```python
1 import pandas as pd
2 import matplotlib.pyplot as plt
```

```python
1 df=pd.read_csv('Amazon Prime RT.csv')
```

```python
1 df.columns
```

```
Index(['show_id', 'type', 'title', 'director', 'cast', 'country', 'date_added',
       'release_year', 'rating', 'duration', 'listed_in', 'description'],
      dtype='object')
```

```python
1 df.head(10)
```

| | show_id | type | title | director | cast | country | date_added | release_year | rating | duration | listed_in | description |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | s1 | Movie | The Grand Seduction | Don McKellar | Brendan Gleeson, Taylor Kitsch, Gordon Pinsent | Canada | March 30, 2021 | 2014 | NaN | 113 min | Comedy, Drama | A small fishing village must procure a local d... |
| 1 | s2 | Movie | Take Care Good Night | Girish Joshi | Mahesh Manjrekar, Abhay Mahajan, Sachin Khedekar | India | March 30, 2021 | 2018 | 13+ | 110 min | Drama, International | A Metro Family decides to fight a Cyber Crimin... |
| 2 | s3 | Movie | Secrets of Deception | Josh Webber | Tom Sizemore, Lorenzo Lamas, Robert LaSardo, R... | United States | March 30, 2021 | 2017 | NaN | 74 min | Action, Drama, Suspense | After a man discovers his wife is cheating on ... |
| 3 | s4 | Movie | Pink: Staying True | Sonia Anderson | Interviews with: Pink, Adele, Beyoncé, Britney... | United States | March 30, 2021 | 2014 | NaN | 69 min | Documentary | Pink breaks the mold once again, bringing her ... |
| | | | | | Harry | | | | | | | |

```python
1 df.shape
```

```
(9668, 12)
```

```python
1 df.isnull().any()
```

| | 0 |
|---|---|
| show_id | False |
| type | False |
| title | False |
| director | True |
| cast | True |
| country | True |
| date_added | True |
| release_year | False |
| rating | True |
| duration | False |
| listed_in | False |
| description | False |

dtype: bool

```python
1 df.isnull().sum()
```

⇄

| | 0 |
|---|---|
| **show_id** | 0 |
| **type** | 0 |
| **title** | 0 |
| **director** | 2083 |
| **cast** | 1233 |
| **country** | 8996 |
| **date_added** | 9513 |
| **release_year** | 0 |
| **rating** | 337 |
| **duration** | 0 |
| **listed_in** | 0 |
| **description** | 0 |

dtype: int64

```
1 df.drop_duplicates
```

⇄

**pandas.core.frame.DataFrame.drop_duplicates**
```
def drop_duplicates(subset: Hashable | Sequence[Hashable] | None=None, *, keep: DropKeep='first',
inplace: bool=False, ignore_index: bool=False) -> DataFrame | None
```

Return DataFrame with duplicate rows removed.

Considering certain columns is optional. Indexes, including time indexes
are ignored.

Parameters

```
1 df['director'].fillna('Unknown',inplace=True)
2 df['cast'].fillna('Unkown',inplace=True)
3 df['country'].fillna('Unkown',inplace=True)
4 df['date_added'].fillna('Unkown',inplace=True)
```

⇄   `<ipython-input-10-d7aebea3dda6>:1: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained ass`
  The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting

  For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col]

    df['director'].fillna('Unknown',inplace=True)
  `<ipython-input-10-d7aebea3dda6>:2: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained ass`
  The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting

  For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col]

    df['cast'].fillna('Unkown',inplace=True)
  `<ipython-input-10-d7aebea3dda6>:3: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained ass`
  The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting

  For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col]

    df['country'].fillna('Unkown',inplace=True)
  `<ipython-input-10-d7aebea3dda6>:4: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained ass`
  The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting

  For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col]

    df['date_added'].fillna('Unkown',inplace=True)

```
1 #use mode for categorical value otherwise use mean
2 mode_value = df['rating'].mode()[0]
3 df['rating'].fillna(mode_value)
```

| | rating |
|---|---|
| 0 | 13+ |
| 1 | 13+ |
| 2 | 13+ |
| 3 | 13+ |
| 4 | 13+ |
| ... | ... |
| 9663 | 7+ |
| 9664 | 13+ |
| 9665 | R |
| 9666 | TV-MA |
| 9667 | R |

