# Netflix Data Analysis

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

```python
df = pd.read_csv('mymoviedb.csv', lineterminator= '\n')
```

```python
df.head()
```

| | Release_Date | Title | Overview | Popularity | Vote_Count | Vote_Average | Original_Language | Genre | |
|---|---|---|---|---|---|---|---|---|---|
| **0** | 2021-12-15 | Spider-Man: No Way Home | Peter Parker is unmasked and no longer able to... | 5083.954 | 8940 | 8.3 | en | Action, Adventure, Science Fiction | https://image.tmdb.org/t/p/o |
| **1** | 2022-03-01 | The Batman | In his second year of fighting crime, Batman u... | 3827.658 | 1151 | 8.1 | en | Crime, Mystery, Thriller | https://image.tmdb.org/t/p/or |
| **2** | 2022-02-25 | No Exit | Stranded at a rest stop in the mountains durin... | 2618.087 | 122 | 6.3 | en | Thriller | https://image.tmdb.org/t/p/orig |
| **3** | 2021-11-24 | Encanto | The tale of an extraordinary family, the Madri... | 2402.201 | 5076 | 7.7 | en | Animation, Comedy, Family, Fantasy | https://image.tmdb.org/t/p/ori |
| **4** | 2021-12-22 | The King's Man | As a collection of history's worst tyrants and... | 1895.511 | 1793 | 7.0 | en | Action, Adventure, Thriller, War | https://image.tmdb.org/t/p/ori |

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 9827 entries, 0 to 9826
Data columns (total 9 columns):
 #   Column            Non-Null Count  Dtype
---  ------            --------------  -----
 0   Release_Date      9827 non-null   object
 1   Title             9827 non-null   object
 2   Overview          9827 non-null   object
 3   Popularity        9827 non-null   float64
 4   Vote_Count        9827 non-null   int64
 5   Vote_Average      9827 non-null   float64
 6   Original_Language 9827 non-null   object
 7   Genre             9827 non-null   object
 8   Poster_Url        9827 non-null   object
dtypes: float64(2), int64(1), object(6)
memory usage: 691.1+ KB
```

In [27]: `df['Genre'].head()`

Out[27]:
```
0    Action, Adventure, Science Fiction
1               Crime, Mystery, Thriller
2                               Thriller
3    Animation, Comedy, Family, Fantasy
4       Action, Adventure, Thriller, War
Name: Genre, dtype: object
```

In [28]: `df.duplicated().sum()`

Out[28]: `np.int64(0)`

In [29]: `df.describe()`

|       | Popularity  | Vote_Count  | Vote_Average |
|-------|-------------|-------------|--------------|
| count | 9827.000000 | 9827.000000 | 9827.000000  |
| mean  | 40.326088   | 1392.805536 | 6.439534     |
| std   | 108.873998  | 2611.206907 | 1.129759     |
| min   | 13.354000   | 0.000000    | 0.000000     |
| 25%   | 16.128500   | 146.000000  | 5.900000     |
| 50%   | 21.199000   | 444.000000  | 6.500000     |
| 75%   | 35.191500   | 1376.000000 | 7.100000     |
| max   | 5083.954000 | 31077.000000| 10.000000    |

• Exploration summary

• Dataframe contains 9827 rows and 9 columns. • Our dataset looks a bit tidy with no NaNs nor duplicated values. • Release_Date column needs to be casted into date time and to extract only the year values. • Overview, Original_Language and Poster-Url not needed, so we'll drop them, • There is noticeable outliers n Popularity column. • Vote_Average better be categorised for proper analysis. • Genre column has coma separated values and white spaces that need to be handled and casted into categories.

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

| | Release_Date | Title | Overview | Popularity | Vote_Count | Vote_Average | Original_Language | Genre | |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 2021-12-15 | Spider-Man: No Way Home | Peter Parker is unmasked and no longer able to... | 5083.954 | 8940 | 8.3 | en | Action, Adventure, Science Fiction | https://image.tmdb.org/t/p/o |
| 1 | 2022-03-01 | The Batman | In his second year of fighting crime, Batman u... | 3827.658 | 1151 | 8.1 | en | Crime, Mystery, Thriller | https://image.tmdb.org/t/p/or |
| 2 | 2022-02-25 | No Exit | Stranded at a rest stop in the mountains durin... | 2618.087 | 122 | 6.3 | en | Thriller | https://image.tmdb.org/t/p/orig |
| 3 | 2021-11-24 | Encanto | The tale of an extraordinary family, the Madri... | 2402.201 | 5076 | 7.7 | en | Animation, Comedy, Family, Fantasy | https://image.tmdb.org/t/p/ori |
| 4 | 2021-12-22 | The King's Man | As a collection of history's worst tyrants and... | 1895.511 | 1793 | 7.0 | en | Action, Adventure, Thriller, War | https://image.tmdb.org/t/p/ori |

In [31]:
```python
df['Release_Date'] = pd.to_datetime(df['Release_Date'])
print(df['Release_Date'].dtypes)
```

datetime64[ns]

In [32]:
```python
df['Release_Date'] = df['Release_Date'].dt.year
```

