

```
In [31]: # IMPORTING VARIOUS LIBRARIES
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

In [61]: # UPLOADING DATA AND SEE TOP 5 ROWS
df= pd.read_csv("netflix.csv")
df.head(5)

Out [61]:
show_id type title director cast country date_added release_year rating duration listed_in description
0 s1 Movie Dick Johnson Is Dead Kirsten Johnson NaN United States September 25, 2021 2020 PG-13 90 min Documentaries As her father nears the end of his life, film...
1 s2 TV Show Blood & Water NaN Ama Qamata, Khosi Ngema, Gail Mablane, Thabani... South Africa September 24, 2021 2021 TV-MA 2 Seasons International TV Shows, TV Dramas, TV Mysteries After crossing paths at a party, a Cape Town L...
2 s3 TV Show Ganglands Julien Leclercq Sami Bouajila, Tracy Gotosas, Samuel Jouy, Nabil... NaN NaN September 24, 2021 2021 TV-MA 1 Season Crime TV Shows, International TV Shows, TV Act... To protect his family from a powerful drug lo...
3 s4 TV Show Jailbirds New Orleans NaN NaN NaN NaN September 24, 2021 2021 TV-MA 1 Season Docuseries, Reality TV Feuds, flirtations and toilet talk go down amo...
4 s5 TV Show Kota Factory NaN Mayor More, Jitendra Kumar, Ranjan Raj, Alam K... India September 24, 2021 2021 TV-MA 2 Seasons International TV Shows, TV Dramas, TV Mysteries To protect his family from a powerful drug lo...

DATA INFO

In [81]: 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 8787 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 [771]: df.shape # THIS SHOWS DATA CONTAINS 8807 ROWS AND 12 COLUMNS
(8807, 12)

Out [771]:

Analyse for Null values

In [781]: df.isna().sum() # This shows data contains null values at director, cast, country, date_added, rating and duration column

Out [781]:
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
dtypes: int64

As the data shows - director, cast , date_added ,rating and duration column consist of null values. lets check for null values

In [781]: df[df['director']!=isna()]

Out [781]:
show_id type title director cast country date_added release_year rating duration listed_in description
1 s2 TV Show Blood & Water NaN Ama Qamata, Khosi Ngema, Gail Mablane, Thabani... South Africa September 24, 2021 2021 TV-MA 2 Seasons International TV Shows, TV Dramas, TV Mysteries After crossing paths at a party, a Cape Town L...
3 s4 TV Show Jailbirds New Orleans NaN NaN NaN September 24, 2021 2021 TV-MA 1 Season Docuseries, Reality TV Feuds, flirtations and toilet talk go down amo...
4 s5 TV Show Kota Factory NaN Mayor More, Jitendra Kumar, Ranjan Raj, Alam K... India September 24, 2021 2021 TV-MA 2 Seasons International TV Shows, Romantic TV Shows, TV ... In a city of coaching centers known to train l...

2634 rows x 12 columns

In [801]: df['director']= df['director'].fillna('unknown_director') # filling null values at director column with "unknown_director"
df

Out [801]:
show_id type title director cast country date_added release_year rating duration listed_in description
0 s1 Movie Dick Johnson Is Dead Kirsten Johnson NaN United States September 25, 2021 2020 PG-13 90 min Documentaries As her father nears the end of his life, film...
1 s2 TV Show Blood & Water unknown_director Ama Qamata, Khosi Ngema, Gail Mablane, Thabani... NaN South Africa September 24, 2021 2021 TV-MA 2 Seasons International TV Shows, TV Dramas, TV Mysteries After crossing paths at a party, a Cape Town L...
2 s3 TV Show Ganglands Julien Leclercq Sami Bouajila, Tracy Gotosas, Samuel Jouy, Nabil... NaN NaN September 24, 2021 2021 TV-MA 1 Season Crime TV Shows, International TV Shows, TV Act... To protect his family from a powerful drug lo...