9668 rows × 1 columns

dtype: object

```
1 df
```

| | show_id | type | title | director | cast | country | date_added | release_year | rating | duration | listed_in | description |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | s1 | Movie | The Grand Seduction | Don McKellar | Brendan Gleeson, Taylor Kitsch, Gordon Pinsent | Canada | March 30, 2021 | 2014 | NaN | 113 min | Comedy, Drama | A small fishing village must procure a local d... |
| 1 | s2 | Movie | Take Care Good Night | Girish Joshi | Mahesh Manjrekar, Abhay Mahajan, Sachin Khedekar | India | March 30, 2021 | 2018 | 13+ | 110 min | Drama, International | A Metro Family decides to fight a Cyber Crimin... |
| 2 | s3 | Movie | Secrets of Deception | Josh Webber | Tom Sizemore, Lorenzo Lamas, Robert LaSardo, R... | United States | March 30, 2021 | 2017 | NaN | 74 min | Action, Drama, Suspense | After a man discovers his wife is cheating on ... |
| 3 | s4 | Movie | Pink: Staying True | Sonia Anderson | Interviews with: Pink, Adele, Beyoncé, Britney... | United States | March 30, 2021 | 2014 | NaN | 69 min | Documentary | Pink breaks the mold once again, bringing her ... |
| 4 | s5 | Movie | Monster Maker | Giles Foster | Harry Dean Stanton, Kieran O'Brien, George Cos... | United Kingdom | March 30, 2021 | 1989 | NaN | 45 min | Drama, Fantasy | Teenage Matt Banting wants to work with a famo... |

```
1 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 9668 entries, 0 to 9667
Data columns (total 12 columns):
 #   Column        Non-Null Count  Dtype
---  ------        --------------  -----
 0   show_id       9668 non-null   object
 1   type          9668 non-null   object
 2   title         9668 non-null   object
 3   director      9668 non-null   object
 4   cast          9668 non-null   object
 5   country       9668 non-null   object
 6   date_added    9668 non-null   object
 7   release_year  9668 non-null   int64
 8   rating        9331 non-null   object
 9   duration      9668 non-null   object
 10  listed_in     9668 non-null   object
 11  description   9668 non-null   object
dtypes: int64(1), object(11)
memory usage: 906.5+ KB
```

```python
1 total_shows= df['show_id'].count()
2 total_shows #This also the total number of titles
```

```
np.int64(9668)
```

```python
1 show_types=df['type'].value_counts()
2 show_types
```

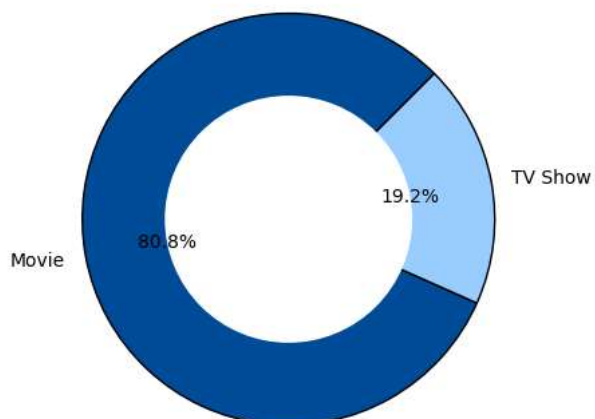|           | count |
|-----------|-------|
| **type**  |       |
| **Movie** | 7814  |
| **TV Show** | 1854 |

dtype: int64

```python
 1 plt.figure(figsize=(5,5))
 2 plt.pie(
 3     show_types,
 4     labels=show_types.index,
 5     colors=['#004C99','#99CCFF'],
 6     autopct='%1.1f%%',
 7     startangle=45,
 8     wedgeprops={'edgecolor': 'black', 'linewidth': 1}
 9 )
10 # Add a white circle at the center to create the donut effect
11 centre_circle = plt.Circle((0, 0), 0.60, fc='white')
12 plt.gca().add_artist(centre_circle)
13 plt.title('Percentage of TV Shows and Movies')
14 plt.show()
```



Percentage of TV Shows and Movies

```
1 #total number of countries
2 total_countries=df['country'].nunique()
3 total_countries
```

87

```
1 df['country'].value_counts().head(5)
```

| country | count |
| --- | --- |
| Unkown | 8996 |
| United States | 253 |
| India | 229 |
| United Kingdom | 28 |
| Canada | 16 |