```python
df['Release_Date'].dtypes
```

Out[32]: dtype('int32')

In [33]: 
```python
df.head(3)
```

Out[33]:

| | Release_Date | Title | Overview | Popularity | Vote_Count | Vote_Average | Original_Language | Genre | |
|---|---|---|---|---|---|---|---|---|---|
| **0** | 2021 | Spider-Man: No Way Home | Peter Parker is unmasked and no longer able to... | 5083.954 | 8940 | 8.3 | en | Action, Adventure, Science Fiction | https://image.tmdb.org/t/p/origi |
| **1** | 2022 | The Batman | In his second year of fighting crime, Batman u... | 3827.658 | 1151 | 8.1 | en | Crime, Mystery, Thriller | https://image.tmdb.org/t/p/origi |
| **2** | 2022 | No Exit | Stranded at a rest stop in the mountains durin... | 2618.087 | 122 | 6.3 | en | Thriller | https://image.tmdb.org/t/p/origina |

Dropping the columns

In [34]: 
```python
cols = ['Overview', 'Original_Language', 'Poster_Url']
```

In [35]: 
```python
df.drop(cols, axis=1, inplace = True)
df.columns
```

```
Out[35]: Index(['Release_Date', 'Title', 'Popularity', 'Vote_Count', 'Vote_Average',
                'Genre'],
               dtype='object')
```

In [36]: `df.head(3)`

Out[36]:

| | 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 |

# Categorizing Vote_Average

We would cut the Vote_Average values and make 4 categories :popular average below_avg not_popular to describe it more using categorize_col() function

```python
In [37]: def categorize_col(df, col, labels):

             edges = [df[col].describe()['min'],
                     df[col].describe()['25%'],
                     df[col].describe()['50%'],
                     df[col].describe()['75%'],
                     df[col].describe()['max']]
             df[col] = pd.cut(df[col], edges, labels = labels, duplicates = 'drop')
             return df
```

```python
In [38]: labels = ['not_popular', 'below_avg', 'average', 'popular']

         categorize_col(df, 'Vote_Average', labels)
         df['Vote_Average'].unique()
```

```
Out[38]: ['popular', 'below_avg', 'average', 'not_popular', NaN]
         Categories (4, object): ['not_popular' < 'below_avg' < 'average' < 'popular']
```

```
In [39]: df.head()
```

Out[39]:

| | 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_avg | 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 [40]: df['Vote_Average'].value_counts()
```

Out[40]:
```
Vote_Average
not_popular    2467
popular        2450
average        2412
below_avg      2398
Name: count, dtype: int64
```

```
In [42]: df.dropna(inplace= True)

         df.isna().sum()
```

Out[42]:
```
Release_Date    0
Title           0
Popularity      0
Vote_Count      0
Vote_Average    0
Genre           0
dtype: int64
```

we'd split genres into a list and then explode our dataframe to have only one genre per row for each movie

```
In [44]: df['Genre'] = df['Genre'].str.split(', ')
         df = df.explode('Genre').reset_index(drop= True)
         df.head()
```

