

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

import plotly.graph_objects as go

from termcolor import colored
from plotly.offline import init_notebook_mode, ipio

import warnings
warnings.filterwarnings("ignore")

In [81]: # read dataset in jupyter notebook
dataset1 = pd.read_csv("dataset1.csv")

In [82]: dataset1["Netflix Release Date"] = pd.to_datetime(dataset1["Netflix Release Date"])

In [83]: #
dataset1.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 2212 entries, 0 to 2212
Data columns (total 37 columns):
 #   Column                                Non-Null Count  Dtype
---  -
 0   title                                 2212 non-null  object
 1   genre                                 2212 non-null  object
 2   languages                            2212 non-null  object
 3   series or movie                      2212 non-null  object
 4   country availability                 2212 non-null  object
 5   runtime                             2212 non-null  object
 6   director                            2212 non-null  object
 7   writer                              2212 non-null  object
 8   IMDb votes                          2212 non-null  float64
 9   View Rating                         2212 non-null  float64
10   IMDb Score                          2212 non-null  float64
11   Rotten Tomatoes Score               2212 non-null  float64
12   Metacritic Score                   2212 non-null  float64
13   Awards Received                     2212 non-null  object
14   Awards Nominated For               2212 non-null  object
15   Bouffice                            2212 non-null  object
16   Release Date                       2212 non-null  object
17   Netflix Release Date                2212 non-null  object
18   Production House                   2212 non-null  object
19   Netflix Release Year                2212 non-null  object
20   Comedy                             1471 non-null  object
21   Crime                              1811 non-null  object
22   Drama                              1992 non-null  object
23   Crime, Adventure, Fantasy, Horror, Romance 1784 non-null  object
24   Crime, Drama, Thriller              1534 non-null  object
25   Horror                              1572 non-null  object
26   Romance                             1764 non-null  object
27   Thriller                            1534 non-null  object
28   Action                              1572 non-null  object
29   Adventure                           1786 non-null  object
30   Sci-Fi                              1886 non-null  object
31   Mystery                             1502 non-null  object
32   Biography                           1591 non-null  object
33   History                             2181 non-null  object
34   Movie                               47 non-null   object
35   Series                              8 non-null   float64
36   director_name                       2212 non-null  object
dtypes: datetime64[ns](3), float64(2), int64(6), object(28)
memory usage: 641.9+ KB

In [84]: dataset1["Netflix Release Date"].dt.year

Out[84]:
0    2021
1    2021
2    2021
3    2021
4    2021
...
2212 2015
2213 2015
2214 2015
2215 2015
2216 2015
Name: Netflix Release Date, Length: 2212, dtype: int64
dataset1["Netflix Release Year"] = dataset1["Netflix Release Date"].dt.year

In [85]: dataset1.isnull().sum()

In [86]:
dataset1

In [87]:
dataset1.duplicated().sum()

In [88]: # to count how many type of shows
dataset1["Series or Movie"].value_counts()

Out[88]:
Movie    2179
Series    47
Name: Series or Movie, dtype: int64

In [89]: #count listed in shows
dataset1["Genre"].value_counts()

Out[89]:
Drama    99
Comedy, Drama, Romance    83
Drama, Romance    78
Comedy    65
...
Adventure, Drama, Western    1
Action, Adventure, Biography, Drama, History    1
Animation, Adventure, Comedy, Family, Music, Musical, Romance    1
Adventure, Family, Fantasy, Musical    1
Biography, Crime, Drama, Thriller, War    1
Name: Genre, Length: 365, dtype: object

In [90]: #import libs for data visualisation
import seaborn as sns
import matplotlib
import matplotlib.pyplot as plt
import matplotlib
import matplotlib

sns.set_style('darkgrid')
matplotlib.rcParams['font.size'] = 14
matplotlib.rcParams['figure.figsize'] = (9, 5)
matplotlib.rcParams['figure.facecolor'] = "#f8cbad"

In [91]: #firstly the difference in types of shows available was analysed
sns.countplot(x="Series or Movie", data=dataset1)
plt.title("Type of the Show on Netflix")

In [92]: #we will now look at what Series of movie account for the most on Netflix using a bar graph.
plt.figure(figsize=(15,6))
dataset1[dataset1["Series or Movie"]=="Movie"]["genre"].value_counts()[:10].plot(kind="barh",color="red")
plt.title("Top 10 Genres of Movies",size=18)