dtype: int64

```
1 #Top 10 countries except Unknown
2 x=df['country'].value_counts().keys()[1:11]
3 y=df['country'].value_counts()[1:11]
4 #plot
5 plt.barh(x[::-1],y[::-1],
6 color=['#000099'],
7 edgecolor='black')
8 plt.title('Top 10 countries')
9 plt.show
```
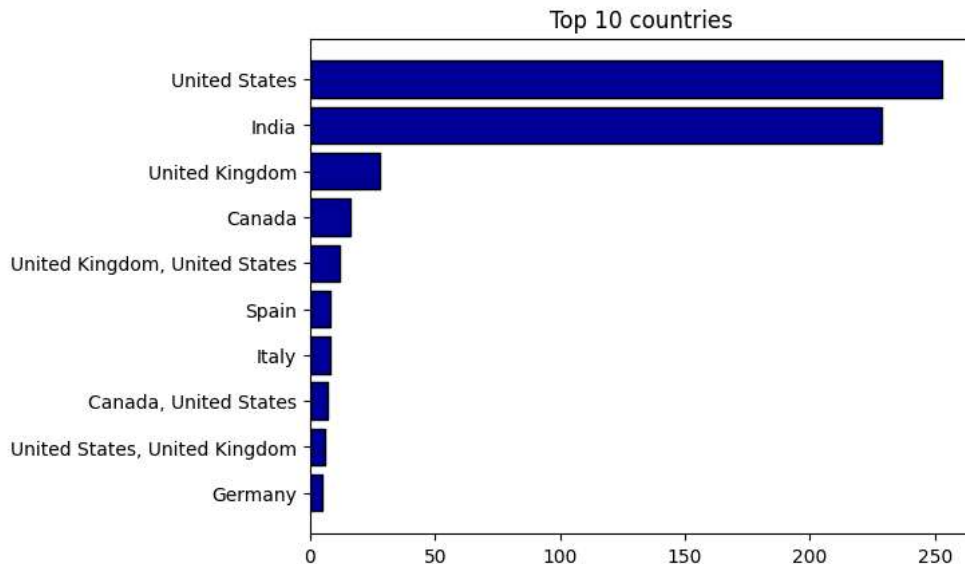
```
matplotlib.pyplot.show
def show(*args, **kwargs) -> None

Display all open figures.

Parameters
----------
block : bool, optional
    Whether to wait for all figures to be closed before returning.
```