Out[44]:

| | 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 [46]: #casting column into category

         df['Genre'] = df['Genre'].astype('category')
         df['Genre'].dtypes
```

Out[46]: CategoricalDtype(categories=['Action', 'Adventure', 'Animation', 'Comedy', 'Crime',
                         'Documentary', 'Drama', 'Family', 'Fantasy', 'History',
                         'Horror', 'Music', 'Mystery', 'Romance', 'Science Fiction',
                         'TV Movie', 'Thriller', 'War', 'Western'],
         , ordered=False, categories_dtype=object)

```
In [47]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 25552 entries, 0 to 25551
Data columns (total 6 columns):
 #   Column        Non-Null Count  Dtype
---  ------        --------------  -----
 0   Release_Date  25552 non-null  int32
 1   Title         25552 non-null  object
 2   Popularity    25552 non-null  float64
 3   Vote_Count    25552 non-null  int64
 4   Vote_Average  25552 non-null  category
 5   Genre         25552 non-null  category
dtypes: category(2), float64(1), int32(1), int64(1), object(1)
memory usage: 749.6+ KB
```

In [48]: `df.nunique()`

Out[48]:
```
Release_Date     100
Title           9415
Popularity      8088
Vote_Count      3265
Vote_Average       4
Genre             19
dtype: int64
```
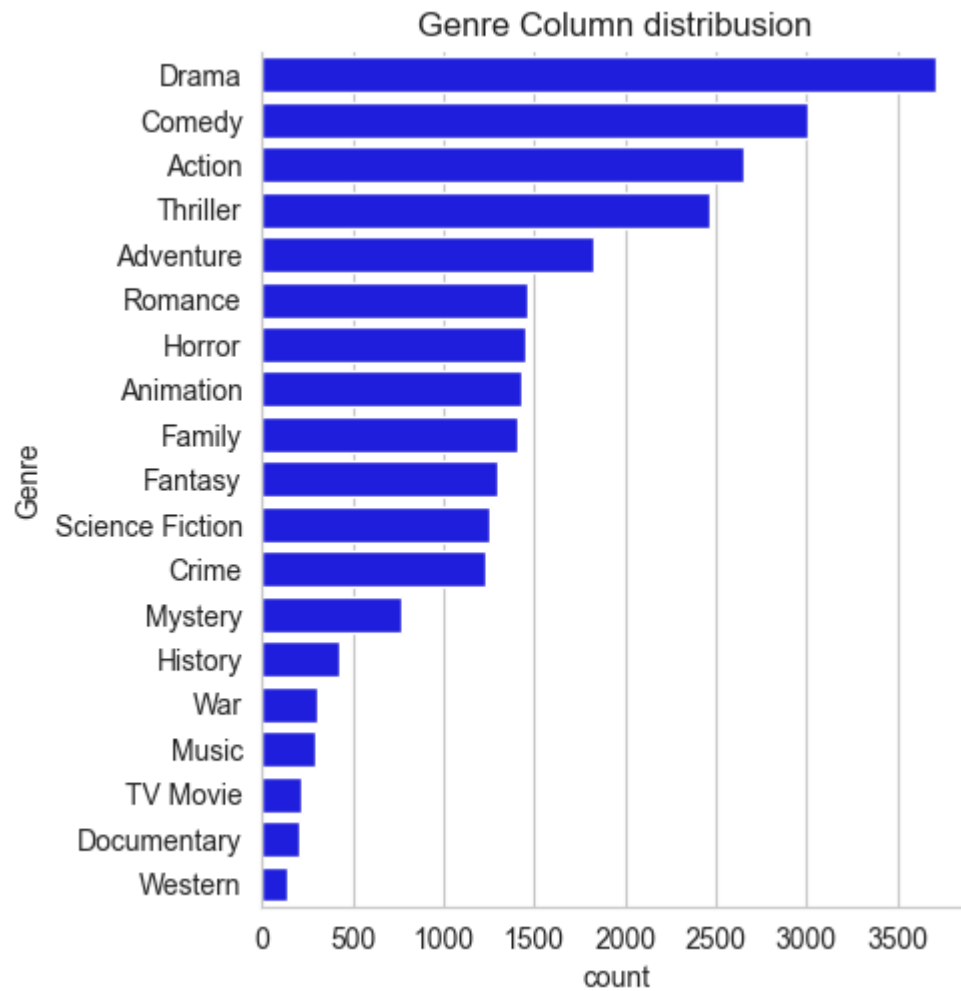
## Data Visualization

In [49]: `sns.set_style('whitegrid')`

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

In [54]: `df['Genre'].describe()`