3 s4 TV Show Jailbirds New Orleans unknown_director NaN NaN NaN September 24, 2021 2021 TV-MA 1 Season Docuseries, Reality TV Feuds, flirtations and toilet talk go down amo...
4 s5 TV Show Kota Factory unknown_director NaN NaN NaN September 24, 2021 2021 TV-MA 2 Seasons International TV Shows, TV Dramas, TV Mysteries To protect his family from a powerful drug lo...

8807 rows x 12 columns

In [811]: df.isna().sum()# director column now showing 0 null, as we have filled null values with "unknown"

Out [811]:
show_id 0
type 0
title 0
director 0
cast 825
country 831
date_added 10
release_year 0
rating 4
duration 3
listed_in 0
description 0
dtypes: int64

In [821]: for a in ['date_added','rating','duration']:
df[a] = df[a].fillna(df[a].mode()[0]) #filling missing values with mode

In [831]: df['country']=df['country'].fillna('unknown_country')
df['cast']=df['cast'].fillna('unknown_cast')

In [841]: df.isna().sum() #as we have filled all the null values now data shows no null values and it is very much cleaned

Out [841]:
show_id 0
type 0
title 0
director 0
cast 0
country 0
date_added 0
release_year 0
rating 0
duration 0
listed_in 0
description 0
dtypes: int64

Checking for unique values

In [851]: df['director'].nunique() #director column contains 4529 unique director entries

Out [851]:
4529

In [861]: df['title'].nunique() # title column contains 8807 unique entries

Out [861]:
8807

In [871]: df['cast'].nunique() # cast column has 7693 unique entries

Out [871]:
7693

In [881]: df['type'].unique() # type column has two unique entries as Movie and Tv Show

Out [881]:
array(['Movie', 'TV Show'], dtype=object)

Checking for Data types

In [891]: df.dtypes #all the data type seems correct except date_added. it should be timestamp

Out [891]:
show_id object
type object
title object
director object
cast object
country object
date_added object
release_year int64
rating object
duration object
listed_in object
description object
dtypes: object

In [901]: df['date_added'] = df['date_added'].str.strip()
df['date_added'] = pd.to_datetime(df['date_added'], format='%b %d, %Y') #Changing date_added datatype to datetime
df.dtypes

Out [901]:
show_id object
type object
title object
director object
cast object
country object
date_added datetime64[ns]
release_year int64
rating object
duration object
listed_in object
description object
dtypes: object

In [911]: #making of new columns day, month, year for better analyzing
df['date']=df['date_added'].dt.day
df['month']=df['date_added'].dt.month
df['year']=df['date_added'].dt.year
df.dtypes

Out [911]:
show_id object
type object
title object
director object
cast object
country object
date_added datetime64[ns]
release_year int64
rating object
duration object
listed_in object
date int64
month int64
year int64
dtypes: object

In [921]: #checking for value counts
df['director'].value_counts()

Out [921]:
unknown_director 2634
Rajiv Chilaka 19
Ravi Campos, Jan Suter 18
Subas Kaday 16
Marcus Raboy 16
Raymie Muizquiz, Stu Livingston ... 1
Joe Menendez 1
Eric Brosa 1
Will Eisenberg 1
Moxez Singh 1
Name: director, Length: 4529, dtype: int64

In [931]: df['cast'].value_counts()

Out [931]:
unknown_cast 825
David Attenborough 19
Vatsal Dubey, Julie Tejwani, Rupa Bhimani, Zigna Bhardwaj, Rajesh Kava, Mousam, Swapnil 14
Samuel West 10
Jeff Dunham 7
...