Top 10 countries

```
1 #Total types of rating
2 df['rating'].nunique()
```
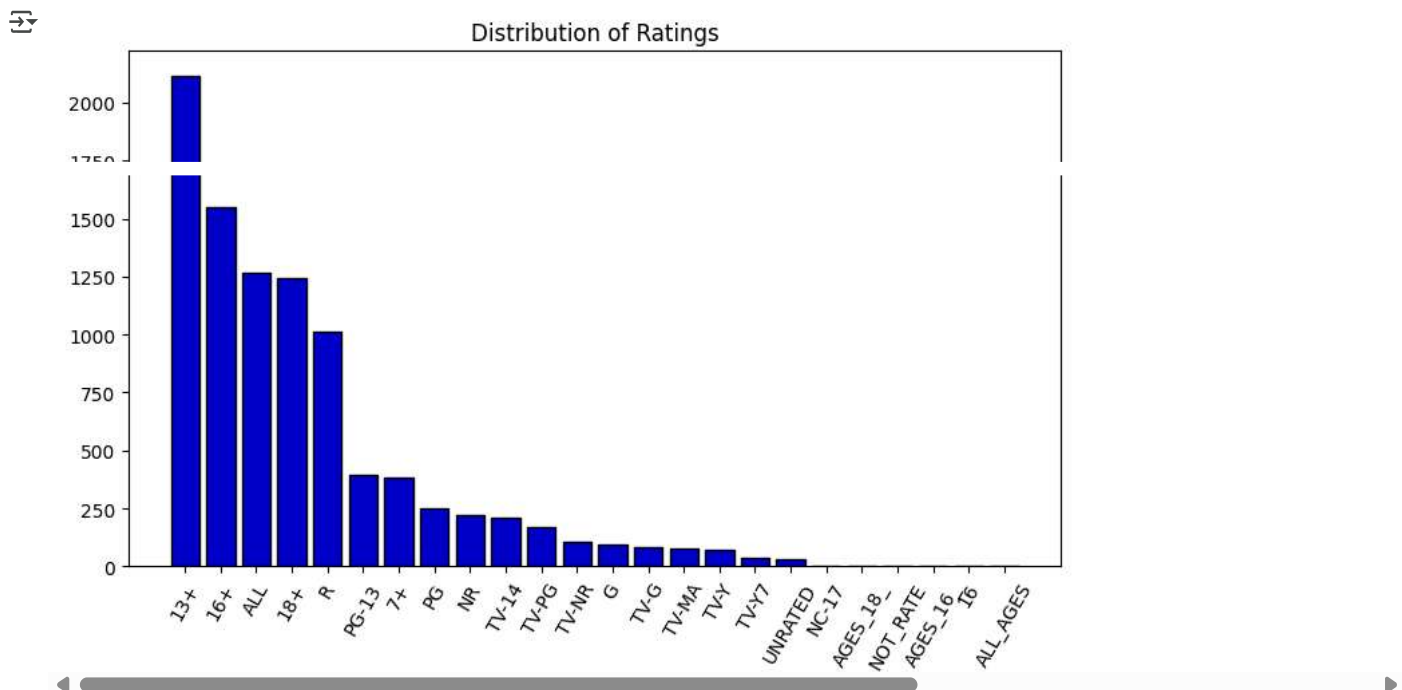
24

```
1 #Bar chart for Distribution of Ratings
2 plt.figure(figsize=(9,5))
3 bars=plt.bar(df['rating'].value_counts().keys(),df['rating'].value_counts(),color='#0000CC',edgecolor
4 plt.title('Distribution of Ratings')
5 plt.xticks(rotation=60)
```

```
 6 #for bar in bars:
 7 #  yval = bar.get_height()
 8 #plt.scatter(bar.get_x() + bar.get_width() / 2, yval, color='black', s=20, zorder=5)
 9
10 plt.show()
```
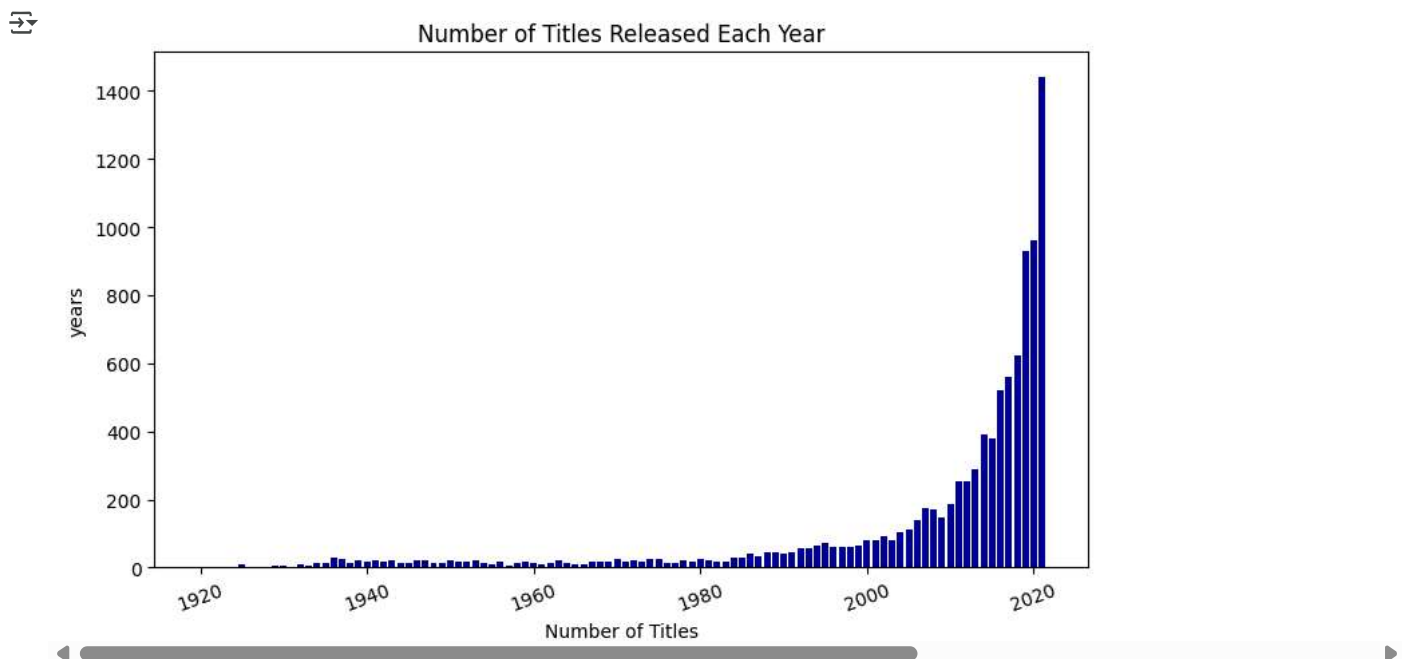


Distribution of Ratings

```
1 #Bar chart of titles released over the years
2 plt.figure(figsize=(9,5))
3 plt.bar(df['release_year'].value_counts().index,df['release_year'].value_counts(),color='#000099')
4 plt.title('Number of Titles Released Each Year')
5 plt.xlabel('Number of Titles')
6 plt.ylabel('years')
7 plt.xticks(rotation=20)
8 plt.show()
```



