

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

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

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

```
Out[ ]: Release_Date Title Overview Popularity Vote_Count Vote_Average Original_
```

0	2021-12-15	Spider-Man: No Way Home	Peter Parker is unmasked and no longer able to...	5083.954	8940	8.3
---	------------	-------------------------	---	----------	------	-----

1	2022-03-01	The Batman	In his second year of fighting crime, Batman u...	3827.658	1151	8.1
---	------------	------------	---	----------	------	-----

2	2022-02-25	No Exit	Stranded at a rest stop in the mountains durin...	2618.087	122	6.3
---	------------	---------	---	----------	-----	-----

3	2021-11-24	Encanto	The tale of an extraordinary family, the Madri...	2402.201	5076	7.7
---	------------	---------	---	----------	------	-----

4	2021-12-22	The King's Man	As a collection of history's worst tyrants and...	1895.511	1793	7.0
---	------------	----------------	---	----------	------	-----



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

```
Out[ ]: 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 [ ]: df.duplicated().sum()
```

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

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

```
Out[ ]:
```

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

### # Explain Summary

We have a dataframe consisting of 9827 rows and 9 columns.

**Our dataset looks bit tidy with no NaNs nor duplicated values.**

**Release\_Date** column needs to be casted into datetime and to extract only the year value.

**Overview, Original\_language** and **poster-url** wouldn't be so useful during analysis, so we'll drop them.

There are noticeable outliers in Popularity column.

Vote\_Average better be categorised for proper analysis.

Genre column has comma separated values and white spaces that need to be handled and casted into category.

### Exploration Summary

## Data Pre-processing Start ....

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

datetime64[ns]

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

Out[ ]: dtype('int32')

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

```
Out[ ]: 
```

	Release_Date	Title	Overview	Popularity	Vote_Count	Vote_Average	Original_
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	



Dropping the columns

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

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

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

```
Out[ ]:
```

	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 would cut Vote\_Average Values and make 4 categories: popular,average ,below\_avg, not\_popular to describe it more categorize fashion.

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

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

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

Out [ ]:

	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 [ ]: `df['Vote_Average'].value_counts()`

Out [ ]: Vote\_Average  
not\_popular 2467  
popular 2450  
average 2412  
below\_avg 2398  
Name: count, dtype: int64

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

Out [ ]: Release\_Date 0  
Title 0  
Popularity 0  
Vote\_Count 0  
Vote\_Average 0  
Genre 0  
dtype: int64

In [ ]: `df.head()`

Out [ ]:

	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 then exploded our dataframe to have only one genre per row for each movie**

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

```
Out [ ]: 
```

	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 [ ]: #Casting column into category
df['Genre']=df['Genre'].astype('category')
df['Genre'].dtypes
```

```
Out [ ]: 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 [ ]: 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 [ ]: df.nunique()
```

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

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

```
Out[ ]:
```

	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 pre-processing Done.

## Data Visualization

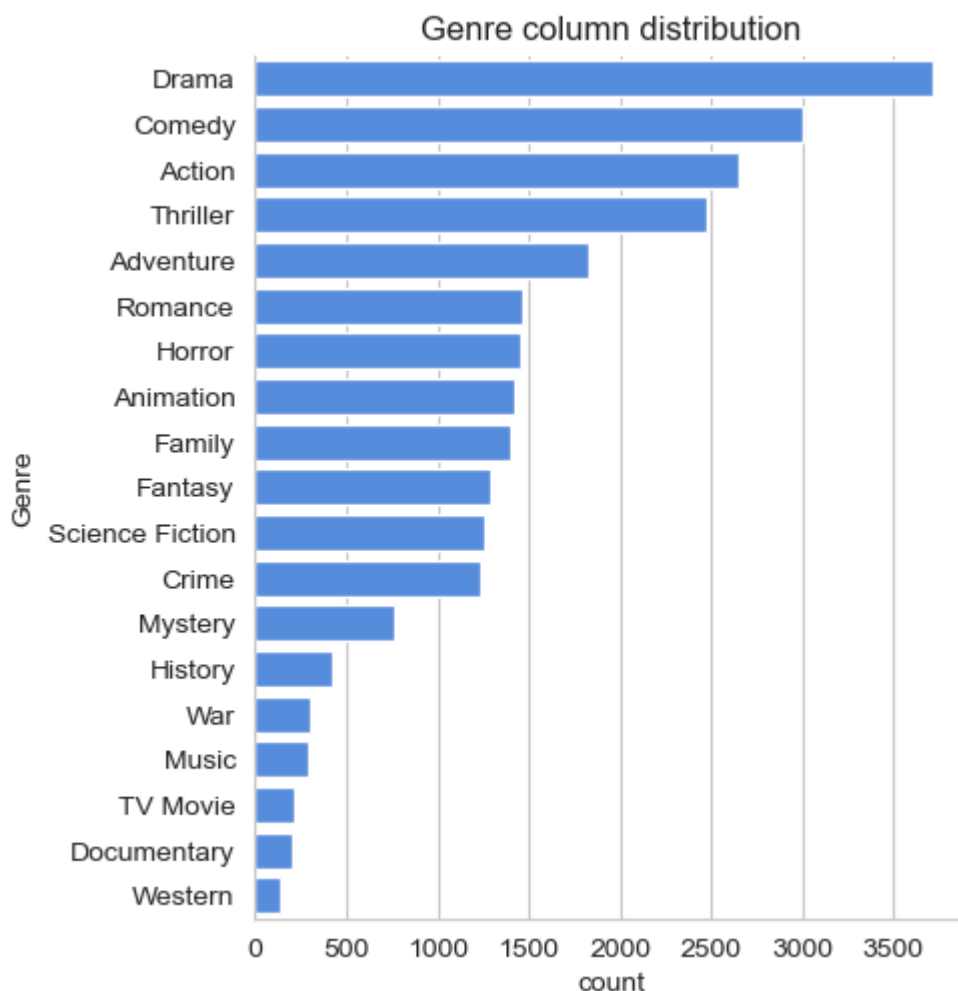
```
In [ ]: sns.set_style('whitegrid')
```

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

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

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

```
In [ ]: 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?

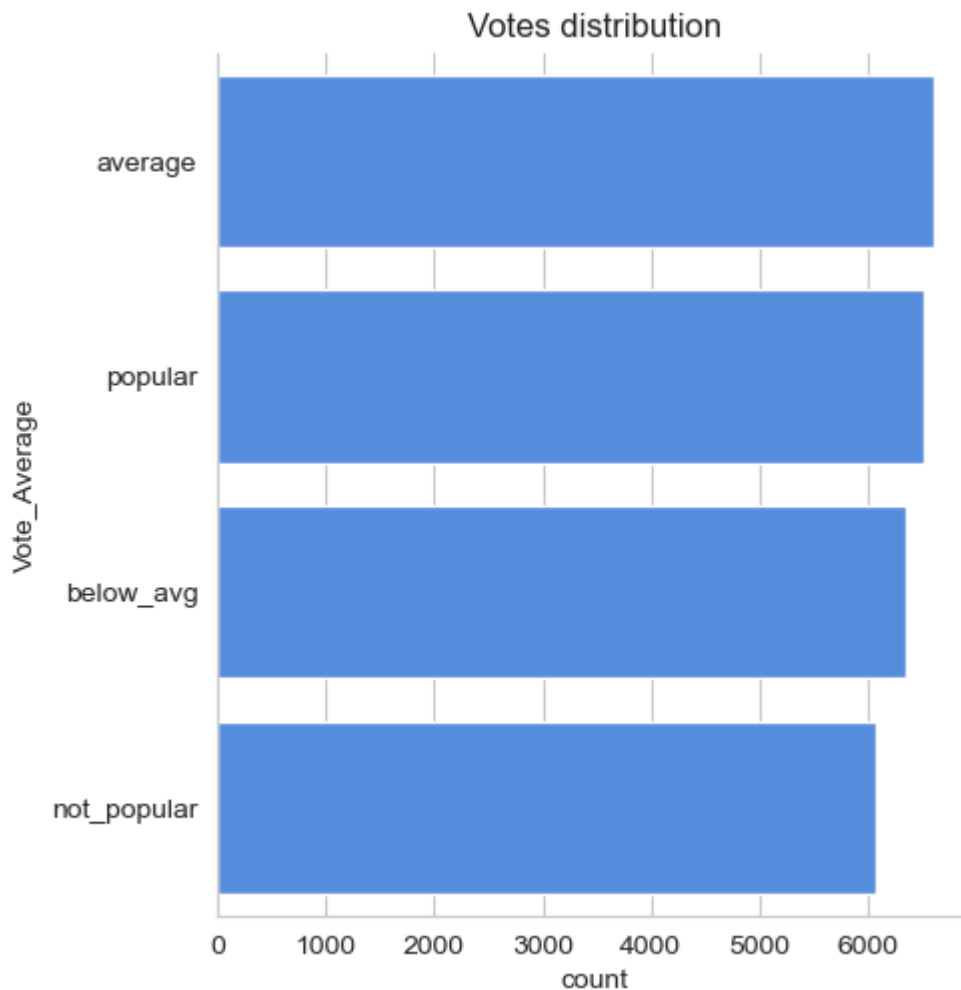
```
In [ ]: df.head()
```

```
Out [ ]:
```

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

Out [ ]:	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 [ ]: df[df['Popularity']==df['Popularity'].max()]
```

Out [ ]:

	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

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

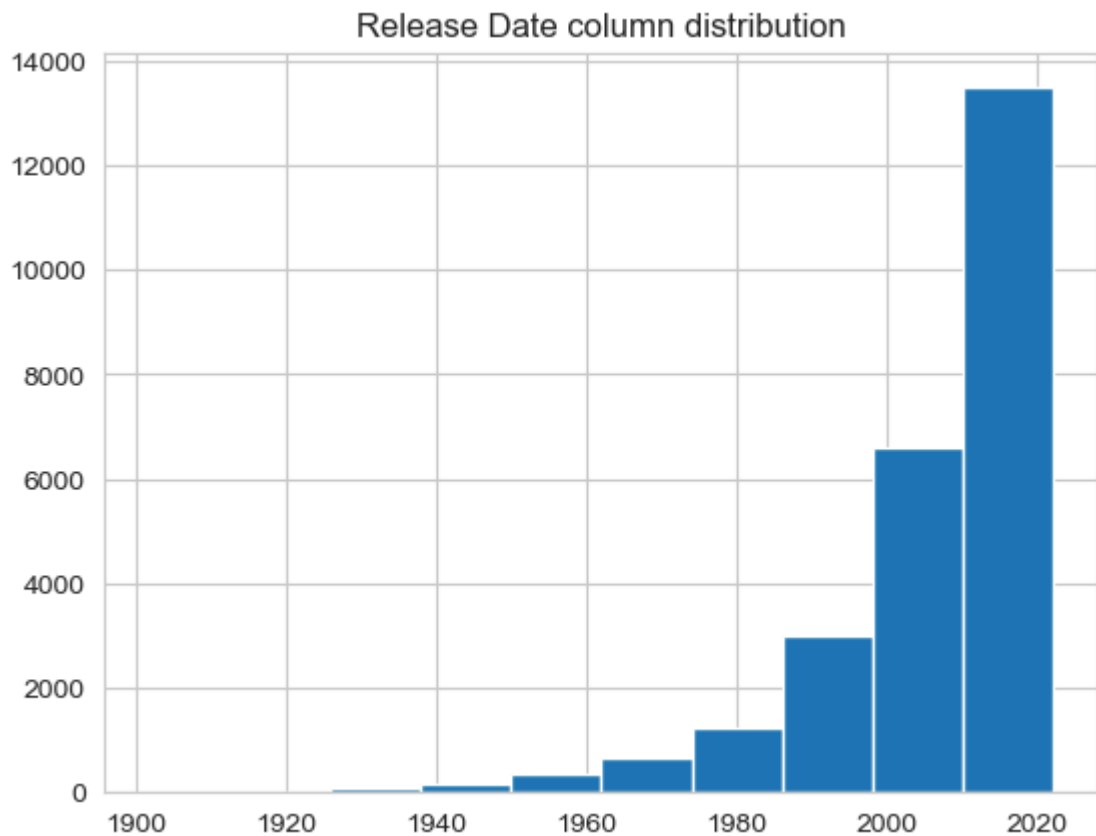
In [ ]: `df[df['Popularity']==df['Popularity'].min()]`

Out [ ]:

	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

## Which year has the most filmed movies?

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



In [ ]: Conclusion

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

Drama genre **is** the most frequent genre **in** our dataest **and** has appeared more than

Q2: Which has highest votes **in** vote avg column?

Average category received the highest votes (**~6700**), showing it's **the most commo**

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

Spider-Man: No Way Home had the highest popularity (**5083.954**), belonging to Acti

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

Movies **with** the lowest popularity (**13.354**) include The United States vs. Billie spanning genres like Music, Drama, War, **and** Sci-Fi.

Q5: Which year has the most filmed movies?

Most movies were filmed **in** the **2020s**, **with** over **14,000** releases recorded.

In [ ]: