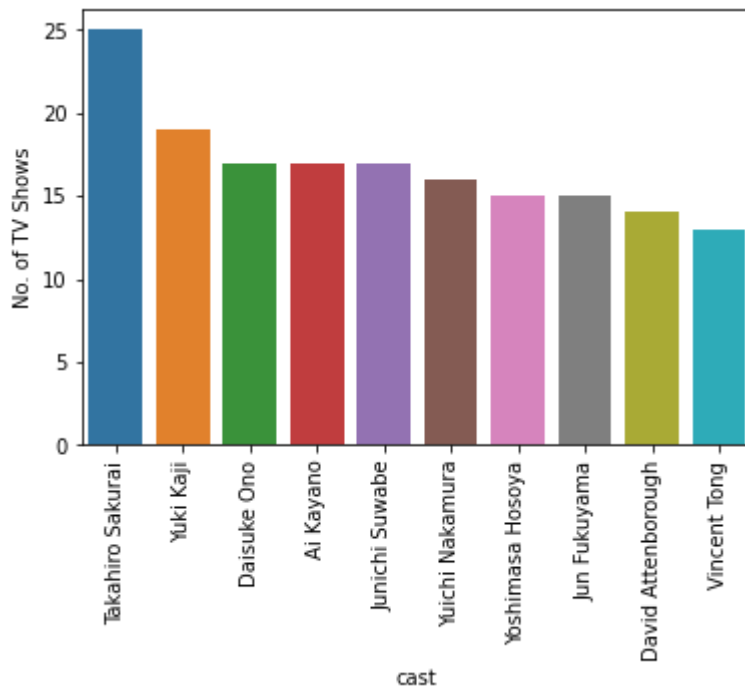
```
top_show_cast.reset_index(1,drop=True,inplace=True)  
top_show_cast.reset_index(inplace = True)  
top_show_cast
```

Out[33]:

	cast	No. of TV Shows
0	Takahiro Sakurai	25
1	Yuki Kaji	19
2	Daisuke Ono	17
3	Ai Kayano	17
4	Junichi Suwabe	17
5	Yuichi Nakamura	16
6	Yoshimasa Hosoya	15
7	Jun Fukuyama	15
8	David Attenborough	14
9	Vincent Tong	13

In [34]:

```
sns.barplot(data=top_show_cast, x= 'cast', y='No. of TV Shows')  
plt.xticks(rotation = 90)  
plt.show()
```



In []:

In [35]:

```
temp2 = df[['title', 'year_added']]
temp2
```

Out[35]:

	title	year_added
0	Dick Johnson Is Dead	2021.0
1	Blood & Water	2021.0
2	Ganglands	2021.0
3	Jailbirds New Orleans	2021.0
4	Kota Factory	2021.0
...
8802	Zodiac	2019.0
8803	Zombie Dumb	2019.0
8804	Zombieland	2019.0
8805	Zoom	2020.0
8806	Zubaan	2019.0

8807 rows × 2 columns

In [36]:

```
cast = cast_expand.reset_index(1, drop=True)
```

In [37]:

```
df2= pd.merge(cast , temp2 , left_index=True , right_on='title')
df2
```

Out[37]:

	cast	title	year_added
0	nan	Dick Johnson Is Dead	2021.0
1	Ama Qamata	Blood & Water	2021.0
1	Khosi Ngema	Blood & Water	2021.0
1	Gail Mabalane	Blood & Water	2021.0
1	Thabang Molaba	Blood & Water	2021.0
...
8806	Manish Chaudhary	Zubaan	2019.0
8806	Meghna Malik	Zubaan	2019.0
8806	Malkeet Rauni	Zubaan	2019.0
8806	Anita Shabdish	Zubaan	2019.0
8806	Chittaranjan Tripathy	Zubaan	2019.0

64951 rows × 3 columns

In [38]:

```
df2 = df2.sort_index()[['title', 'year_added', 'cast']]
df2
```

Out[38]:

	title	year_added	cast
0	Dick Johnson Is Dead	2021.0	nan
1	Blood & Water	2021.0	Ama Qamata
1	Blood & Water	2021.0	Khosi Ngema
1	Blood & Water	2021.0	Gail Mabalane
1	Blood & Water	2021.0	Thabang Molaba
...
8806	Zubaan	2019.0	Manish Chaudhary
8806	Zubaan	2019.0	Meghna Malik
8806	Zubaan	2019.0	Malkeet Rauni
8806	Zubaan	2019.0	Anita Shabdish
8806	Zubaan	2019.0	Chittaranjan Tripathy

64951 rows × 3 columns

In [39]:

```
grp1 = df2.groupby(by='cast')['year_added'].value_counts().sort_values(ascending = False)
```

In [40]:

```
grp1 = pd.DataFrame(grp1)
grp1.head()
```

Out[40]:

	year_added	
cast	year_added	
nan	2019.0	165
	2020.0	155
	2021.0	150
	2018.0	150
	2017.0	140

In [41]:

```
grp1 = grp1.drop(index= 'nan')
```

C:\Users\dgoya\anaconda3\lib\site-packages\pandas\core\generic.py:4150: PerformanceWarning: dropping on a non-lexsorted multi-index without a level parameter may impact performance.

```
obj = obj._drop_axis(labels, axis, level=level, errors=errors)
```

