

Netflix Movies Data Analysis

```
In [1]: # Importing the libraries

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

```
In [2]: # importing our dataset

data = pd.read_csv(r'G:\dataset files\mymoviedb.csv', lineterminator = '\n')
```

```
In [3]: data.head()
```

```
Out[3]:
```

	Release_Date	Title	Overview	Popularity	Vote_Count	Vote_Average	Original_Language
0	2021-12-15	Spider-Man: No Way Home	Peter Parker is unmasked and no longer able to...	5083.954	8940	8.3	en
1	2022-03-01	The Batman	In his second year of fighting crime, Batman u...	3827.658	1151	8.1	en
2	2022-02-25	No Exit	Stranded at a rest stop in the mountains durin...	2618.087	122	6.3	en
3	2021-11-24	Encanto	The tale of an extraordinary family, the Madri...	2402.201	5076	7.7	en
4	2021-12-22	The King's Man	As a collection of history's worst tyrants and...	1895.511	1793	7.0	en

```
In [4]: data.info()
```



```
Out[8]: 0    2021-12-15
        1    2022-03-01
        2    2022-02-25
        3    2021-11-24
        4    2021-12-22
        Name: Release_Date, dtype: datetime64[ns]
```

```
In [9]: data['Release_Date'] = data['Release_Date'].dt.year
```

```
In [10]: data.head()
```

```
Out[10]:
```

	Release_Date	Title	Overview	Popularity	Vote_Count	Vote_Average	Original_Language
0	2021	Spider-Man: No Way Home	Peter Parker is unmasked and no longer able to...	5083.954	8940	8.3	€
1	2022	The Batman	In his second year of fighting crime, Batman u...	3827.658	1151	8.1	€
2	2022	No Exit	Stranded at a rest stop in the mountains durin...	2618.087	122	6.3	€
3	2021	Encanto	The tale of an extraordinary family, the Madri...	2402.201	5076	7.7	€
4	2021	The King's Man	As a collection of history's worst tyrants and...	1895.511	1793	7.0	€

```
In [11]: data['Release_Date'].dtypes
```

```
Out[11]: dtype('int64')
```

- here 'Overview', 'Original_Language', 'Poster_Url' wont be so useful, so we will drop them.

```
In [12]: x = ['Overview', 'Original_Language', 'Poster_Url']
```

```
In [13]: x
```

```
Out[13]: ['Overview', 'Original_Language', 'Poster_Url']
```

```
In [14]: data.drop(x, axis=1, inplace=True)
```

```
In [15]: data.head()
```

Out[15]:	Release_Date	Title	Popularity	Vote_Count	Vote_Average	Genre
0	2021	Spider-Man: No Way Home	5083.954	8940	8.3	Action, Adventure, Science Fiction
1	2022	The Batman	3827.658	1151	8.1	Crime, Mystery, Thriller
2	2022	No Exit	2618.087	122	6.3	Thriller
3	2021	Encanto	2402.201	5076	7.7	Animation, Comedy, Family, Fantasy
4	2021	The King's Man	1895.511	1793	7.0	Action, Adventure, Thriller, War

- categorizing 'Vote_Average' column
- we will cut the 'Vote_Average' values and make 4 categories: 'not popular', 'below Average', 'Average', 'Popular' to describing it .

```
In [16]: def categorical_cols(data, col, labels):
    edges = [data[col].describe()['min'],
              data[col].describe()['25%'],
              data[col].describe()['50%'],
              data[col].describe()['75%'],
              data[col].describe()['max'],]

    data[col] = pd.cut(data[col], edges, labels = labels, duplicates = 'drop')
    return data
```

- creating Labels

```
In [17]: labels = ['not popular', 'below Average', 'Average', 'Popular']
```

```
In [18]: categorical_cols(data, 'Vote_Average', labels)
```

```
data['Vote_Average'].unique()
```

```
Out[18]: ['Popular', 'below Average', 'Average', 'not popular', NaN]
Categories (4, object): ['not popular' < 'below Average' < 'Average' < 'Popular']
```