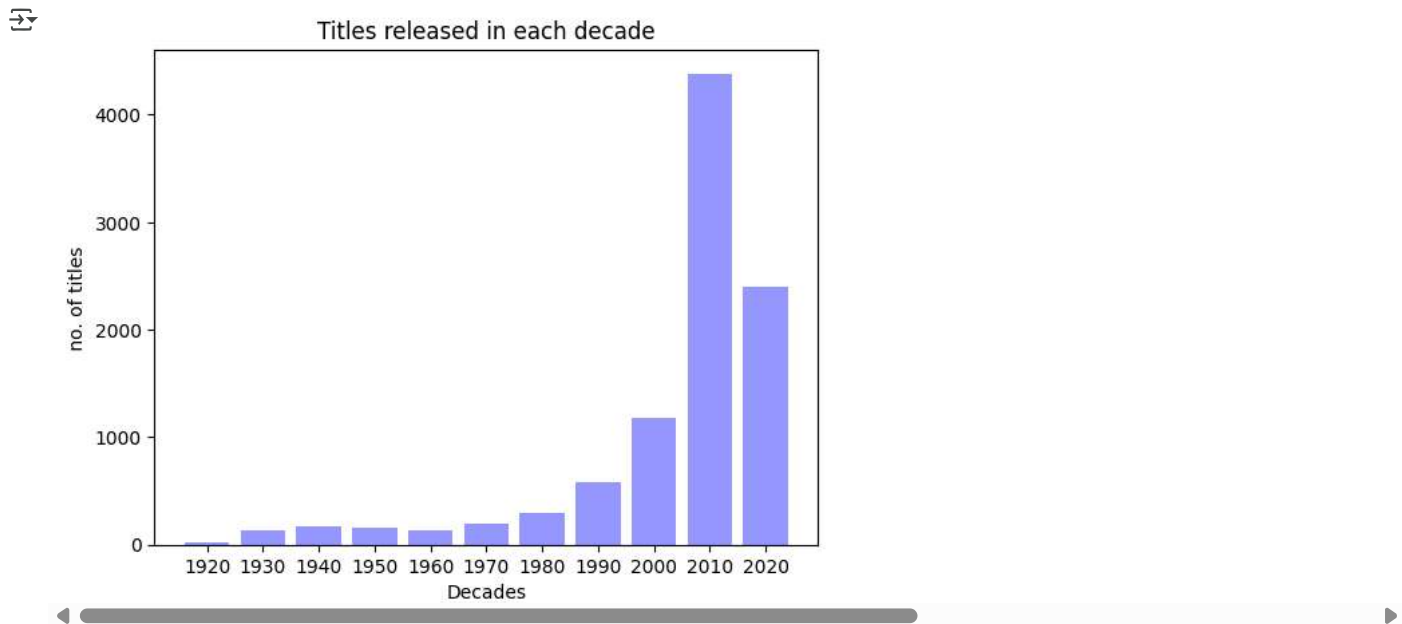Number of Titles Released Each Year

```
1 #Titles released each decade
2 #2005 // 10 computes to 200 and *10 makes it 2000 . Hence grouping it to 2000s
3 #Appends the character 's' after year like 2000s 2010s
4 df['decade']=(df['release_year'] // 10 * 10).astype(str)
5 decade_counts=df['decade'].value_counts().sort_index()
6 plt.bar(decade_counts.index,decade_counts.values,color='#9999FF')
7 plt.xlabel('Decades')
8 plt.ylabel(' no. of titles ')
```

```
 9 plt.title('Titles released in each decade')
10 plt.show()
```