In [93]: #we will now look at what Genres of movie account for the most on Netflix using a bar graph.
plt.figure(figsize=(15,6))
dataset1[dataset1["Series or Movie"]=="Series"]["genre"].value_counts()[:10].plot(kind="barh",color="black")
plt.title("Top 10 Genres of Series",size=18)

In [94]: #types of shows by ratings
plt.figure(figsize=(15,18))
sns.countplot(x="View Rating", data=dataset1)
plt.title("View Rating")

In [95]: #
plt.figure(figsize=(15,18))
sns.countplot(x="Series or Movie", hue="View Rating", data=dataset1)
plt.title("Movies vs Series Rating")

In [96]:
dataset1["Comedy"] = dataset1["Genre"].str.contains("Comedy")
dataset1["Crime"] = dataset1["Genre"].str.contains("Crime", False, NaN=True)
dataset1["Crime"] = dataset1["Genre"].str.contains("Crime", False, NaN=True)
dataset1["Crime"] = dataset1["Genre"].str.contains("Crime", False, NaN=True)
dataset1["Drama"] = dataset1["Genre"].str.contains("Drama")
dataset1["Drama"] = dataset1["Genre"].str.contains("Drama")
dataset1["Drama"] = dataset1["Genre"].str.contains("Drama")
dataset1["Fantasy"] = dataset1["Genre"].str.contains("Fantasy")
dataset1["Fantasy"] = dataset1["Genre"].str.contains("Fantasy")
dataset1["Fantasy"] = dataset1["Genre"].str.contains("Fantasy")
dataset1["Horror"] = dataset1["Genre"].str.contains("Horror")
dataset1["Horror"] = dataset1["Genre"].str.contains("Horror")
dataset1["Horror"] = dataset1["Genre"].str.contains("Horror")
dataset1["Romance"] = dataset1["Genre"].str.contains("Romance")
dataset1["Romance"] = dataset1["Genre"].str.contains("Romance")
dataset1["Romance"] = dataset1["Genre"].str.contains("Romance")
dataset1["Thriller"] = dataset1["Genre"].str.contains("Thriller")
dataset1["Thriller"] = dataset1["Genre"].str.contains("Thriller")
dataset1["Thriller"] = dataset1["Genre"].str.contains("Thriller")
dataset1["Action"] = dataset1["Genre"].str.contains("Action", False, NaN=True)
dataset1["Action"] = dataset1["Genre"].str.contains("Action", False, NaN=True)
dataset1["Action"] = dataset1["Genre"].str.contains("Action", False, NaN=True)
dataset1["Adventure"] = dataset1["Genre"].str.contains("Adventure", False, NaN=True)
dataset1["Adventure"] = dataset1["Genre"].str.contains("Adventure", False, NaN=True)
dataset1["Adventure"] = dataset1["Genre"].str.contains("Adventure", False, NaN=True)
dataset1["Fantasy"] = dataset1["Genre"].str.contains("Fantasy")
dataset1["Fantasy"] = dataset1["Genre"].str.contains("Fantasy")
dataset1["Fantasy"] = dataset1["Genre"].str.contains("Fantasy")
dataset1["Sci-Fi"] = dataset1["Genre"].str.contains("Sci-Fi", False, NaN=True)
dataset1["Sci-Fi"] = dataset1["Genre"].str.contains("Sci-Fi", False, NaN=True)
dataset1["Sci-Fi"] = dataset1["Genre"].str.contains("Sci-Fi", False, NaN=True)
dataset1["Mystery"] = dataset1["Genre"].str.contains("Mystery")
dataset1["Mystery"] = dataset1["Genre"].str.contains("Mystery")
dataset1["Mystery"] = dataset1["Genre"].str.contains("Mystery")
dataset1["Biography"] = dataset1["Genre"].str.contains("Biography")
dataset1["Biography"] = dataset1["Genre"].str.contains("Biography")
dataset1["Biography"] = dataset1["Genre"].str.contains("Biography")
dataset1["History"] = dataset1["Genre"].str.contains("History")
dataset1["History"] = dataset1["Genre"].str.contains("History")
dataset1["History"] = dataset1["Genre"].str.contains("History")

In [97]:
plt.figure(figsize=(10,5))
sns.countplot(x="Series or Movie", data=dataset1)
plt.title("Series or Movie")
plt.show()

In [98]:
from wordcloud import WordCloud, ImageColorGenerator

In [99]:
text = " ".join(str(each) for each in dataset1.Title)
# Create and generate a word cloud image:
wordcloud = WordCloud(max_words=200, background_color="gray").generate(text)
plt.figure(figsize=(10,6))
plt.figure(figsize=(10,10))
# Display the generated image:
plt.imshow(wordcloud, interpolation='bilinear')
plt.title("Most Popular Title", fontsize = 30)
plt.axis('off')
plt.show()