In [42]:

```
grp1.columns = ['No. of Movies/Shows']  
grp1.reset_index(1)
```

Out[42]:

	year_added	No. of Movies/Shows
cast		
Julie Tejjwani	2021.0	22
Rupa Bhimani	2021.0	22
Rajesh Kava	2021.0	21
Anupam Kher	2018.0	19
Jigna Bhardwaj	2021.0	19
...
Ibrahim Suleiman	2021.0	1
Ibrahima Gueye	2020.0	1
Ibrahima Mbaye	2019.0	1
Ibrahima Traore	2019.0	1
Şopé Dirisù	2020.0	1

50860 rows × 2 columns

In [43]:

```
per_year = grp1.loc[top10_cast['cast']]
per_year
```

Out[43]:

No. of Movies/Shows		
cast	year_added	
Anupam Kher	2018.0	19
	2020.0	10
	2021.0	5
	2017.0	5
	2019.0	4
Shah Rukh Khan	2017.0	14
	2018.0	11
	2020.0	5
	2019.0	3
	2021.0	2
Julie Tejawani	2021.0	22
	2019.0	9
	2020.0	2
Naseeruddin Shah	2019.0	10
	2018.0	9
	2020.0	6
	2017.0	4
	2021.0	3
Takahiro Sakurai	2019.0	11
	2020.0	7
	2021.0	4
	2016.0	4
	2017.0	3
	2018.0	3
Rupa Bhimani	2021.0	22
	2019.0	9
Om Puri	2018.0	14
	2019.0	7
	2020.0	6
	2017.0	3
Akshay Kumar	2018.0	11
	2020.0	8
	2019.0	5

No. of Movies/Shows

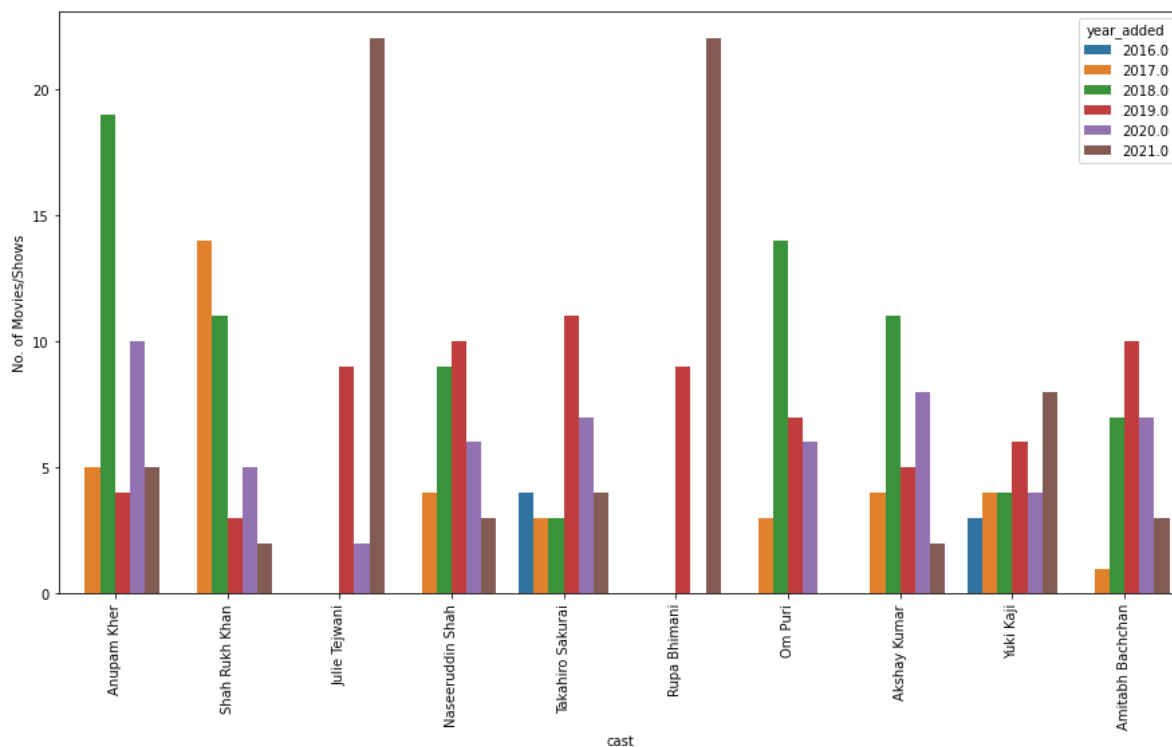
cast	year_added	
Yuki Kaji	2017.0	4
	2021.0	2
	2021.0	8
	2019.0	6
	2017.0	4
	2020.0	4
	2018.0	4
	2016.0	3
Amitabh Bachchan	2019.0	10
	2020.0	7
	2018.0	7
	2021.0	3
	2017.0	1

In [44]:

```
per_year.reset_index(inplace=True)
```

In [45]:

```
plt.figure(figsize=(15,8))
sns.barplot(data=per_year, x='cast', y='No. of Movies/Shows', hue='year_added')
plt.xticks(rotation=90)
plt.show()
```



In []:

In [46]:

df.head()

Out[46]:

	show_id	type	title	director	cast	country	date_added	release_year	rating
0	s1	Movie	Dick Johnson Is Dead	Kirsten Johnson	NaN	United States	September 25, 2021	2020	PG-13
1	s2	TV Show	Blood & Water	NaN	Ama Qamata, Khosi Ngema, Gail Mabalane, Thaban...	South Africa	September 24, 2021	2021	TV-MA
2	s3	TV Show	Ganglands	Julien Leclercq	Sami Bouajila, Tracy Gotoas, Samuel Jouy, Nabi...	NaN	September 24, 2021	2021	TV-MA
3	s4	TV Show	Jailbirds New Orleans	NaN	NaN	NaN	September 24, 2021	2021	TV-MA
4	s5	TV Show	Kota Factory	NaN	Mayur More, Jitendra Kumar, Ranjan Raj, Alam K...	India	September 24, 2021	2021	TV-MA