Titles released in each decade



```
1 df
```

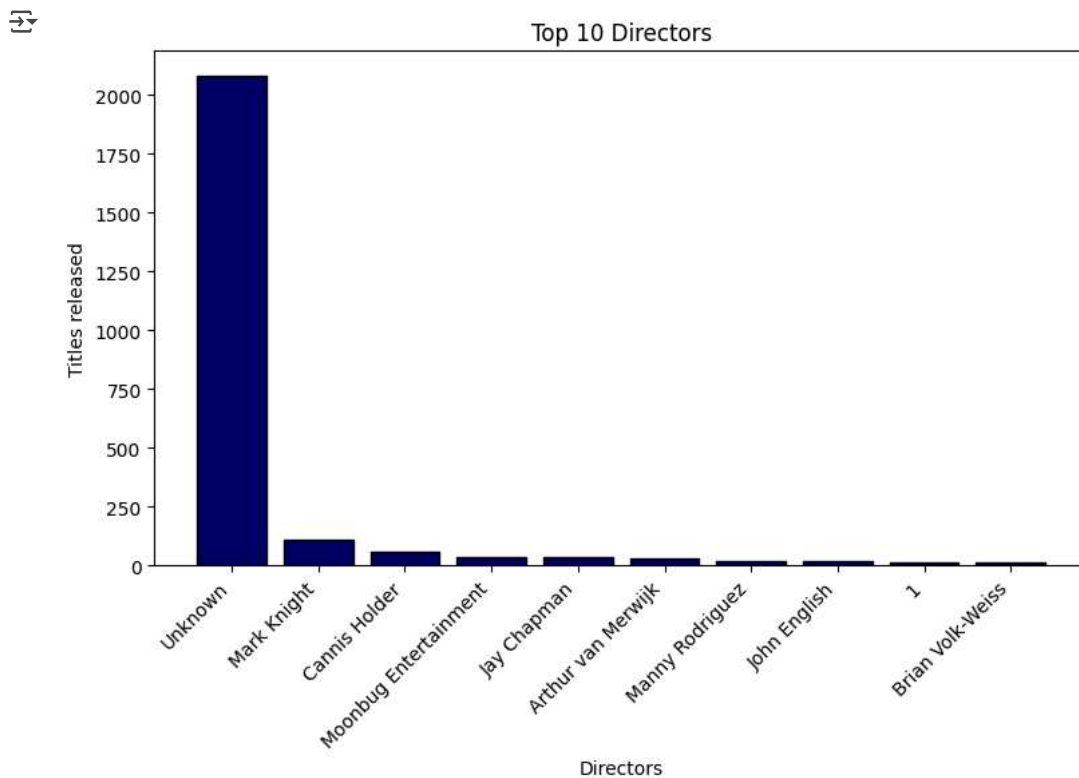| | show_id | type | title | director | cast | country | date_added | release_year | rating | duration | listed_in | description |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | s1 | Movie | The Grand Seduction | Don McKellar | Brendan Gleeson, Taylor Kitsch, Gordon Pinsent | Canada | March 30, 2021 | 2014 | NaN | 113 min | Comedy, Drama | A small fishing village must procure a local d... |
| 1 | s2 | Movie | Take Care Good Night | Girish Joshi | Mahesh Manjrekar, Abhay Mahajan, Sachin Khedekar | India | March 30, 2021 | 2018 | 13+ | 110 min | Drama, International | A Metro Family decides to fight a Cyber Crimin... |
| 2 | s3 | Movie | Secrets of Deception | Josh Webber | Tom Sizemore, Lorenzo Lamas, Robert LaSardo, R... | United States | March 30, 2021 | 2017 | NaN | 74 min | Action, Drama, Suspense | After a man discovers his wife is cheating on ... |
| 3 | s4 | Movie | Pink: Staying True | Sonia Anderson | Interviews with: Pink, Adele, Beyoncé, Britney... | United States | March 30, 2021 | 2014 | NaN | 69 min | Documentary | Pink breaks the mold once again, bringing her ... |
| 4 | s5 | Movie | Monster Maker | Giles Foster | Harry Dean Stanton, Kieran O'Brien, George Cos... | United Kingdom | March 30, 2021 | 1989 | NaN | 45 min | Drama, Fantasy | Teenage Matt Banting wants to work with a famo... |

```
1 df['director'].value_counts().head(10)
```

|  | count |
| --- | --- |
| **director** | |
| **Unknown** | 2083 |
| **Mark Knight** | 113 |
| **Cannis Holder** | 61 |
| **Moonbug Entertainment** | 37 |
| **Jay Chapman** | 34 |
| **Arthur van Merwijk** | 30 |
| **Manny Rodriguez** | 22 |
| **John English** | 20 |
| **1** | 16 |
| **Brian Volk-Weiss** | 15 |

dtype: int64

```
1 #show .nunique() to find no. of unique values instead of unique which shows all unique values
2 df['director'].nunique()
```

5774

```
1 #Top 10 Directors
2 plt.figure(figsize=(9,5))
3 plt.bar(df['director'].value_counts().keys()[:10],df['director'].value_counts()[:10],color='#000066',
4 plt.title('Top 10 Directors')
5 plt.xlabel('Directors')
6 plt.ylabel('Titles released')
7 plt.xticks(rotation=45,ha='right')
8 plt.show()
```



```
1 #Top 10 geners after splitting data considering','
2 #Use lambda to remove spaces before and afer the values
3 #Use .explode() to create
4 #all_genres = df['listed_in'].str.split(',').apply(lambda x:[genre.strip() for genre in x]).explode()
```

```
1 df['listed_in'].nunique()
```

