



# Netflix\_Case\_Study

## Importing the libraries

```
In [70]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

## Loading the dataset

```
In [71]: !gdown https://d2beiqkhq929f0.cloudfront.net/public_assets/assets/000/000/940/original/netflix.csv
Downloading...
From: https://d2beiqkhq929f0.cloudfront.net/public_assets/assets/000/000/940/original/netflix.csv
To: /content/netflix.csv
100% 3.40M/3.40M [00:01<00:00, 1.95MB/s]
```

```
In [72]: # Importing the dataset
df = pd.read_csv('netflix.csv')
df
```

Out[72]:

	show_id	type	title	director	cast	country	date_added	rel
<b>0</b>	s1	Movie	Dick Johnson Is Dead	Kirsten Johnson	NaN	United States	September 25, 2021	
<b>1</b>	s2	TV Show	Blood & Water	NaN	Ama Qamata, Khosi Ngema, Gail Mabalane, Thaban...	South Africa	September 24, 2021	
<b>2</b>	s3	TV Show	Ganglands	Julien Leclercq	Sami Bouajila, Tracy Gotoas, Samuel Jouy, Nabi...	NaN	September 24, 2021	
<b>3</b>	s4	TV Show	Jailbirds New Orleans	NaN	NaN	NaN	September 24, 2021	
<b>4</b>	s5	TV Show	Kota Factory	NaN	Mayur More, Jitendra Kumar, Ranjan Raj, Alam K...	India	September 24, 2021	
...	...	...	...	...	...	...	...	
<b>8802</b>	s8803	Movie	Zodiac	David Fincher	Mark Ruffalo, Jake Gyllenhaal, Robert Downey J...	United States	November 20, 2019	
<b>8803</b>	s8804	TV Show	Zombie Dumb	NaN	NaN	NaN	July 1, 2019	
<b>8804</b>	s8805	Movie	Zombieland	Ruben Fleischer	Jesse Eisenberg, Woody Harrelson, Emma Stone, ...	United States	November 1, 2019	

	show_id	type	title	director	cast	country	date_added	rel
8805	s8806	Movie	Zoom	Peter Hewitt	Tim Allen, Courteney Cox, Chevy Chase, Kate Ma...	United States	January 11, 2020	
8806	s8807	Movie	Zubaan	Mozez Singh	Vicky Kaushal, Sarah-Jane Dias, Raaghav Chanan...	India	March 2, 2019	

8807 rows × 12 columns

```
In [73]: df.shape
```

```
Out[73]: (8807, 12)
```

```
In [74]: # Checking for missing values in the dataset
df.isnull().sum()
```

```
Out[74]:
```

	0
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

we have missing values for the following  
columns(director,cast,country,date\_added,rating and duration)

# Handling Missing Values in Specific Columns

Filling NaN value with the column Cast, Director, Country, Rating and Duration to the initial dataframe

```
In [75]: country_mode = df['country'].mode()[0]  
country_mode
```

```
Out[75]: 'United States'
```

```
In [76]: movie_rating = df.loc[df['type'] == 'Movie', 'rating'].mode()[0]  
movie_rating
```

```
Out[76]: 'TV-MA'
```

```
In [77]: df["new"] = df.loc[df["type"] == "movie", "rating"]
```

```
In [78]: tv_rating = df.loc[df['type'] == 'TV Show', 'rating'].mode()[0]  
tv_rating
```

```
Out[78]: 'TV-MA'
```

```
In [79]: movie_duration_mode = df.loc[df['type'] == 'Movie', 'duration'].mode()[0]  
movie_duration_mode
```

```
Out[79]: '90 min'
```

```
In [80]: tv_duration_mode = df.loc[df['type'] == 'TV Show', 'duration'].mode()[0]  
tv_duration_mode
```

```
Out[80]: '1 Season'
```

- For the 'director' and 'cast' columns, we replace missing values with 'unknown director' and "unknown cast" to maintain data integrity and avoid any bias in the analysis.
- In the 'country' column, we fill in missing values with the mode (most frequently occurring value) to ensure consistency and minimize data loss.
- For the 'rating' column, we fill in missing values based on the 'type' of the show. We assign the mode of 'rating' for movies and TV shows separately.
- For the 'duration' column, we fill in missing values based on the 'type' of the show. We assign the mode of 'duration' for movies and TV shows

separately.