In [47]:

```
# Separating the directors into columns for each title
constraint2 = df['director'].apply(lambda x: str(x).split(", ")).to_list()
director_expand = pd.DataFrame(constraint2, index=df['title'])
```

In [48]:

```
#stacking all director columns into rows for each title
director_expand = director_expand.stack()
director_expand = pd.DataFrame(director_expand, columns=['director'])
```

In [49]:

director_expand

Out[49]:

		director
title		
Dick Johnson Is Dead	0	Kirsten Johnson
Blood & Water	0	nan
Ganglands	0	Julien Leclercq
Jailbirds New Orleans	0	nan
Kota Factory	0	nan
...
Zodiac	0	David Fincher
Zombie Dumb	0	nan
Zombieland	0	Ruben Fleischer
Zoom	0	Peter Hewitt
Zubaan	0	Mozez Singh

9612 rows × 1 columns

In [50]:

```
# top 10 directors with most content
director_expand.value_counts().head(11)
```

Out[50]:

```
director
nan                2634
Rajiv Chilaka      22
Jan Suter           21
Raúl Campos        19
Suhas Kadav        16
Marcus Raboy       16
Jay Karas           15
Cathy Garcia-Molina 13
Martin Scorsese     12
Youssef Chahine     12
Jay Chapman         12
dtype: int64
```

In [51]:

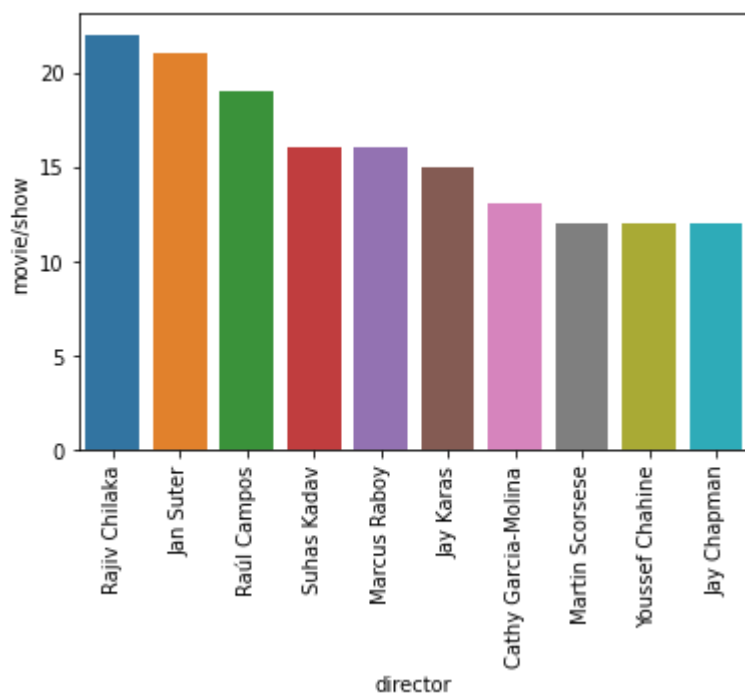
```
top10_director = director_expand.value_counts()[1:11]
top10_director = pd.DataFrame(top10_director, columns = ['movie/show']).reset_index()
top10_director
```

Out[51]:

	director	movie/show
0	Rajiv Chilaka	22
1	Jan Suter	21
2	Raúl Campos	19
3	Suhas Kadav	16
4	Marcus Raboy	16
5	Jay Karas	15
6	Cathy Garcia-Molina	13
7	Martin Scorsese	12
8	Youssef Chahine	12
9	Jay Chapman	12

In [52]:

```
sns.barplot(data= top10_director, x= 'director' , y='movie/show')
plt.xticks(rotation = 90)
plt.show()
```



In []:

In [53]:

```
categorised_director = pd.merge(director_expand, temp, left_index=True, right_index= True)
categorised_director
```

Out[53]:

		director	type
title			
Dick Johnson Is Dead	0	Kirsten Johnson	Movie
Blood & Water	0	nan	TV Show
Ganglands	0	Julien Leclercq	TV Show
Jailbirds New Orleans	0	nan	TV Show
Kota Factory	0	nan	TV Show
...
Zodiac	0	David Fincher	Movie
Zombie Dumb	0	nan	TV Show
Zombieland	0	Ruben Fleischer	Movie
Zoom	0	Peter Hewitt	Movie
Zubaan	0	Mozez Singh	Movie

9612 rows × 2 columns

In [54]:

```
top_movie_dir = pd.DataFrame(categorised_director[categorised_director['type'] == 'Movie'].
top_movie_dir
```

Out[54]:

No. of Movies		
director	type	
Rajiv Chilaka	Movie	22
Jan Suter	Movie	21
Raúl Campos	Movie	19
Suhas Kadav	Movie	16
Marcus Raboy	Movie	15
Jay Karas	Movie	15
Cathy Garcia-Molina	Movie	13
Youssef Chahine	Movie	12
Martin Scorsese	Movie	12
Jay Chapman	Movie	12

In [55]:

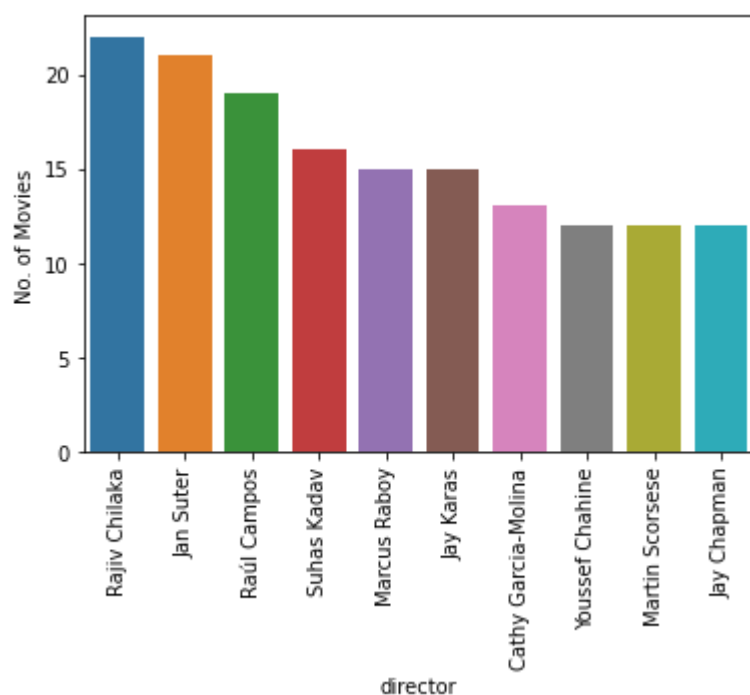
```
top_movie_dir.reset_index(1,drop=True,inplace=True)
top_movie_dir.reset_index(inplace = True)
top_movie_dir
```

Out[55]:

	director	No. of Movies
0	Rajiv Chilaka	22
1	Jan Suter	21
2	Raúl Campos	19
3	Suhas Kadav	16
4	Marcus Raboy	15
5	Jay Karas	15
6	Cathy Garcia-Molina	13
7	Youssef Chahine	12
8	Martin Scorsese	12
9	Jay Chapman	12

In [56]:

```
sns.barplot(data=top_movie_dir, x= 'director', y='No. of Movies')
plt.xticks(rotation = 90)
plt.show()
```



In [57]:

```
top_show_dir = pd.DataFrame(categorised_director[categorised_director['type'] == 'TV Show'])  
top_show_dir
```

Out[57]:

No. of TV Shows		
director	type	
Ken Burns	TV Show	3
Alastair Fothergill	TV Show	3
Jung-ah Im	TV Show	2
Joe Berlinger	TV Show	2
Hsu Fu-chun	TV Show	2
Stan Lathan	TV Show	2
Gautham Vasudev Menon	TV Show	2
Lynn Novick	TV Show	2
Shin Won-ho	TV Show	2
Iginio Straffi	TV Show	2

In [58]:

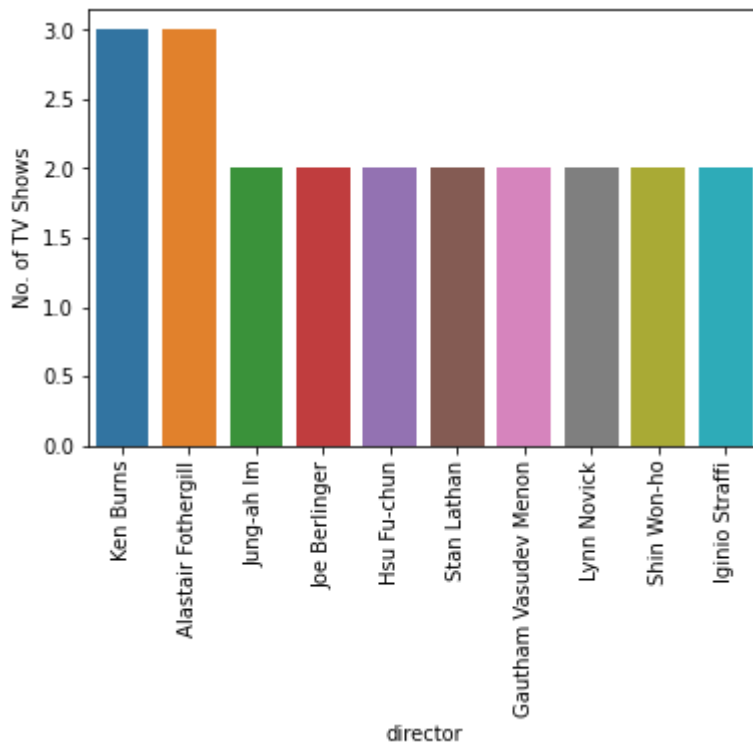
```
top_show_dir.reset_index(1,drop=True,inplace=True)  
top_show_dir.reset_index(inplace = True)  
top_show_dir
```

Out[58]:

	director	No. of TV Shows
0	Ken Burns	3
1	Alastair Fothergill	3
2	Jung-ah Im	2
3	Joe Berlinger	2
4	Hsu Fu-chun	2
5	Stan Lathan	2
6	Gautham Vasudev Menon	2
7	Lynn Novick	2
8	Shin Won-ho	2
9	Iginio Straffi	2

In [59]:

```
sns.barplot(data=top_show_dir, x= 'director', y='No. of TV Shows')  
plt.xticks(rotation = 90)  
plt.show()
```