```
In [19]: data.head()
```

Out[19]:

	Release_Date	Title	Popularity	Vote_Count	Vote_Average	Genre
0	2021	Spider-Man: No Way Home	5083.954	8940	Popular	Action, Adventure, Science Fiction
1	2022	The Batman	3827.658	1151	Popular	Crime, Mystery, Thriller
2	2022	No Exit	2618.087	122	below Average	Thriller
3	2021	Encanto	2402.201	5076	Popular	Animation, Comedy, Family, Fantasy
4	2021	The King's Man	1895.511	1793	Average	Action, Adventure, Thriller, War

In [20]: `data['Vote_Average'].value_counts()`

Out[20]:

```
not popular    2467
Popular        2450
Average        2412
below Average  2398
Name: Vote_Average, dtype: int64
```

In [21]: `data.dropna(inplace=True)`
`data.isna().sum()`

Out[21]:

```
Release_Date    0
Title           0
Popularity      0
Vote_Count      0
Vote_Average    0
Genre           0
dtype: int64
```

In [22]: `data.head()`

Out[22]:

	Release_Date	Title	Popularity	Vote_Count	Vote_Average	Genre
0	2021	Spider-Man: No Way Home	5083.954	8940	Popular	Action, Adventure, Science Fiction
1	2022	The Batman	3827.658	1151	Popular	Crime, Mystery, Thriller
2	2022	No Exit	2618.087	122	below Average	Thriller
3	2021	Encanto	2402.201	5076	Popular	Animation, Comedy, Family, Fantasy
4	2021	The King's Man	1895.511	1793	Average	Action, Adventure, Thriller, War

- now split Genre into a list and then explode our dataset to have only one Genre per row for each movie.

In [23]: `data['Genre'] = data['Genre'].str.split(',')`
`data = data.explode('Genre').reset_index(drop=True)`
`data.head()`

Out[27]:	Release_Date	Title	Popularity	Vote_Count	Vote_Average	Genre
0	2021	Spider-Man: No Way Home	5083.954	8940	Popular	Action
1	2021	Spider-Man: No Way Home	5083.954	8940	Popular	Adventure
2	2021	Spider-Man: No Way Home	5083.954	8940	Popular	Science Fiction
3	2022	The Batman	3827.658	1151	Popular	Crime
4	2022	The Batman	3827.658	1151	Popular	Mystery

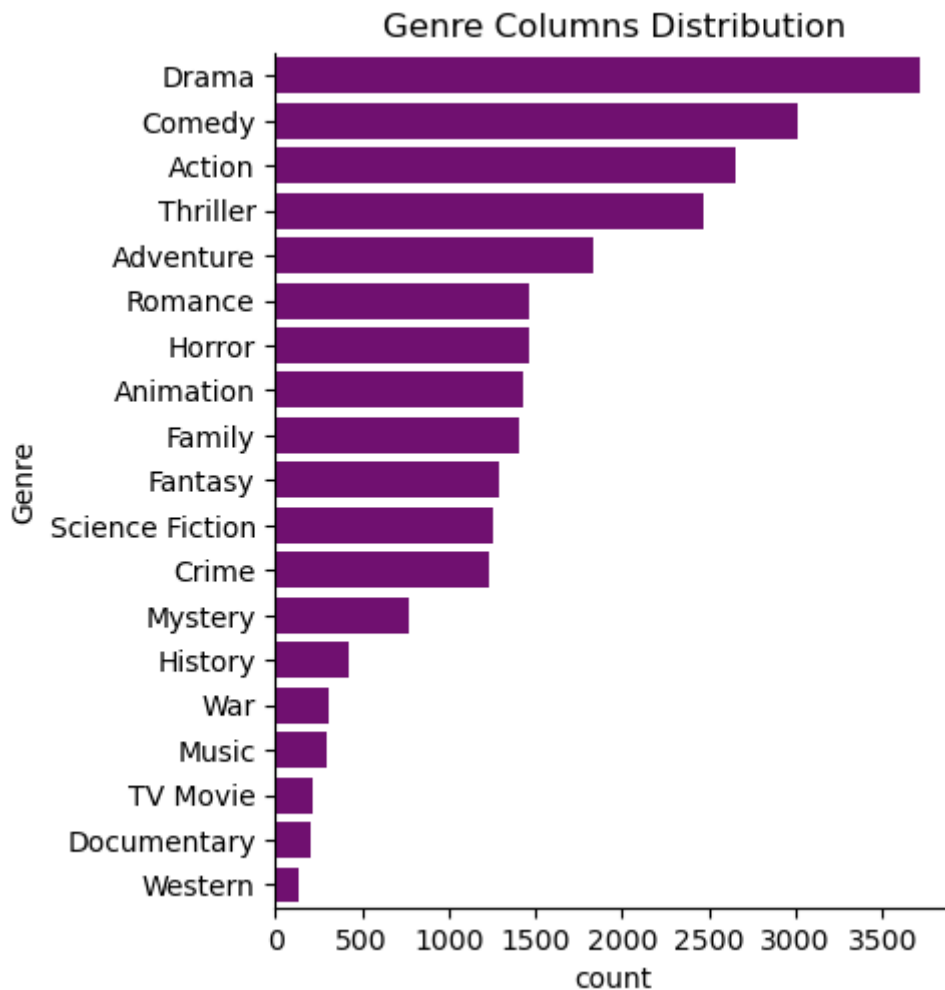
Here is the endup of DATA PREPROSSEING
PROCESS

1. What is the most frequent genre of
movies released on Netflix?