518

```
1 #Start year in the dataset
2 df['release_year'].min()
```

1920

```
1 #End year in the dataset
2 df['release_year'].max()
```
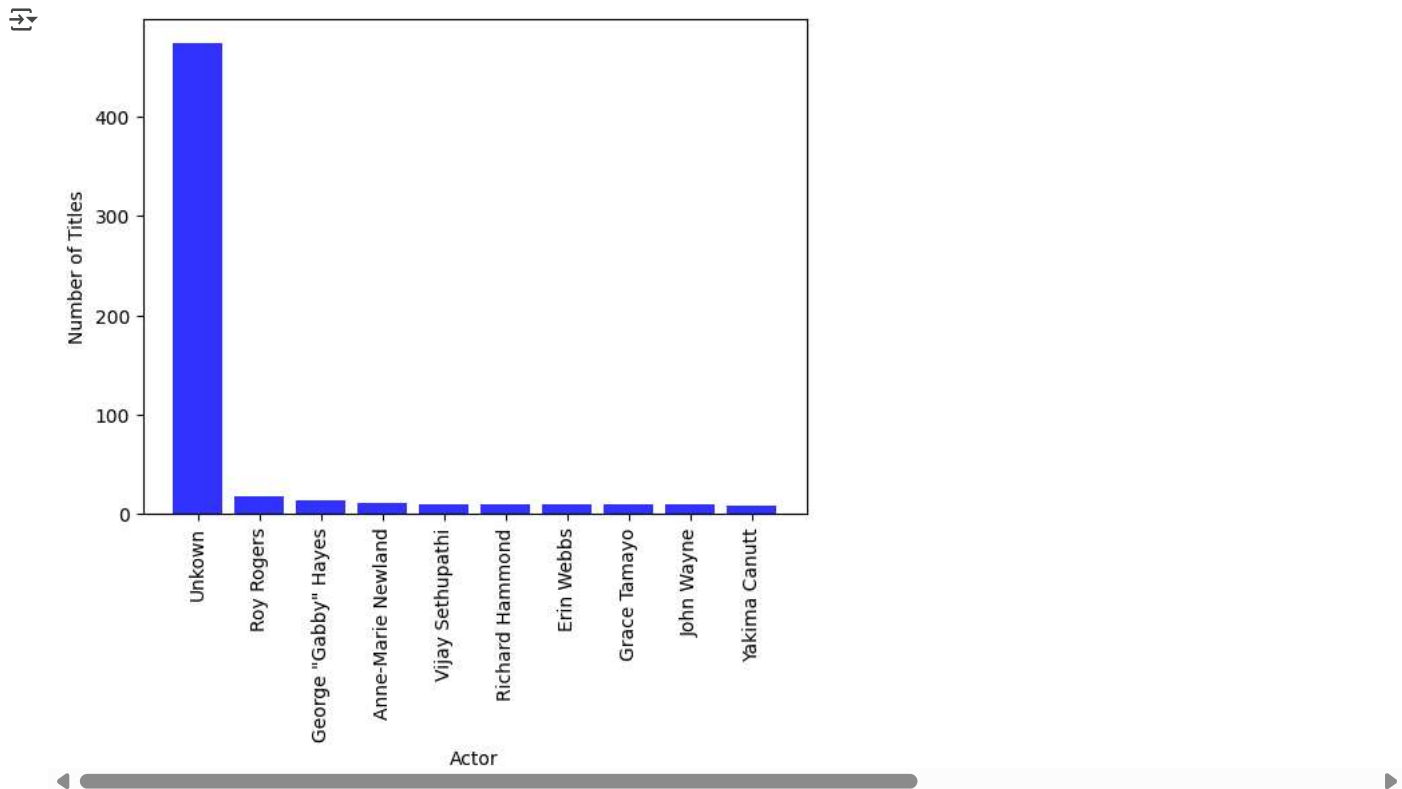
2021

```
1 df['duration'].value_counts(5)
```

|          | proportion |
|----------|------------|
| duration |            |
| 1 Season  | 0.138808   |
| 2 Seasons | 0.023686   |
| 90 min    | 0.023066   |
| 91 min    | 0.022755   |
| 92 min    | 0.020997   |
| ...       | ...        |
| 192 min   | 0.000103   |
| 207 min   | 0.000103   |
| 269 min   | 0.000103   |
| 15 Seasons| 0.000103   |
| 191 min   | 0.000103   |