In []:

In [60]:

```
dire = director_expand.reset_index(1,drop=True)
```

In [61]:

```
df3= pd.merge(dire , temp2 , left_index=True , right_on='title')
df3
```

Out[61]:

	director	title	year_added
0	Kirsten Johnson	Dick Johnson Is Dead	2021.0
1	nan	Blood & Water	2021.0
2	Julien Leclercq	Ganglands	2021.0
3	nan	Jailbirds New Orleans	2021.0
4	nan	Kota Factory	2021.0
...
8802	David Fincher	Zodiac	2019.0
8803	nan	Zombie Dumb	2019.0
8804	Ruben Fleischer	Zombieland	2019.0
8805	Peter Hewitt	Zoom	2020.0
8806	Mozes Singh	Zubaan	2019.0

9612 rows × 3 columns

In [62]:

```
df3 = df3.sort_index()[['title', 'year_added' , 'director']]
df3
```

Out[62]:

	title	year_added	director
0	Dick Johnson Is Dead	2021.0	Kirsten Johnson
1	Blood & Water	2021.0	nan
2	Ganglands	2021.0	Julien Leclercq
3	Jailbirds New Orleans	2021.0	nan
4	Kota Factory	2021.0	nan
...
8802	Zodiac	2019.0	David Fincher
8803	Zombie Dumb	2019.0	nan
8804	Zombieland	2019.0	Ruben Fleischer
8805	Zoom	2020.0	Peter Hewitt
8806	Zubaan	2019.0	Mozes Singh

9612 rows × 3 columns

In [63]:

```
grp2 = df3.groupby(by='director')['year_added'].value_counts().sort_values(ascending = False)
```

In [64]:

```
grp2 = pd.DataFrame(grp2)  
grp2.head()
```

Out[64]:

		year_added
director	year_added	
nan	2019.0	598
	2020.0	564
	2021.0	470
	2018.0	435
	2017.0	334

In [65]:

```
grp2 = grp2.drop(index= 'nan')
```

C:\Users\dgoya\anaconda3\lib\site-packages\pandas\core\generic.py:4150: PerformanceWarning: dropping on a non-lexsorted multi-index without a level parameter may impact performance.

```
obj = obj._drop_axis(labels, axis, level=level, errors=errors)
```

In [66]:

```
grp2.columns = ['No. of Movies/Shows']  
grp2.reset_index(1)
```

Out[66]:

	year_added	No. of Movies/Shows
director		
Rajiv Chilaka	2021.0	17
Suhas Kadav	2021.0	15
Raúl Campos	2018.0	12
Jan Suter	2018.0	12
Youssef Chahine	2020.0	11
...
Huang Hsin-Yao	2021.0	1
Hua Shan	2018.0	1
Hsu Chih-yen	2021.0	1
Hsu Chih-yen	2019.0	1
Şenol Sönmez	2021.0	1

5982 rows × 2 columns

In [67]:

```
per_year1 = grp2.loc[top10_director['director']]
per_year1
```

Out[67]:

No. of Movies/Shows		
director	year_added	
Rajiv Chilaka	2021.0	17
	2019.0	3
	2020.0	2
Jan Suter	2018.0	12
	2017.0	5
	2016.0	4
Raúl Campos	2018.0	12
	2017.0	4
	2016.0	3
Suhas Kadav	2021.0	15
	2017.0	1
Marcus Raboy	2017.0	6
	2018.0	4
	2019.0	3
	2020.0	2
	2016.0	1
Jay Karas	2016.0	4
	2018.0	3
	2015.0	2
	2017.0	2
	2019.0	2
	2020.0	1
Cathy Garcia-Molina	2014.0	1
	2019.0	7
	2020.0	6
Martin Scorsese	2019.0	7
	2021.0	3
	2020.0	2
Youssef Chahine	2020.0	11
	2021.0	1
Jay Chapman	2017.0	7
	2019.0	2
	2018.0	1

No. of Movies/Shows

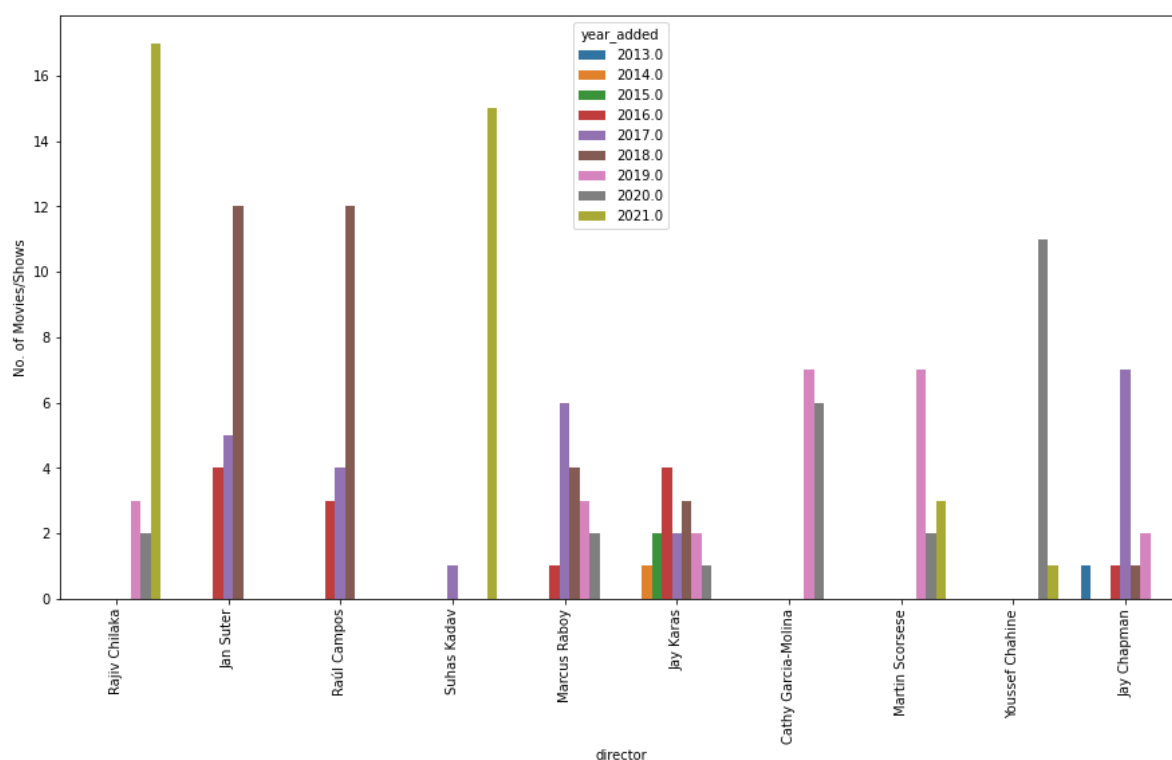
director	year_added
	2016.0
	1
	2013.0
	1