```
In [81]: df['director'].fillna('Unknown Director',inplace = True)
df['cast'].fillna('Unknown cast',inplace = True)
df['rating'] = df.apply(lambda x: movie_rating if x['type'] == 'Movie' and pd.
                        else tv_rating if x['type'] == 'TV Show' and pd.isna(x
                        else x['rating'], axis=1)
df['duration'] = df.apply(lambda x: movie_duration_mode if x['type'] == 'Movie
                        and pd.isna(x['duration'])
                        else tv_duration_mode if x['type'] == 'TV Show'
                        and pd.isna(x['duration'])
                        else x['duration'], axis=1)

df['country'].fillna( country_mode, inplace=True)

df
```

```
/tmp/ipython-input-1812104441.py:1: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment using an inplace method.
```

The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values always behaves as a copy.

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col].method(value) instead, to perform the operation inplace on the original object.

```
df['director'].fillna('Unknown Director',inplace = True)
```

```
/tmp/ipython-input-1812104441.py:2: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment using an inplace method.
```

The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values always behaves as a copy.

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col].method(value) instead, to perform the operation inplace on the original object.

```
df['cast'].fillna('Unknown cast',inplace = True)
```

```
/tmp/ipython-input-1812104441.py:13: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment using an inplace method.
```

The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values always behaves as a copy.

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col].method(value) instead, to perform the operation inplace on the original object.

```
df['country'].fillna( country_mode, inplace=True)
```

Out[81]:

	show_id	type	title	director	cast	country	date_added	rel
<b>0</b>	s1	Movie	Dick Johnson Is Dead	Kirsten Johnson	Unknown cast	United States	September 25, 2021	
<b>1</b>	s2	TV Show	Blood & Water	Unknown Director	Ama Qamata, Khosi Ngema, Gail Mabalane, Thaban...	South Africa	September 24, 2021	
<b>2</b>	s3	TV Show	Ganglands	Julien Leclercq	Sami Bouajila, Tracy Gotoas, Samuel Jouy, Nabi...	United States	September 24, 2021	
<b>3</b>	s4	TV Show	Jailbirds New Orleans	Unknown Director	Unknown cast	United States	September 24, 2021	
<b>4</b>	s5	TV Show	Kota Factory	Unknown Director	Mayur More, Jitendra Kumar, Ranjan Raj, Alam K...	India	September 24, 2021	
...	...	...	...	...	...	...	...	
<b>8802</b>	s8803	Movie	Zodiac	David Fincher	Mark Ruffalo, Jake Gyllenhaal, Robert Downey J...	United States	November 20, 2019	
<b>8803</b>	s8804	TV Show	Zombie Dumb	Unknown Director	Unknown cast	United States	July 1, 2019	
<b>8804</b>	s8805	Movie	Zombieland	Ruben Fleischer	Jesse Eisenberg, Woody Harrelson, Emma Stone, ...	United States	November 1, 2019	

	show_id	type	title	director	cast	country	date_added	rel
8805	s8806	Movie	Zoom	Peter Hewitt	Tim Allen, Courteney Cox, Chevy Chase, Kate Ma...	United States	January 11, 2020	
8806	s8807	Movie	Zubaan	Mozez Singh	Vicky Kaushal, Sarah-Jane Dias, Raaghav Chanan...	India	March 2, 2019	

8807 rows × 13 columns

In [82]: `df.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 8807 entries, 0 to 8806
Data columns (total 13 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        8807 non-null   object
4   cast            8807 non-null   object
5   country         8807 non-null   object
6   date_added      8797 non-null   object
7   release_year    8807 non-null   int64
8   rating          8807 non-null   object
9   duration        8807 non-null   object
10  listed_in       8807 non-null   object
11  description      8807 non-null   object
12  new             0 non-null      object
dtypes: int64(1), object(12)
memory usage: 894.6+ KB
```