219 rows × 1 columns

dtype: float64

```
 1 #Reset the index to ensure unique index values. drop=True to discard old index
 2 # ignore_index=True will prevent the error by creating a new index for the exploded column
 3 df = df.reset_index(drop=True)
 4 df['actors']=df['cast'].str.split(',').apply(lambda x: [actor.strip() for actor in x]).explode(ignore
 5
 6
 7 plt.bar(df['actors'].value_counts().keys()[:10],df['actors'].value_counts()[:10],color='#3333FF')
 8 plt.ylabel('Number of Titles')
 9 plt.xlabel('Actor')
10 plt.xticks(rotation=90)
11 plt.show()
```
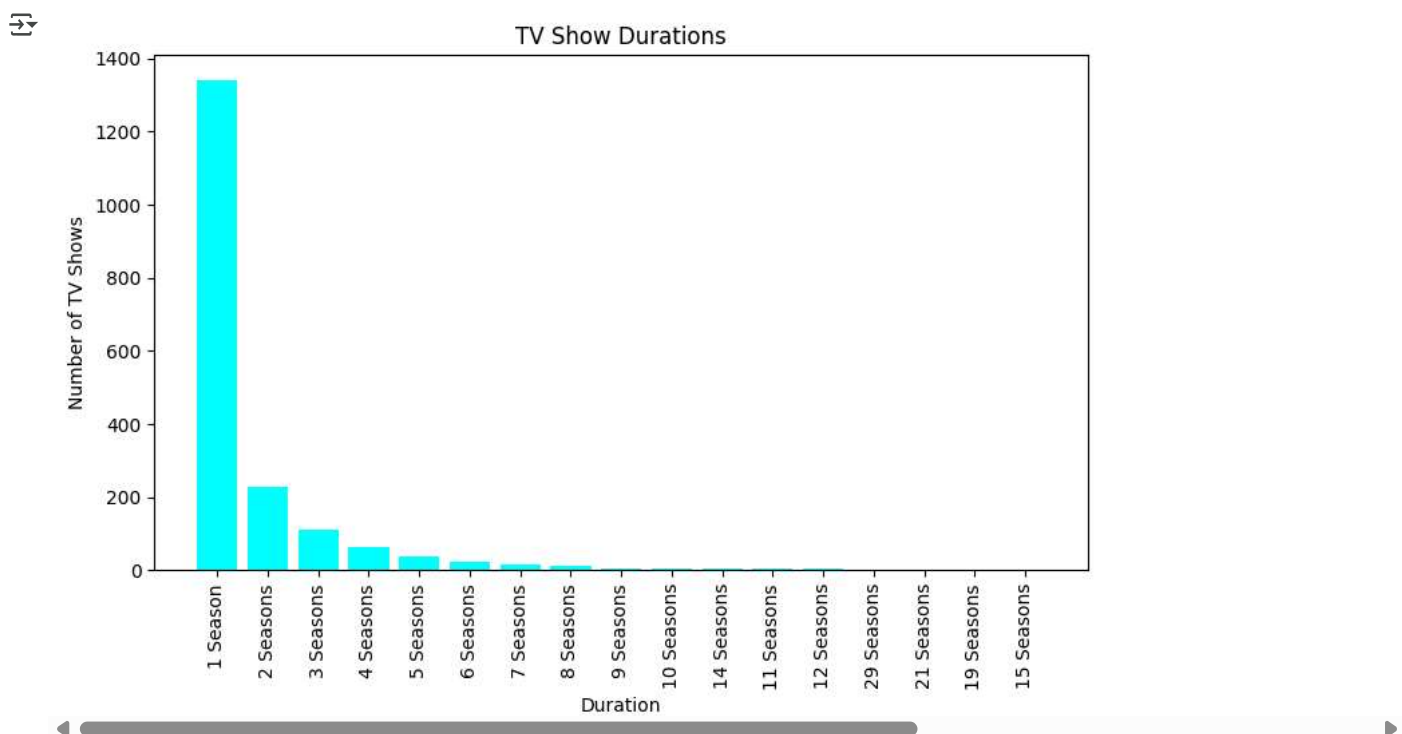
```
1 #TV shows duration
2 # Filter for TV Shows and remove rows with NaN values in the 'duration' column
3 tv_shows = df[df['type'] == 'TV Show'].dropna(subset=['duration'])
4 duration_counts = tv_shows['duration'].value_counts()
```

```
1 # Plot the distribution of TV Show durations
2 plt.figure(figsize=(9, 5))
3 plt.bar(duration_counts.index, duration_counts.values, color='#00FFFF')
4 plt.title('TV Show Durations')
5 plt.xlabel('Duration')
6 plt.ylabel('Number of TV Shows')
7 plt.xticks(rotation=90)
8 plt.show()
9
```

```
1 #Filter for movies and remove the NaN values in the 'Duration' column
2 movie= df[df['type']=='Movie'].dropna(subset=['duration'])
3
4 # Extract numeric duration in minutes from the 'duration' column
5 movie['duration_minutes']=movie['duration'].str.extract('(\d+)').astype(int)
6
7 # Define bins and labels for movie durations
8 bins = [0, 60, 90, 120, 150, 180, 300]
9 labels = ['0-60 min', '61-90 min', '91-120 min', '121-150 min', '151-180 min', '180+ min']
10
11 # Categorize movies into duration groups
12 movie['duration_group'] = pd.cut(movie['duration_minutes'], bins=bins, labels=labels)
13
14 duration_group_counts = movie['duration_group'].value_counts().sort_index()
```

```
1 plt.bar(duration_group_counts.index, duration_group_counts.values, color='blue')
2 plt.title('Movie Durations')
3 plt.xlabel('Duration Group')
4 plt.ylabel('Number of Movies')
5 plt.xticks(rotation=90)
6 plt.show()
```