In [68]:

```
per_year1.reset_index(inplace=True)
```

In [69]:

```
plt.figure(figsize=(15,8))
sns.barplot(data=per_year1, x='director', y='No. of Movies/Shows', hue='year_added')
plt.xticks(rotation=90)
plt.show()
```



In []:

In []:

Working on Listed in

In [70]:

df.head()

Out[70]:

	show_id	type	title	director	cast	country	date_added	release_year	rating
0	s1	Movie	Dick Johnson Is Dead	Kirsten Johnson	NaN	United States	September 25, 2021	2020	PG-13
1	s2	TV Show	Blood & Water	NaN	Ama Qamata, Khosi Ngema, Gail Mabalane, Thaban...	South Africa	September 24, 2021	2021	TV-MA
2	s3	TV Show	Ganglands	Julien Leclercq	Sami Bouajila, Tracy Gotoas, Samuel Jouy, Nabi...	NaN	September 24, 2021	2021	TV-MA
3	s4	TV Show	Jailbirds New Orleans	NaN	NaN	NaN	September 24, 2021	2021	TV-MA
4	s5	TV Show	Kota Factory	NaN	Mayur More, Jitendra Kumar, Ranjan Raj, Alam K...	India	September 24, 2021	2021	TV-MA

In [71]:

```
# Separating the genres into columns for each title
constraint3 = df['listed_in'].apply(lambda x: str(x).split(", ")).to_list()
genre_expand = pd.DataFrame(constraint3 , index=df['title'])
```

In [72]:

```
#stacking all genre columns into rows for each title
genre_expand = genre_expand.stack()
genre_expand = pd.DataFrame(genre_expand, columns=['genre'])
```

In [73]:

```
genre_expand
```

Out[73]:

		genre
title		
Dick Johnson Is Dead	0	Documentaries
Blood & Water	0	International TV Shows
	1	TV Dramas
	2	TV Mysteries
Ganglands	0	Crime TV Shows
...
Zoom	0	Children & Family Movies
	1	Comedies
Zubaan	0	Dramas
	1	International Movies
	2	Music & Musicals

19323 rows × 1 columns

In [74]:

```
genre_expand.value_counts().head(10)
```

Out[74]:

genre	
International Movies	2752
Dramas	2427
Comedies	1674
International TV Shows	1351
Documentaries	869
Action & Adventure	859
TV Dramas	763
Independent Movies	756
Children & Family Movies	641
Romantic Movies	616
dtype: int64	

In [75]:

```
(genre_expand.value_counts(normalize=True)*100)[:10]
```

Out[75]:

```
genre
International Movies    14.242095
Dramas                 12.560161
Comedies               8.663251
International TV Shows  6.991668
Documentaries          4.497231
Action & Adventure     4.445479
TV Dramas              3.948662
Independent Movies     3.912436
Children & Family Movies 3.317290
Romantic Movies        3.187911
dtype: float64
```

In [76]:

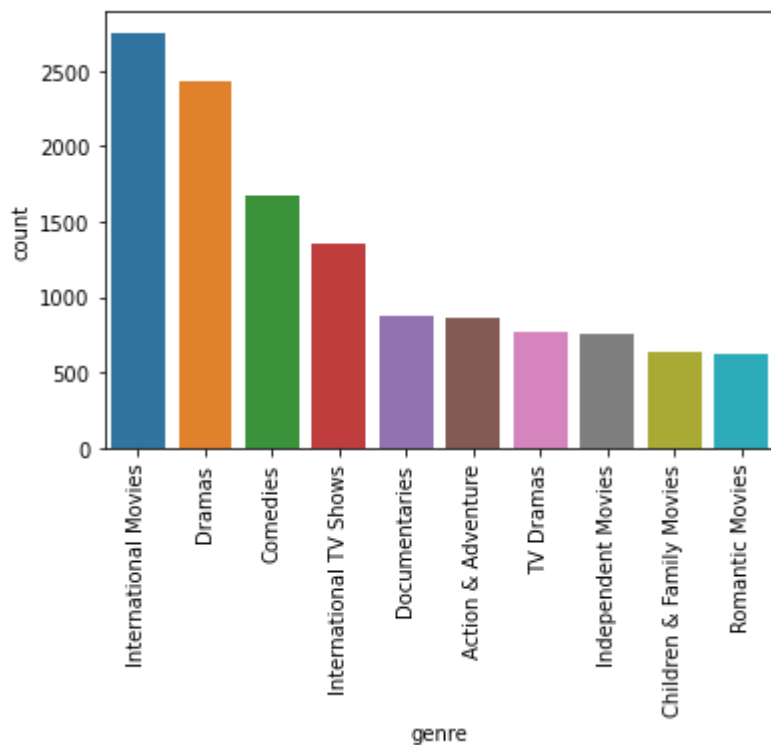
```
top_genre = pd.DataFrame(genre_expand.value_counts()[:10], columns=['count'])
top_genre.reset_index(inplace=True)
top_genre
```