## Un-nesting the dataset

unnesting column of director with titles

```
In [83]: split_director = df['director'].str.split(',', expand=True)
result_df = pd.concat([df['title'], split_director], axis=1)
dfdirector=pd.melt(result_df, id_vars=["title"], value_name="director")
dfdirector.drop(['variable'], axis=1, inplace=True)
unique_director=dfdirector.drop_duplicates()
```



```
dfunique_dr= unique_director.dropna()
dfunique_dr
```

Out[83]:

	<b>title</b>	<b>director</b>
<b>0</b>	Dick Johnson Is Dead	Kirsten Johnson
<b>1</b>	Blood & Water	Unknown Director
<b>2</b>	Ganglands	Julien Leclercq
<b>3</b>	Jailbirds New Orleans	Unknown Director
<b>4</b>	Kota Factory	Unknown Director
...	...	...
<b>95585</b>	Movie 43	Rusty Cundieff
<b>102764</b>	Walt Disney Animation Studios Short Films Coll...	Mike Gabriel
<b>103787</b>	HALO Legends	Hiroshi Yamazaki
<b>104392</b>	Movie 43	James Gunn
<b>111571</b>	Walt Disney Animation Studios Short Films Coll...	Mark Henn

9612 rows × 2 columns

## Unnesting column of cast with **titles**

```
In [84]: split_cast = df['cast'].str.split(',', expand=True)
result_df = pd.concat([df['title'], split_cast], axis=1)
dcast=pd.melt(result_df, id_vars=["title"], value_name="cast")
dcast.drop(['variable'], axis=1, inplace=True)
unique_cast=dcast.drop_duplicates()
dcastunique= unique_cast.dropna()
dcastunique
```

Out[84]:

	title	cast
0	Dick Johnson Is Dead	Unknown cast
1	Blood & Water	Ama Qamata
2	Ganglands	Sami Bouajila
3	Jailbirds New Orleans	Unknown cast
4	Kota Factory	Mayur More
...	...	...
417703	Black Mirror	Jon Hamm
424590	Social Distance	Ayize Ma'at
426510	Black Mirror	Oona Chaplin
433397	Social Distance	Lovie Simone
435317	Black Mirror	Rafe Spall

64949 rows × 2 columns

## Unnesting column of country with **titles**

```
In [85]: split_country = df['country'].str.split(',', expand=True)
result_df = pd.concat([df['title'], split_country], axis=1)
dfcountry=pd.melt(result_df, id_vars=["title"], value_name="country")
dfcountry.drop(['variable'], axis=1, inplace=True)
unique_country=dfcountry.drop_duplicates()
dfcountryunique= unique_country.dropna()
dfcountryunique
```

Out[85]:

	title	country
0	Dick Johnson Is Dead	United States
1	Blood & Water	South Africa
2	Ganglands	United States
3	Jailbirds New Orleans	United States
4	Kota Factory	India
...	...	...
78859	The Look of Silence	Germany
85496	Barbecue	Sweden
87666	The Look of Silence	Netherlands
94303	Barbecue	United States
103110	Barbecue	Uruguay

10850 rows × 2 columns

```
In [86]: split_listed_in = df['listed_in'].str.split(',', expand=True)
result_df = pd.concat([df['title'], split_listed_in], axis=1)
dflist=pd.melt(result_df, id_vars=["title"], value_name="listed_in")
dflist.drop(['variable'], axis=1, inplace=True)
dfunique_list=dflist.drop_duplicates()
dfunique_list= dfunique_list.dropna()
dfunique_list
```

Out[86]:

	title	listed_in
0	Dick Johnson Is Dead	Documentaries
1	Blood & Water	International TV Shows
2	Ganglands	Crime TV Shows
3	Jailbirds New Orleans	Docuseries
4	Kota Factory	International TV Shows
...	...	...
26414	Zindagi Gulzar Hai	TV Dramas
26415	Zinzana	Thrillers
26416	Zodiac	Thrillers
26417	Zombie Dumb	TV Comedies
26420	Zubaan	Music & Musicals

