

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

In [7]: df = pd.read_csv('mymovieDb.csv',lineterminator = '\n')

In [9]: df.head()

Out [9]:
```

| | Release_Date | Title | Overview | Popularity | Vote_Count | Vote_Average | Original_Language | Genre | Poster_Url |
|---|--------------|-------------------------|---|------------|------------|--------------|-------------------|------------------------------------|--|
| 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.tmbd.org/vp/original/1g0dhy1q4i... |
| 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.tmbd.org/vp/original/74xTEg17R3... |
| 2 | 2022-02-25 | No Exit | Stranded at a rest stop in the mountains durin... | 2618.087 | 122 | 6.3 | en | Thriller | https://image.tmbd.org/vp/original/VD4hLnQWK3... |
| 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.tmbd.org/vp/original/4QPNHkMI5... |
| 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.tmbd.org/vp/original/aaq4Pw5Xeu... |

```
In [11]: 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: 631.1+ KB

In [17]: df['Genre'].head()

Out [17]:
```

| | |
|---|------------------------------------|
| 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 [21]: df.duplicated().sum()

Out [21]: 0

In [23]: df.describe()

Out [23]:
```

| | Popularity | Vote_Count | Vote_Average |
|-------|-------------|--------------|--------------|
| count | 9827.000000 | 9827.000000 | 9827.000000 |
| mean | 40.326088 | 1392.805536 | 6.439634 |
| std | 106.873998 | 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% | 36.191500 | 1376.000000 | 7.100000 |
| max | 5083.954000 | 31077.000000 | 10.000000 |

```
In [25]: df['Release_Date'] = pd.to_datetime(df['Release_Date'])
print(df['Release_Date'].dtypes)

datetime64[ns]

In [27]: df['Release_Date'] = df['Release_Date'].dt.year
df['Release_Date'].dtypes

Out [27]: dtype('int32')
```

Dropping the columns

```
In [30]: cols = ['Overview','Original_Language','Poster_Url']

In [32]: df.drop(cols, axis=1, inplace = True)

Out [32]: Index(['Release_Date', 'Title', 'Popularity', 'Vote_Count', 'Vote_Average',
              'Genre'],
              dtype='object')

In [34]: df.head()

Out [34]:
```

| | 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 |

```
In [48]: 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

In [50]: labels = ['not_popular', 'below_avg', 'average', 'popular']
categorize_col(df, 'Vote_Average', labels)
df['Vote_Average'].unique()

Out [50]: ['popular', 'below_avg', 'average', 'not_popular', NaN]
Categories (4, object): ['not_popular' < 'below_avg' < 'average' < 'popular']

In [52]: df.head()

Out [52]:
```

| | 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 [54]: df['Vote_Average'].value_counts()

Out [54]: Vote_Average
not_popular    2467
popular        2450
average        2412
below_avg      2398
Name: count, dtype: int64

In [56]: df.dropna(inplace = True)
df.isna().sum()

Out [56]: Release_Date    0
Title                0
Popularity           0
Vote_Count           0
Vote_Average         0
Genre                0
dtype: int64

In [58]: df.head()

Out [58]:
```

| | 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 |

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

```
In [61]: df['Genre'] = df['Genre'].str.split(',')
df = df.explode('Genre').reset_index(drop = True)
df.head()

Out [61]:
```

| | 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 [63]: # casting column into category
df['Genre'] = df['Genre'].astype('category')
df['Genre'].dtypes

Out [63]: 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 [65]: 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 [67]: df.nunique()

Out [67]: Release_Date    100
Title                9415
Popularity           8088
Vote_Count           3265
Vote_Average         4
Genre                19
dtype: int64

In [69]: df.head()

Out [69]:
```