Out[76]:

	genre	count
0	International Movies	2752
1	Dramas	2427
2	Comedies	1674
3	International TV Shows	1351
4	Documentaries	869
5	Action & Adventure	859
6	TV Dramas	763
7	Independent Movies	756
8	Children & Family Movies	641
9	Romantic Movies	616

In [77]:

```
sns.barplot(data=top_genre, x="genre", y="count")  
plt.xticks(rotation = 90)  
plt.show()
```



In []:

In [78]:

```
genre = genre_expand.reset_index(1,drop=True)
```


In [79]:

```
df4= pd.merge(genre , temp2 , left_index=True , right_on='title')
df4
```

Out[79]:

	genre	title	year_added
0	Documentaries	Dick Johnson Is Dead	2021.0
1	International TV Shows	Blood & Water	2021.0
1	TV Dramas	Blood & Water	2021.0
1	TV Mysteries	Blood & Water	2021.0
2	Crime TV Shows	Ganglands	2021.0
...
8805	Children & Family Movies	Zoom	2020.0
8805	Comedies	Zoom	2020.0
8806	Dramas	Zubaan	2019.0
8806	International Movies	Zubaan	2019.0
8806	Music & Musicals	Zubaan	2019.0

19323 rows × 3 columns

In [80]:

```
df4.head()
```

Out[80]:

	genre	title	year_added
0	Documentaries	Dick Johnson Is Dead	2021.0
1	International TV Shows	Blood & Water	2021.0
1	TV Dramas	Blood & Water	2021.0
1	TV Mysteries	Blood & Water	2021.0
2	Crime TV Shows	Ganglands	2021.0

In [81]:

```
grp3 = df4.groupby(by='genre')['year_added'].value_counts().sort_values(ascending = False)
grp3
```

Out[81]:

genre	year_added	
International Movies	2018.0	668
	2019.0	610
	2020.0	575
Dramas	2019.0	564
	2020.0	535
		...
Classic & Cult TV	2016.0	1
	2015.0	1
	2014.0	1
Children & Family Movies	2012.0	1
Thrillers	2011.0	1

Name: year_added, Length: 331, dtype: int64

In [82]:

```
grp3 = pd.DataFrame(grp3)
grp3.head()
```

Out[82]:

		year_added
genre	year_added	
International Movies	2018.0	668
	2019.0	610
	2020.0	575
Dramas	2019.0	564
	2020.0	535

In []:

In [83]:

```
grp3.columns = ['No. of Movies/Shows']
grp3.reset_index(1)
```

Out[83]:

	year_added	No. of Movies/Shows
genre		
International Movies	2018.0	668
International Movies	2019.0	610
International Movies	2020.0	575
Dramas	2019.0	564
Dramas	2020.0	535
...
Classic & Cult TV	2016.0	1
Classic & Cult TV	2015.0	1
Classic & Cult TV	2014.0	1
Children & Family Movies	2012.0	1
Thrillers	2011.0	1

331 rows × 2 columns

In [84]:

```
top_genre
```

Out[84]:

	genre	count
0	International Movies	2752
1	Dramas	2427
2	Comedies	1674
3	International TV Shows	1351
4	Documentaries	869
5	Action & Adventure	859
6	TV Dramas	763
7	Independent Movies	756
8	Children & Family Movies	641
9	Romantic Movies	616

In [85]:

```
per_year2 = grp3.loc[top_genre['genre']]
per_year2
```

Out[85]:

No. of Movies/Shows		
genre	year_added	
International Movies	2018.0	668
	2019.0	610
	2020.0	575
	2021.0	408
	2017.0	395
...
Romantic Movies	2021.0	114
	2018.0	108
	2017.0	63
	2016.0	7
	2015.0	1