19323 rows × 2 columns

Merging and create the final table

```
In [87]: d=pd.merge(dfunique_dr, dcastunique,on="title", how="outer")
d1=pd.merge(d, dfcountryunique, on="title", how="outer" )
d2=pd.merge(d1, dfuniqueelist, on="title", how="outer" )
d2
```

Out[87]:

		title	director	cast	country	listed_in
<b>0</b>		#Alive	Cho Il	Yoo Ah-in	South Korea	Horror Movies
<b>1</b>		#Alive	Cho Il	Yoo Ah-in	South Korea	International Movies
<b>2</b>		#Alive	Cho Il	Yoo Ah-in	South Korea	Thrillers
<b>3</b>		#Alive	Cho Il	Park Shin-hye	South Korea	Horror Movies
<b>4</b>		#Alive	Cho Il	Park Shin-hye	South Korea	International Movies
<b>...</b>		...	...	...	...	...
<b>202053</b>	최강전사 미니특공대 : 영웅의 탄생		Young Jun Lee	Yang Jeong-hwa	United States	Children & Family Movies
<b>202054</b>	최강전사 미니특공대 : 영웅의 탄생		Young Jun Lee	Jeon Tae-yeol	United States	Children & Family Movies
<b>202055</b>	최강전사 미니특공대 : 영웅의 탄생		Young Jun Lee	Shin Yong-woo	United States	Children & Family Movies
<b>202056</b>	최강전사 미니특공대 : 영웅의 탄생		Young Jun Lee	Lee So-young	United States	Children & Family Movies
<b>202057</b>	최강전사 미니특공대 : 영웅의 탄생		Young Jun Lee	So-yeon	United States	Children & Family Movies

202058 rows × 5 columns

```
In [88]: d3=df[["title", "show_id", "type", "date_added", "release_year", "rating", "du
d3
```

Out[88]:

	title	show_id	type	date_added	release_year	rating	duration
<b>0</b>	Dick Johnson Is Dead	s1	Movie	September 25, 2021	2020	PG-13	90 min
<b>1</b>	Blood & Water	s2	TV Show	September 24, 2021	2021	TV-MA	2 Seasons
<b>2</b>	Ganglands	s3	TV Show	September 24, 2021	2021	TV-MA	1 Season
<b>3</b>	Jailbirds New Orleans	s4	TV Show	September 24, 2021	2021	TV-MA	1 Season
<b>4</b>	Kota Factory	s5	TV Show	September 24, 2021	2021	TV-MA	2 Seasons
...	...	...	...	...	...	...	...
<b>8802</b>	Zodiac	s8803	Movie	November 20, 2019	2007	R	158 min
<b>8803</b>	Zombie Dumb	s8804	TV Show	July 1, 2019	2018	TV-Y7	2 Seasons
<b>8804</b>	Zombieland	s8805	Movie	November 1, 2019	2009	R	88 min
<b>8805</b>	Zoom	s8806	Movie	January 11, 2020	2006	PG	88 min
<b>8806</b>	Zubaan	s8807	Movie	March 2, 2019	2015	TV-14	111 min

8807 rows × 7 columns

```
In [89]: df_final=pd.merge(d2, d3, on="title", how="inner")
df_final
```

Out[89]:									
	title	director	cast	country	listed_in	show_id	type	date_added	
0	#Alive	Cho Il	Yoo Ah-in	South Korea	Horror Movies	s2037	Movie	September 8, 2020	
1	#Alive	Cho Il	Yoo Ah-in	South Korea	International Movies	s2037	Movie	September 8, 2020	
2	#Alive	Cho Il	Yoo Ah-in	South Korea	Thrillers	s2037	Movie	September 8, 2020	
3	#Alive	Cho Il	Park Shin-hye	South Korea	Horror Movies	s2037	Movie	September 8, 2020	
4	#Alive	Cho Il	Park Shin-hye	South Korea	International Movies	s2037	Movie	September 8, 2020	
...	...	...	...	...	...	...	...	...	
202053	최강전사 미니특공대 : 영웅의 탄생	Young Jun Lee	Yang Jeong-hwa	United States	Children & Family Movies	s7109	Movie	September 1, 2020	
202054	최강전사 미니특공대 : 영웅의 탄생	Young Jun Lee	Jeon Tae-yeol	United States	Children & Family Movies	s7109	Movie	September 1, 2020	
202055	최강전사 미니특공대 : 영웅의 탄생	Young Jun Lee	Shin Yong-woo	United States	Children & Family Movies	s7109	Movie	September 1, 2020	
202056	최강전사 미니특공대 : 영웅의 탄생	Young Jun Lee	Lee So-young	United States	Children & Family Movies	s7109	Movie	September 1, 2020	
202057	최강전사 미니특공대 : 영웅의 탄생	Young Jun Lee	So-yeon	United States	Children & Family Movies	s7109	Movie	September 1, 2020	