Nick Lachey, Vanessa Lachey 1
Takeru Sato, Kasumi Arimura, Haru, Kentaro Sakaguchi, Takayuki Yamada, Kendo Kobayashi, Ken Yasuda, Arata Furuta, Suzuki Matsuo, Koichi Yamadera, Arata Furuta, Chikako Kaku, Kotaro Yoshida 1
Toyin Abraham, Sambasa Nzeribe, Chioma Chukwuka Akpotha, Chioma Omeruah, Chivetalu Agu, Dele Odiele, Femi Adebayo, Bayray McNwizu, Biodun Stephen 1
Heeraj Kahl, Geetanjali Kulkarni, Danish Hussain, Sheeba Chaddha, Paras Priyadarshan, Anshul Chauhan, Anud Singh Dhaka, Shririn Sewani, Mihir Ahuja, Vasundhara Rajput 1
Vicky Kaushal, Sarah-Jane Dias, Raaghav Chhanna, Manish Chaudhary, Meghna Malik, Malket Rauni, Anita Shahidish, Chhittaranjan Tripathy 1
Name: cast, Length: 7693, dtype: int64

In [941]: df['type'].value_counts() #This shows Movie type has 6131 entry whereas Tv Show type has 2676 entry

Out [941]:
Movie 6131
TV Show 2676
Name: type, dtype: int64

In [951]: df['title'].value_counts()

Out [951]:
Dick Johnson Is Dead 1
Ip Man 2
Hannibal Buress: Comedy Camisado 1
Tubbo PAT 1
Masha's Tales 1
Love for Sale 2
ROAD TO ROMA 1
Good Time 1
Captain Underpants Epic Choice-o-Rama 1
Zubean 1
Name: title, Length: 8807, dtype: int64

In [961]: df['country'].value_counts() #To check the value counts of country column

Out [961]:
United States 2818
India 972
unknown_country 831
United Kingdom 419
Japan 245
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: 749, dtype: int64

In [971]: df['listed_in'].value_counts()

Out [971]:
Dramas, International Movies 362
Documentaries 359
Stand-Up Comedy 334
Comedies, Dramas, International Movies 274
Dramas, Independent Movies, International Movies 252
Kids' TV, Animation & Adventure, TV Dramas 1
TV Comedies, TV Action, TV Horror 1
Children & Family Movies, Comedies, LGBTQ Movies 1
Kids' TV, Spanish-Language TV Shows, Teen TV Shows 1
Cult Movies, Dramas, Thrillers 1
Name: listed_in, Length: 514, dtype: int64

In [981]: df['rating'].value_counts() # this column has anonyms value as 74 min, 84 min, 66min

Out [981]:
TV-MA 3211
TV-14 2160
TV-PG 863
R 789
PG-13 490
TV-Y 334
TV-Y7 307
PG 287
TV-G 220
NR 80
G 41
TV-Y7-FV 6
NC-17 3
TV-14 3
74 min 1
84 min 1
66 min 1
Name: rating, dtype: int64

In [991]: df['duration'].value_counts()

Out [991]:
1 Season 1796
2 Seasons 425
3 Seasons 199
90 min 152
94 min 146
...
16 min 1
186 min 1
193 min 1
189 min 1
191 min 1
Name: duration, Length: 220, dtype: int64

In [1001]: #cleaning of country column
df['country']= df['country'].astype(str)
df['country']= df['country'].apply(lambda x: x.split(',')[0])

In [1011]: df['country'].value_counts()

Out [1011]:
United States 3210
India 1008
unknown_country 831
United Kingdom 626
Canada 271
...