```
Out[54]:  count     25552
          unique       19
          top       Drama
          freq       3715
          Name: Genre, dtype: object
```

```
In [56]:  sns.catplot(y= 'Genre', data = df, kind = 'count',
                      order = df['Genre'].value_counts().index,
                      color = 'blue')
          plt.title('Genre Column distribusion')
          plt.show()
```
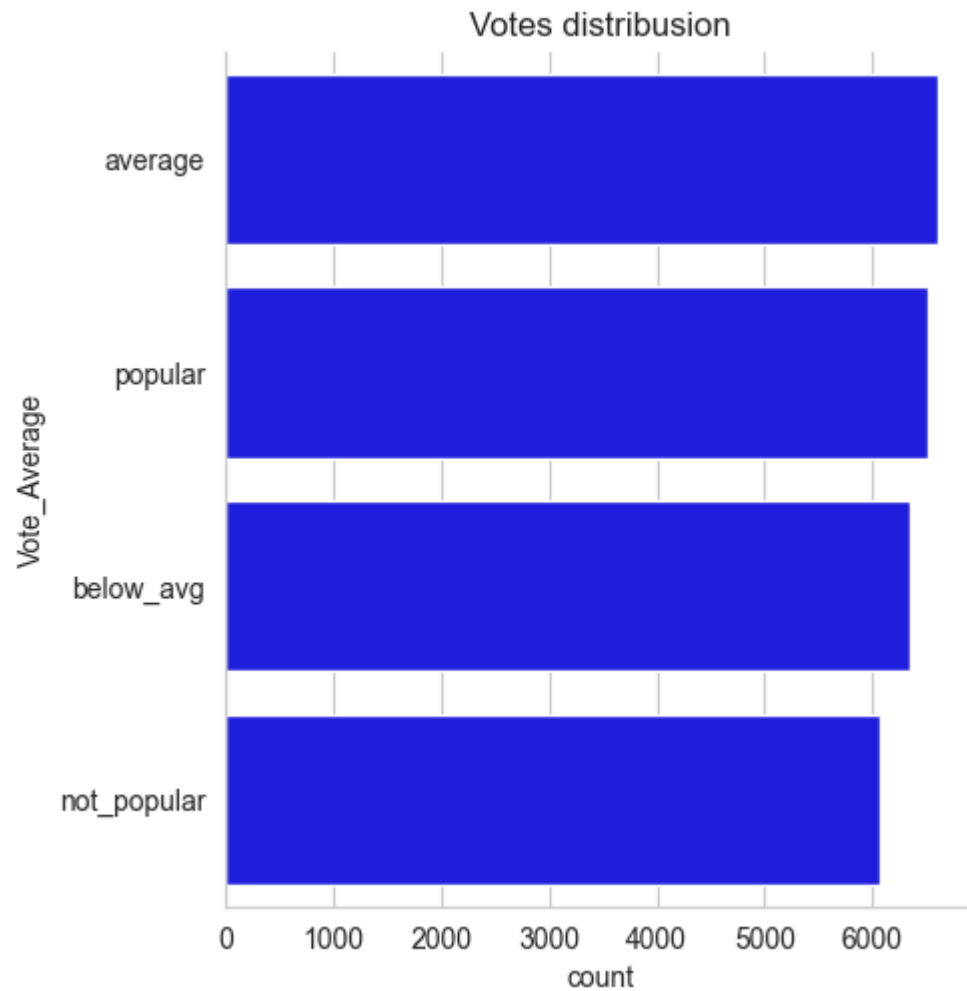
Genre Column distribusion

## 2)Which genres has highest votes in vote avg column?

```
In [57]: df.head()
```

Out[57]:

| | 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 [58]:
```python
sns.catplot(y= 'Vote_Average', data = df, kind = 'count',
            order = df['Vote_Average'].value_counts().index,
            color = 'blue')
plt.title('Votes distribusion')
plt.show()
```

Votes distribusion

## 3)Which movie got the highest popularity?What is its genre?

In [59]: `df.head(2)`

| | 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 |

In [60]:
```python
df[df['Popularity'] == df['Popularity'].max()]
```

| | 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)Which movie got the lowest popularity?What is its genre?

In [61]:
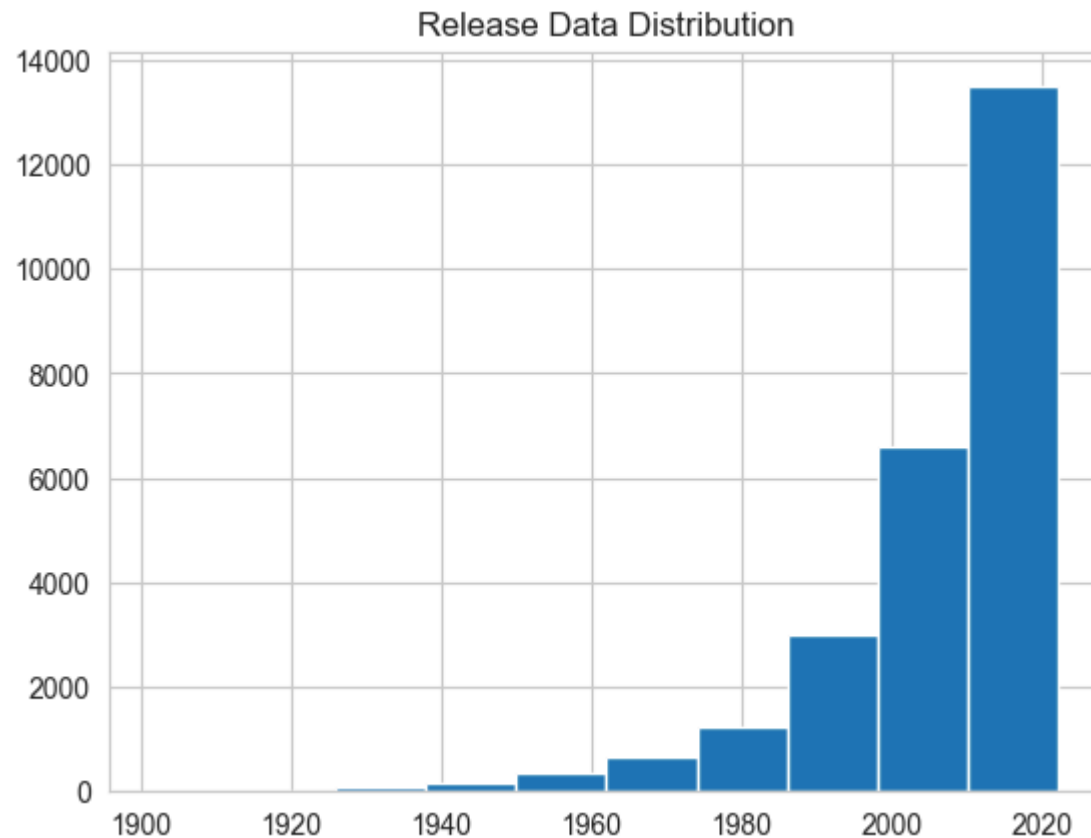```python
df[df['Popularity'] == df['Popularity'].min()]
```

| | 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 most filmmed movies?

```
In [62]:  df['Release_Date'].hist()
          plt.title('Release Data Distribution')
          plt.show()
```

Release Data Distribution



## Conclusion

Q1: What is the most frequent genre in the dataset? Drams genre is the most frequent genre in our dataset and has appeared more than 14% of the times among 19 other genres.

Q2: Which genres has highest votes? We have 25.5% of our dataset with popular vote (6520 rows). Drama again gets the highest popularity.

Q3: Which movie got the highest popularity? What is its genre? Spider-Man: No Way Home has the highest popularity rate in our dataset and it has genres of Axtion, Adventure and Science Fiction

Q4: Which movie got the lowest popularity? what is its genre? The united states, thread has the highest lowest rate in our dataset and it has genres of music, drama, war, sci-fi and history.

Q5: Which year has the most filmmed movies? year 2020 has the highest filmming rate in our dataset.