202058 rows × 11 columns

```
In [90]: df_final
```



Out[90]:

	title	director	cast	country	listed_in	show_id	type	date_added
0	#Alive	Cho Il	Yoo Ah-in	South Korea	Horror Movies	s2037	Movie	September 8, 2020
1	#Alive	Cho Il	Yoo Ah-in	South Korea	International Movies	s2037	Movie	September 8, 2020
2	#Alive	Cho Il	Yoo Ah-in	South Korea	Thrillers	s2037	Movie	September 8, 2020
3	#Alive	Cho Il	Park Shin-hye	South Korea	Horror Movies	s2037	Movie	September 8, 2020
4	#Alive	Cho Il	Park Shin-hye	South Korea	International Movies	s2037	Movie	September 8, 2020
...	...	...	...	...	...	...	...	...
202053	최강전사 미니특공대 : 영웅의 탄생	Young Jun Lee	Yang Jeong-hwa	United States	Children & Family Movies	s7109	Movie	September 1, 2020
202054	최강전사 미니특공대 : 영웅의 탄생	Young Jun Lee	Jeon Tae-yeol	United States	Children & Family Movies	s7109	Movie	September 1, 2020
202055	최강전사 미니특공대 : 영웅의 탄생	Young Jun Lee	Shin Yong-woo	United States	Children & Family Movies	s7109	Movie	September 1, 2020
202056	최강전사 미니특공대 : 영웅의 탄생	Young Jun Lee	Lee So-young	United States	Children & Family Movies	s7109	Movie	September 1, 2020
202057	최강전사 미니특공대 : 영웅의 탄생	Young Jun Lee	So-yeon	United States	Children & Family Movies	s7109	Movie	September 1, 2020

202058 rows × 11 columns

```
In [91]: df_final.dropna(inplace=True)
```

```
In [92]: df_final.duplicated().sum()
```

```
Out[92]: np.int64(0)
```

```
In [93]: df_final
```

Out[93]:

	title	director	cast	country	listed_in	show_id	type	date_added
0	#Alive	Cho Il	Yoo Ah-in	South Korea	Horror Movies	s2037	Movie	September 8, 2020
1	#Alive	Cho Il	Yoo Ah-in	South Korea	International Movies	s2037	Movie	September 8, 2020
2	#Alive	Cho Il	Yoo Ah-in	South Korea	Thrillers	s2037	Movie	September 8, 2020
3	#Alive	Cho Il	Park Shin-hye	South Korea	Horror Movies	s2037	Movie	September 8, 2020
4	#Alive	Cho Il	Park Shin-hye	South Korea	International Movies	s2037	Movie	September 8, 2020
...	...	...	...	...	...	...	...	...
202053	최강전사 미 니특공대 : 영웅의 탄생	Young Jun Lee	Yang Jeong-hwa	United States	Children & Family Movies	s7109	Movie	September 1, 2020
202054	최강전사 미 니특공대 : 영웅의 탄생	Young Jun Lee	Jeon Tae-yeol	United States	Children & Family Movies	s7109	Movie	September 1, 2020
202055	최강전사 미 니특공대 : 영웅의 탄생	Young Jun Lee	Shin Yong-woo	United States	Children & Family Movies	s7109	Movie	September 1, 2020
202056	최강전사 미 니특공대 : 영웅의 탄생	Young Jun Lee	Lee So-young	United States	Children & Family Movies	s7109	Movie	September 1, 2020
202057	최강전사 미 니특공대 : 영웅의 탄생	Young Jun Lee	So-yeon	United States	Children & Family Movies	s7109	Movie	September 1, 2020