Namibia 1
Senegal 1
Cameroon 1
Syria 1
Somalia 1
Name: country, Length: 90, dtype: int64

In [1021]: #cleaning of rating column
df['rating']=df['rating'].unique()

Out [1021]:
array(['PG-13', 'TV-MA', 'PG', 'TV-14', 'TV-PG', 'TV-Y', 'TV-Y7', 'TV-Y7-FV', 'R', 'TV-G', 'G', 'NC-17', '74 min', '84 min', '66 min', 'NR', 'TV-Y7-FV', 'UR'], dtype=object)

In [1101]: df.groupby(['rating', 'type']).size().unstack().fillna(0).nlargest(10, columns=['Movie', 'TV Show'])

Out [1101]:
type Movie TV Show
rating
TV-MA 2064.0 1147.0
TV-14 1427.0 733.0
R 797.0 2.0
TV-PG 540.0 323.0
PG-13 490.0 0.0
PG 287.0 0.0
TV-Y7 129.0 195.0
TV-Y 131.0 176.0
TV-G 126.0 94.0
NR 75.0 5.0

In [1111]: df.groupby(['rating', 'type']).size().unstack().fillna(0).nlargest(10, columns=['Movie', 'TV Show']).plot(kind='bar', stacked=True, colormap='cool')
figsize=(8, 4)
plt.show() # To visualise Top 10 Ratings

Top 10 Ratings: Movies vs. TV Shows

Rating
Count
type
TV-MA 3000
TV-14 2000
R 800
TV-PG 800
PG-13 500
PG 400
TV-Y7 300
TV-Y 200
TV-G 100
NR 50

In [1121]: df['release_year'].value_counts().plot(marker='o', linestyle='-', color='blue', title='Release Year Trends', figsize=(10,4))
plt.xlabel('Release Year')
plt.ylabel('Count')
plt.xticks(rotation=45)
plt.show() # This shows the release year trend and shows after 2000 the most of the content is released

Release Year Trends

Count
Release Year
1940 0
1950 0
1960 0
1970 0
1980 0
1990 0
2000 0
2010 0
2020 1000
2021 1000
2022 1000
2023 1000
2024 1000

In [1131]: sns.scatterplot(data=df, x='type')
plt.title('Distribution of Content Types')
plt.show() # distribution of content shows type Movies are released more than that of Tv Show

Distribution of Content Types

Count
type
Movie 6000
TV Show 2500

In [1132]: # Top 10 countries PRODUCING MOVIE
top_countries = df[df['type'] == 'Movie']
top_10_countries_tv_show = list(df[df['type'] == 'TV Show'].value_counts().index[:10])
print('TOP 10 COUNTRIES:', top_10_countries)
#Below are top 10 countries who have produced more number of movie united states are on the top

TOP 10 COUNTRIES: ('United States', 'India', 'United Kingdom', 'Canada', 'Spain', 'Egypt', 'Nigeria', 'Indonesia', 'Turkey', 'Japan')

In [1133]: df.tv_shows = df[df['type'] == 'TV Show']
top_10_countries_tv_show = list(df.tv_shows[df['country'].value_counts().index[:10])
print('TOP 10 COUNTRIES:', top_10_countries_tv_show)
# Below are list of countries tv show having more number of Tv Show released. United states here also on maintaining the top position

TOP 10 COUNTRIES: ('United States', 'United Kingdom', 'Japan', 'South Korea', 'India', 'Taiwan', 'Canada', 'France', 'Australia', 'Spain')

No of distinct titles on the basis of genre

In [1134]: df.genre = df.groupby(['listed_in']).agg({'rating':'nunique'}).reset_index().sort_values(by=['rating'], ascending = False)[:10]

In [1135]: df.genre #Top 10 genres as per the rating

Out [1135]:
listed_in rating
319 Dramas, Independent Movies, International Movies 10
442 Movies 10
274 Documentaries 9
217 Comedies, International Movies 8
289 Documentaries, Sports Movies 8
125 Children & Family Movies, Comedies 8
230 Comedies, Romantic Movies 8
326 Dramas, International Movies 8
220 Comedies, International Movies, Romantic Movies 8
200 Comedies, Dramas, International Movies 7

In [1136]: plt.figure(figsize=(6,8))
plt.barh(df.genre[:11]['listed_in'],df.genre[:11]['rating'])
plt.xlabel('Genres')
plt.ylabel('Count')
plt.show()

Dramas, Independent Movies, International Movies
Movies
Documentaries
Comedies, International Movies
Documentaries, Sports Movies
Children & Family Movies, Comedies
Comedies, Romantic Movies
Dramas, International Movies
Comedies, International Movies, Romantic Movies
Comedies, Dramas, International Movies
```