<Figure size 720x432 with 0 Axes>

Most Popular Title

In [100]:
order = ['G','PG','PG-13','Not Rated','Unrated','PG-13','PG-13','TV-14','R','NC-17','TV-MA']
dataset1["View Rating"] = dataset1["View Rating"].replace("G", "G").replace("PG", "PG").replace("PG-13", "PG-13").replace("Not Rated", "Not Rated").replace("Unrated", "Unrated").replace("PG-13", "PG-13").replace("TV-14", "TV-14").replace("R", "R").replace("NC-17", "NC-17").replace("TV-MA", "TV-MA")
dataset1["Total Count"] = dataset1["View Rating"].value_counts().reset_index().sort_values("Total Count", ascending=False)
plt.figure(figsize=(10,5))
plt.title("Total Count")
plt.show()

Ratings for Movies & Series

In [101]:
fig, ax = plt.subplots(1,2,figsize=(10,5))
g1 = sns.countplot(dataset1["View Rating"], order=order, palette="Set2", ax=ax[0])
g1.set_title("Ratings for Movies")
g1.set_xlabel("Rating")
g2 = sns.countplot(dataset1["View Rating"], order=order, palette="Set2", ax=ax[1])
g2.set_title("Ratings for Series")
g2.set_xlabel("Rating")
g2.set_ylabel("Total Count")
g2.set_ylabel("Total Count")
fig.show()

Ratings for Movies

Ratings for Series

In [102]:
dataset1["Director_name"] = dataset1["Director"].apply(lambda x: x.replace(",","").replace(" ","").split(","))
Director_count = []
for i in dataset1["Director_name"]:
    Director_count.append(Director_count.count(i))
Director_dict = dict(zip(Director_count, Director_count))
dataset1["Director_name"] = dataset1["Director_name"].apply(lambda x: x.replace(",","").replace(" ","").split(","))

In [103]:
plt.figure(figsize=(15,5))
sns.countplot(x="Director", data=dataset1)
plt.title("Top20 Director on Netflix", fontweight="bold")
plt.show()

Top20 Director on Netflix

In [104]:
dataset1["Series"] = dataset1[dataset1["Series or Movie"] == "Series"]
dataset1["Movie"] = dataset1[dataset1["Series or Movie"] == "Movie"]

In [105]:
dataset1["content"] = dataset1["Netflix Release Year"].value_counts().reset_index().rename(columns={"Netflix Release Year": "count", "index": "Netflix Release Year"})
dataset1["content"] = dataset1["content"].apply(lambda x: 188*x/sum(dataset1["content"])))
dataset1["Series"] = dataset1["Series"].value_counts().reset_index().rename(columns={"Series": "count", "index": "Series"})
dataset1["Series"] = dataset1["Series"].apply(lambda x: 188*x/sum(dataset1["Series"])))
dataset1["Movie"] = dataset1["Movie"].value_counts().reset_index().rename(columns={"Movie": "count", "index": "Movie"})
dataset1["Movie"] = dataset1["Movie"].apply(lambda x: 188*x/sum(dataset1["Movie"])))
t1 = go.Scatter(x=dataset1["Series"], y=dataset1["Series"], name="Series", marker=dict(color="red", size=10))
t2 = go.Scatter(x=dataset1["Movie"], y=dataset1["Movie"], name="Movie", marker=dict(color="blue", size=10))
t3 = go.Scatter(x=dataset1["content"], y=dataset1["content"], name="Total Content", marker=dict(color="green", size=10))
data = [t1, t2, t3]
layout = go.Layout(title="Content added over the years", legend=dict(x=0.1, y=1.1, orientation="h"))
fig = go.Figure(data, layout)
fig.show()

In [106]:
import plotly.graph_objects as go
from plotly.offline import init_notebook_mode, ipio

In [107]:
filtered_countries = dataset1["Country Availability"].value_counts().reset_index().drop(True)
filtered_countries = filtered_countries[filtered_countries["Country Availability"] != "Unavailable"]
plt.plot(filtered_countries["Country Availability"], filtered_countries["Country Availability"])

In [108]:
plt.figure(figsize=(15,5))
sns.countplot(x="Runtime", data=dataset1)
plt.title("Based on Runtime", fontweight="bold")
plt.show()

Based on Runtime
```