201900 rows × 11 columns

```
In [94]: df_final["duration"]=df_final["duration"].str.split(' ').str[0]  
df_final['duration'] = df_final['duration'].astype(int)
```

```
In [95]: df_final
```

Out[95]:

	title	director	cast	country	listed_in	show_id	type	date_added
0	#Alive	Cho Il	Yoo Ah-in	South Korea	Horror Movies	s2037	Movie	September 8, 2020
1	#Alive	Cho Il	Yoo Ah-in	South Korea	International Movies	s2037	Movie	September 8, 2020
2	#Alive	Cho Il	Yoo Ah-in	South Korea	Thrillers	s2037	Movie	September 8, 2020
3	#Alive	Cho Il	Park Shin-hye	South Korea	Horror Movies	s2037	Movie	September 8, 2020
4	#Alive	Cho Il	Park Shin-hye	South Korea	International Movies	s2037	Movie	September 8, 2020
...	...	...	...	...	...	...	...	...
202053	최강전사 미 니특공대 : 영웅의 탄생	Young Jun Lee	Yang Jeong-hwa	United States	Children & Family Movies	s7109	Movie	September 1, 2020
202054	최강전사 미 니특공대 : 영웅의 탄생	Young Jun Lee	Jeon Tae-yeol	United States	Children & Family Movies	s7109	Movie	September 1, 2020
202055	최강전사 미 니특공대 : 영웅의 탄생	Young Jun Lee	Shin Yong-woo	United States	Children & Family Movies	s7109	Movie	September 1, 2020
202056	최강전사 미 니특공대 : 영웅의 탄생	Young Jun Lee	Lee So-young	United States	Children & Family Movies	s7109	Movie	September 1, 2020
202057	최강전사 미 니특공대 : 영웅의 탄생	Young Jun Lee	So-yeon	United States	Children & Family Movies	s7109	Movie	September 1, 2020

201900 rows × 11 columns

We convert the 'date\_added' column to datetime format using `pd.to_datetime()` to enable further analysis based on date-related attributes.

```
In [96]: df_final["date_added"] = pd.to_datetime(df_final['date_added'], format='%B %d',
df_final["date_added"] = pd.to_datetime(df_final['date_added'])
df_final
```

Out[96]:

	title	director	cast	country	listed_in	show_id	type	date_added
<b>0</b>	#Alive	Cho Il	Yoo Ah-in	South Korea	Horror Movies	s2037	Movie	2020-09-01
<b>1</b>	#Alive	Cho Il	Yoo Ah-in	South Korea	International Movies	s2037	Movie	2020-09-01
<b>2</b>	#Alive	Cho Il	Yoo Ah-in	South Korea	Thrillers	s2037	Movie	2020-09-01
<b>3</b>	#Alive	Cho Il	Park Shin-hye	South Korea	Horror Movies	s2037	Movie	2020-09-01
<b>4</b>	#Alive	Cho Il	Park Shin-hye	South Korea	International Movies	s2037	Movie	2020-09-01
...	...	...	...	...	...	...	...	...
<b>202053</b>	최강전사 미 니특공대 : 영웅의 탄생	Young Jun Lee	Yang Jeong-hwa	United States	Children & Family Movies	s7109	Movie	2018-09-01
<b>202054</b>	최강전사 미 니특공대 : 영웅의 탄생	Young Jun Lee	Jeon Tae-yeol	United States	Children & Family Movies	s7109	Movie	2018-09-01
<b>202055</b>	최강전사 미 니특공대 : 영웅의 탄생	Young Jun Lee	Shin Yong-woo	United States	Children & Family Movies	s7109	Movie	2018-09-01
<b>202056</b>	최강전사 미 니특공대 : 영웅의 탄생	Young Jun Lee	Lee So-young	United States	Children & Family Movies	s7109	Movie	2018-09-01
<b>202057</b>	최강전사 미 니특공대 : 영웅의 탄생	Young Jun Lee	So-yeon	United States	Children & Family Movies	s7109	Movie	2018-09-01

201900 rows × 11 columns

## Q1. Defining Problem Statement and Analysing basic metrics

Our primary objective is to identify the most promising types of shows to produce, thereby maximizing our growth potential in the entertainment industry. To achieve this goal, we need to comprehensively analyze the data and derive actionable insights that will enable us to make informed decisions on content creation and business expansion.