| | 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 |

Data Visualization

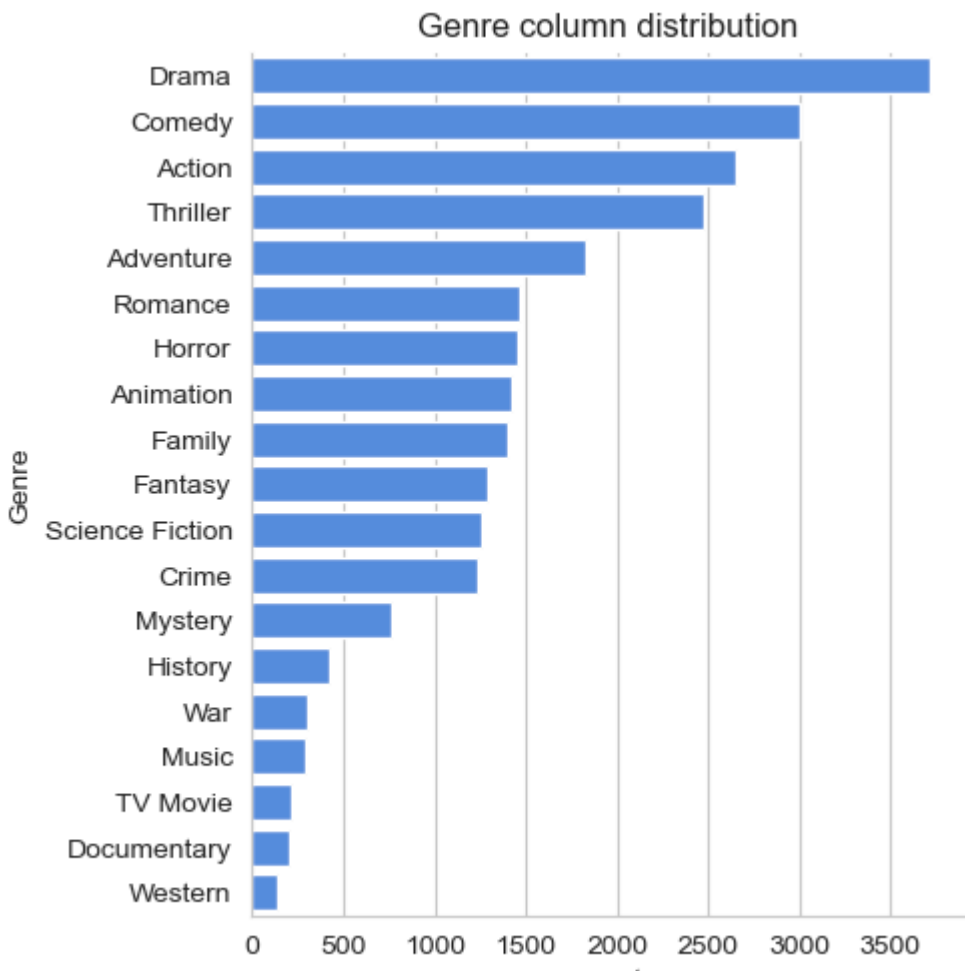
```
In [72]: sns.set_style('whitegrid')

what is the most frequent genre of movies released on Netflix?
```

```
In [75]: df['Genre'].describe()

Out [75]: count      25552
unique        19
top      Drama
freq         3715
Name: Genre, dtype: object

In [79]: sns.catplot(y = 'Genre', data = df, kind = 'count',
                    order = df['Genre'].value_counts().index,
                    color = '#4287f5')
plt.title('Genre column distribution')
plt.show()
```



Which has highest votes in vote avg column?

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

Out [82]:
```

| | 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 [88]: sns.catplot(y = 'Vote_Average', data = df, kind = 'count',
                    order = df['Vote_Average'].value_counts().index,
                    color = '#4287f5')
plt.title('Votes distribution')
plt.show()
```



What movie got the highest popularity? What's its genre

```
In [92]: df.head(2)

Out [92]:
```

| | 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 [94]: df[df['Popularity'] == df['Popularity'].max()]

Out [94]:
```

| | 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 |

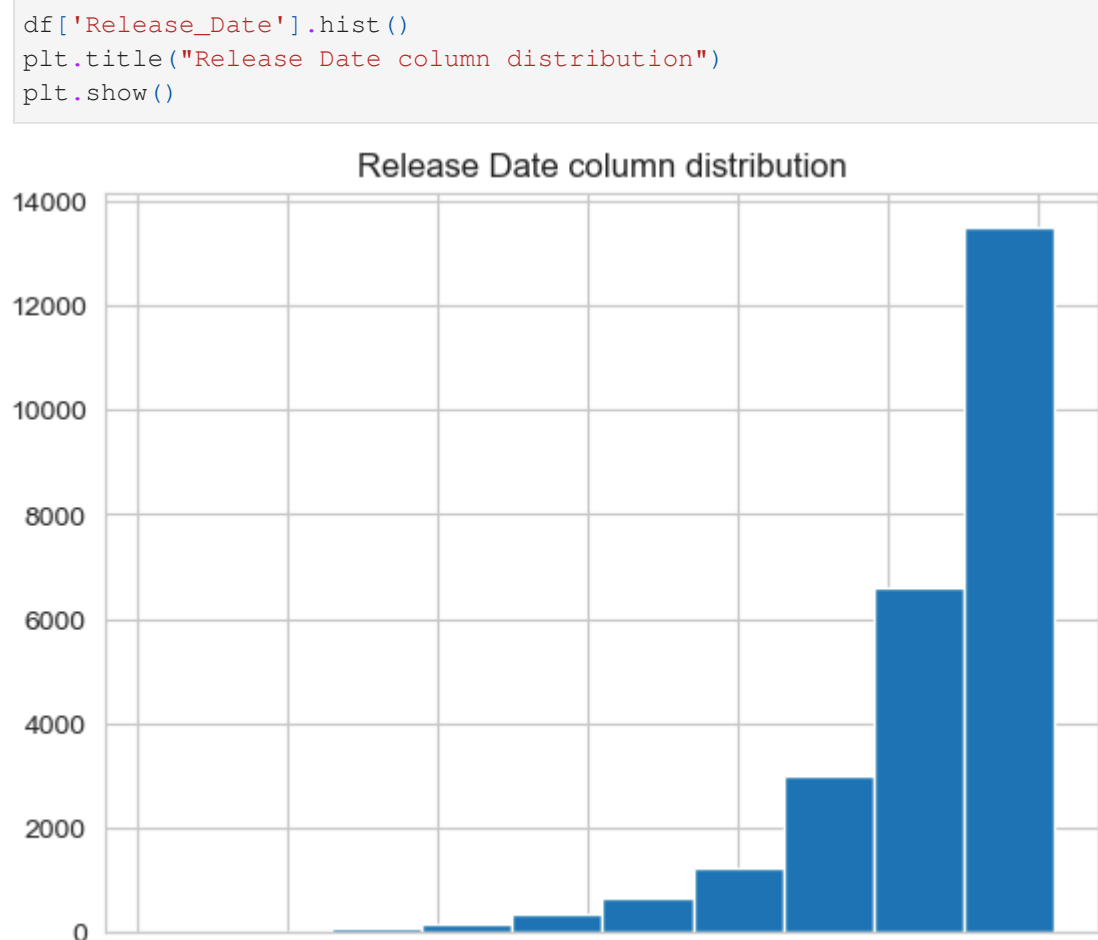
```
In [96]: df[df['Popularity'] == df['Popularity'].min()]

Out [96]:
```

| | 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 | 1964 | Threads | 13.354 | 186 | popular | War |
| 25550 | 1984 | Threads | 13.354 | 186 | popular | Drama |
| 25551 | 1984 | Threads | 13.354 | 186 | popular | Science Fiction |

which year has the most filmed movies?

```
In [99]: df['Release_Date'].hist()
plt.title('Release Date column distribution')
plt.show()
```



```
In [ ]: Conclusion
Q1: What is the most frequent genre in the dataset?
Drama genre is the most frequent genre in our dataset and has appeared more than 14% of the times among 19 other genres.
Q2: What Genres has highest votes?
Q3: What movie get the highest popularity? what's it's genre?
Q4: Which Year has the
```