```
In [28]: data['Genre'].describe()
```

```
Out[28]: count      25552  
unique         19  
top      Drama  
freq         3715  
Name: Genre, dtype: object
```

```
In [29]: sns.catplot(y = 'Genre', data = data,  
                    kind = 'count',  
                    order = data['Genre'].value_counts().index,  
                    color='purple')  
  
plt.title("Genre Columns Distribution")  
plt.show()
```



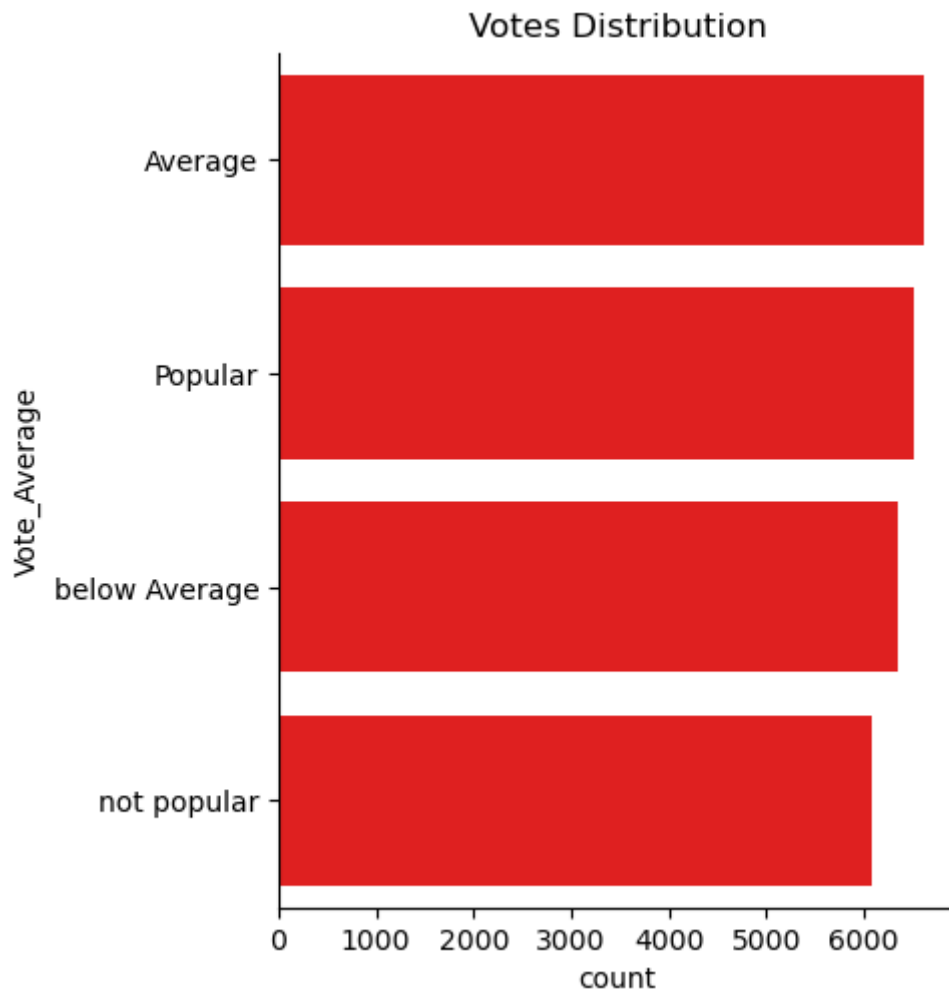
2. Which has highest votes in vote avg column?

```
In [30]: data.head()
```

```
Out[30]:
```

	Release_Date	Title	Popularity	Vote_Count	Vote_Average	Genre
0	2021	Spider-Man: No Way Home	5083.954	8940	Popular	Action
1	2021	Spider-Man: No Way Home	5083.954	8940	Popular	Adventure
2	2021	Spider-Man: No Way Home	5083.954	8940	Popular	Science Fiction
3	2022	The Batman	3827.658	1151	Popular	Crime
4	2022	The Batman	3827.658	1151	Popular	Mystery

```
In [31]: sns.catplot(y = 'Vote_Average', data = data,
                    kind = 'count',
                    order = data['Vote_Average'].value_counts().index,
                    color = 'red')
plt.title("Votes Distribution")
plt.show()
```

3. What movie got the highest popularity? what's its genre?

```
In [32]: data.head()
```

Out[32]:	Release_Date	Title	Popularity	Vote_Count	Vote_Average	Genre
0	2021	Spider-Man: No Way Home	5083.954	8940	Popular	Action
1	2021	Spider-Man: No Way Home	5083.954	8940	Popular	Adventure
2	2021	Spider-Man: No Way Home	5083.954	8940	Popular	Science Fiction
3	2022	The Batman	3827.658	1151	Popular	Crime
4	2022	The Batman	3827.658	1151	Popular	Mystery

```
In [33]: data[data['Popularity'] == data['Popularity'].max()]
```

Out[33]:

	Release_Date	Title	Popularity	Vote_Count	Vote_Average	Genre
0	2021	Spider-Man: No Way Home	5083.954	8940	Popular	Action
1	2021	Spider-Man: No Way Home	5083.954	8940	Popular	Adventure
2	2021	Spider-Man: No Way Home	5083.954	8940	Popular	Science Fiction

4. What movie got the lowest popularity? what's its genre?

```
In [34]: data[data['Popularity'] == data['Popularity'].min()]
```

Out[34]:

	Release_Date	Title	Popularity	Vote_Count	Vote_Average	Genre
25546	2021	The United States vs. Billie Holiday	13.354	152	Average	Music
25547	2021	The United States vs. Billie Holiday	13.354	152	Average	Drama
25548	2021	The United States vs. Billie Holiday	13.354	152	Average	History
25549	1984	Threads	13.354	186	Popular	War
25550	1984	Threads	13.354	186	Popular	Drama
25551	1984	Threads	13.354	186	Popular	Science Fiction

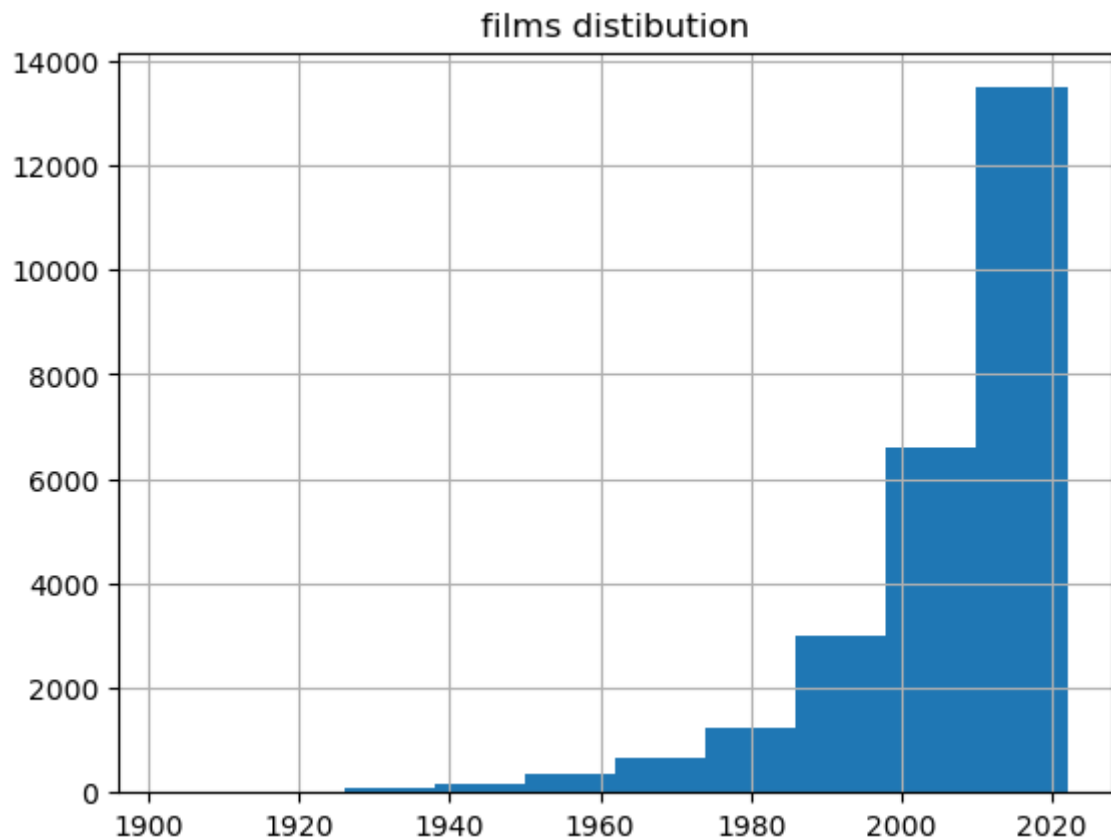
5. Which year has the most filmed movies?

```
In [35]: data.head()
```

Out[35]:

	Release_Date	Title	Popularity	Vote_Count	Vote_Average	Genre
0	2021	Spider-Man: No Way Home	5083.954	8940	Popular	Action
1	2021	Spider-Man: No Way Home	5083.954	8940	Popular	Adventure
2	2021	Spider-Man: No Way Home	5083.954	8940	Popular	Science Fiction
3	2022	The Batman	3827.658	1151	Popular	Crime
4	2022	The Batman	3827.658	1151	Popular	Mystery

```
In [36]: data['Release_Date'].hist()  
plt.title("films distribution")  
plt.show()
```



Conclusion

Q1: What is the most frequent genre of movies released on Netflix?

Ans: Drama is the most frequent genre in our dataset and has appeared more than 14% of the times among 19 other genres.

Q2: Which has highest votes in vote avg column?

Ans: we have 25.5% of our dataset with popular vote (6520 rows). Drama again get the higher popularity among fans by being having more than 18.5% movies.

Q3: What movie got the highest popularity? what's its genre?

Ans: Spider-Man: No Way Home got the highest popularity of 5083.954 Its Genre is Action.

Q4: What movie got the lowest popularity? what's its genre?

Ans: The United States vs. Billie Holiday got the lowest popularity of 13.354 Its Genre is Music.

Q5: Which year has the most filmed movies?

Ans: Year 2020 has the most filmed movies.

The End