## Analysing Basic Metrics

```
In [97]: df_final.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Index: 201900 entries, 0 to 202057
Data columns (total 11 columns):
#   Column          Non-Null Count  Dtype
---  -
0   title           201900 non-null object
1   director        201900 non-null object
2   cast            201900 non-null object
3   country         201900 non-null object
4   listed_in       201900 non-null object
5   show_id         201900 non-null object
6   type            201900 non-null object
7   date_added      200312 non-null datetime64[ns]
8   release_year    201900 non-null int64
9   rating          201900 non-null object
10  duration        201900 non-null int64
dtypes: datetime64[ns](1), int64(2), object(8)
memory usage: 18.5+ MB
```

```
In [98]: unique_counts = df_final.nunique()
print(unique_counts)
```



```
title            8797
director         5121
cast            39261
country          197
listed_in        73
show_id          8797
type             2
date_added       1699
release_year     74
rating           17
duration         210
dtype: int64
```

## **Q2. Observations on the shape of data, data types of all the attributes, conversion of categorical attributes to 'category' (If required), missing value detection, statistical summary**

```
In [99]: df_final.shape
```

```
Out[99]: (201900, 11)
```

```
In [100... df_final.dtypes
```

Out[100...

**0**

<b>title</b>	object
<b>director</b>	object
<b>cast</b>	object
<b>country</b>	object
<b>listed_in</b>	object
<b>show_id</b>	object
<b>type</b>	object
<b>date_added</b>	datetime64[ns]
<b>release_year</b>	int64
<b>rating</b>	object
<b>duration</b>	int64

**dtype:** object

In [101...

```
# Checking for missing values in the dataset
df_final.isnull().sum()
```

Out[101...

**0**

<b>title</b>	0
<b>director</b>	0
<b>cast</b>	0
<b>country</b>	0
<b>listed_in</b>	0
<b>show_id</b>	0
<b>type</b>	0
<b>date_added</b>	1588
<b>release_year</b>	0
<b>rating</b>	0
<b>duration</b>	0

**dtype:** int64

- There is no need of changing any column to categorical column.

- Missing value already detected and fixed my filling it with unknown data and remaining filled with mode values We are going to ignore the missing value of "date\_added" since very few data are missing

```
In [105... df_final.drop_duplicates(inplace=True)
```

```
In [106... #Statistical summary
df_final.describe().round(2)
```

```
Out[106...
```

	date_added	release_year	duration
<b>count</b>	200312	201900.00	201900.00
<b>mean</b>	2019-06-24 18:04:32.167418880	2013.45	77.74
<b>min</b>	2008-01-01 00:00:00	1925.00	1.00
<b>25%</b>	2018-07-01 00:00:00	2012.00	4.00
<b>50%</b>	2019-09-13 00:00:00	2016.00	95.00
<b>75%</b>	2020-09-15 00:00:00	2019.00	112.00
<b>max</b>	2021-09-25 00:00:00	2021.00	312.00
<b>std</b>	NaN	9.02	51.46

```
In [107... df_final.describe(include="datetime")
```

```
Out[107...
```

	date_added
<b>count</b>	200312
<b>mean</b>	2019-06-24 18:04:32.167418880
<b>min</b>	2008-01-01 00:00:00
<b>25%</b>	2018-07-01 00:00:00
<b>50%</b>	2019-09-13 00:00:00
<b>75%</b>	2020-09-15 00:00:00
<b>max</b>	2021-09-25 00:00:00

submitted by Kanimozhi