91 rows × 1 columns

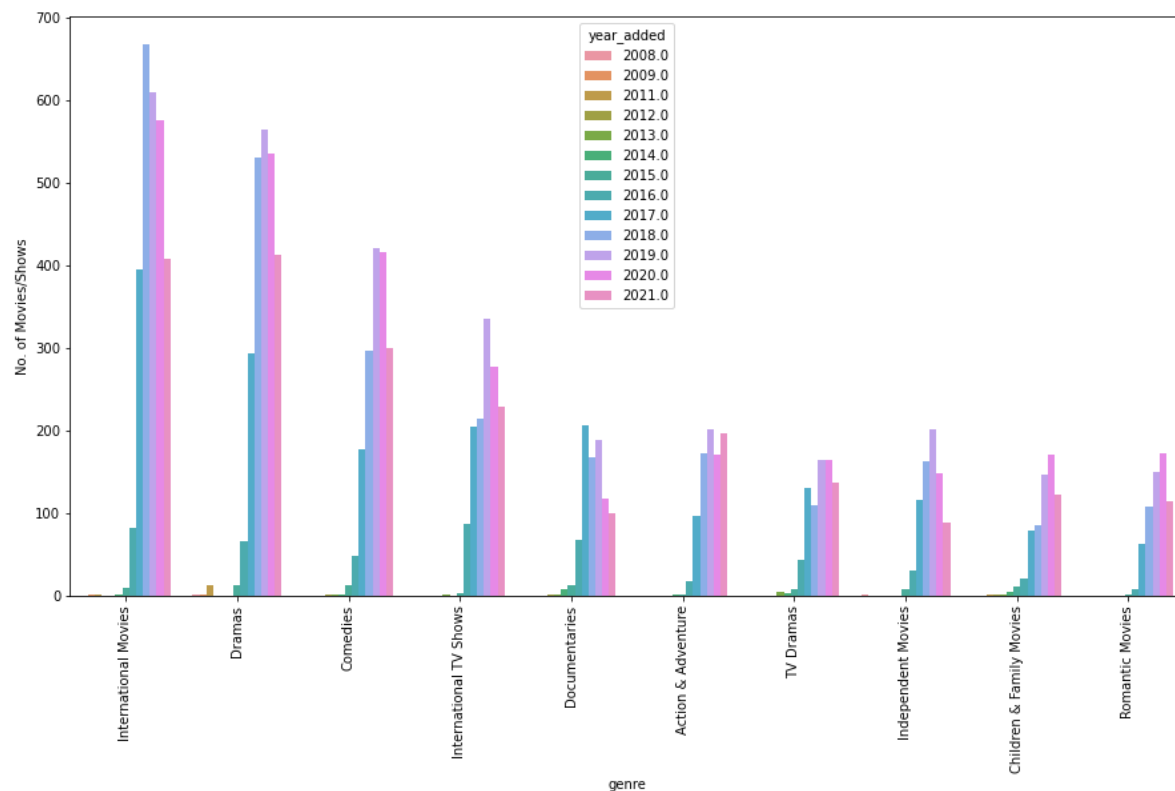
In []:

In [86]:

```
per_year2.reset_index(inplace=True)
```

In [87]:

```
plt.figure(figsize=(15,8))
sns.barplot(data=per_year2, x='genre' , y='No. of Movies/Shows' , hue='year_added')
plt.xticks(rotation=90)
plt.show()
```



In []:

In []:

In []:

Country

In [88]:

```
df.head()
```

Out[88]:

	show_id	type	title	director	cast	country	date_added	release_year	rating	
0	s1	Movie	Dick Johnson Is Dead	Kirsten Johnson	NaN	United States	September 25, 2021	2020	PG-13	
1	s2	TV Show	Blood & Water	NaN	Ama Qamata, Khosi Ngema, Gail Mabalane, Thaban...	South Africa	September 24, 2021	2021	TV-MA	
2	s3	TV Show	Ganglands	Julien Leclercq	Sami Bouajila, Tracy Gotoas, Samuel Jouy, Nabi...	NaN	September 24, 2021	2021	TV-MA	
3	s4	TV Show	Jailbirds New Orleans	NaN	NaN	NaN	September 24, 2021	2021	TV-MA	
4	s5	TV Show	Kota Factory	NaN	Mayur More, Jitendra Kumar, Ranjan Raj, Alam K...	India	September 24, 2021	2021	TV-MA	

In [89]:

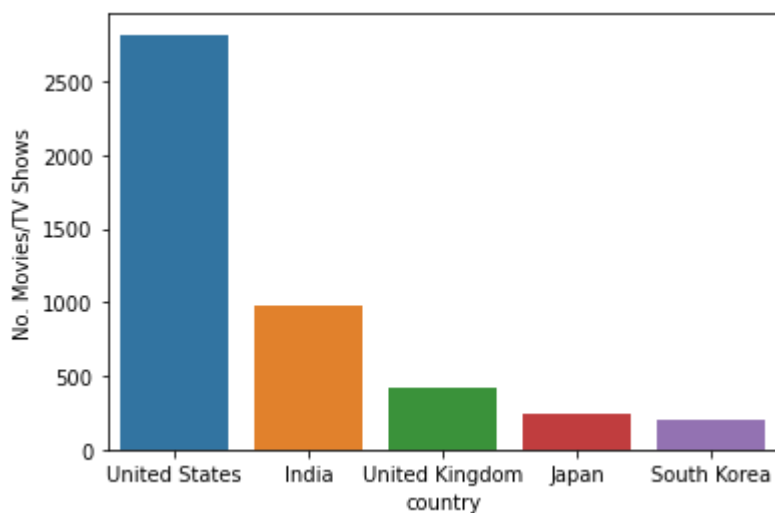
```
# Country with most content on Netflix
country = pd.DataFrame(df['country'].value_counts()[:5])
country.reset_index(inplace = True)
country.columns = ['country', 'No. Movies/TV Shows']
country
```

Out[89]:

	country	No. Movies/TV Shows
0	United States	2818
1	India	972
2	United Kingdom	419
3	Japan	245
4	South Korea	199

In [90]:

```
sns.barplot(data=country, x='country', y = 'No. Movies/TV Shows')
plt.show()
```



In []:

In []:

In []:

In []:

In []: