

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

In [2]: df= pd.read_csv('mymoviedb.csv', lineterminator='\n')

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

Out[3]:

	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	<a href="https://image.tmdb.org">https://image.tmdb.org</a>
1	2022-03-01	The Batman	In his second year of fighting crime, Batman u...	3827.658	1151	8.1	en	Crime, Mystery, Thriller	<a href="https://image.tmdb.org">https://image.tmdb.org</a>
2	2022-02-25	No Exit	Stranded at a rest stop in the mountains durin...	2618.087	122	6.3	en	Thriller	<a href="https://image.tmdb.org">https://image.tmdb.org</a>
3	2021-11-24	Encanto	The tale of an extraordinary family, the Madri...	2402.201	5076	7.7	en	Animation, Comedy, Family, Fantasy	<a href="https://image.tmdb.org">https://image.tmdb.org</a>
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	<a href="https://image.tmdb.org">https://image.tmdb.org</a>

```
In [4]: 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 [5]: df['Genre'].head()

Out[5]: 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 [6]: df.duplicated().sum()
```

Out[6]: 0

```
In [7]: df.describe()
```

Out[7]:

	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

```
In [8]: # making Release into dae format insted of object
df['Release_Date'] = pd.to_datetime(df['Release_Date'])

print(df['Release_Date'].dtype)
```

datetime64[ns]

```
In [9]: # We need only year
df['Release_Date'] = df['Release_Date'].dt.year

df['Release_Date'].dtype
```

Out[9]: dtype('int64')

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

Out[10]:

	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	<a href="https://image.tmdb.org">https://image.tmdb.org</a>
1	2022	The Batman	In his second year of fighting crime, Batman U...	3827.658	1151	8.1	en	Crime, Mystery, Thriller	<a href="https://image.tmdb.org">https://image.tmdb.org</a>
2	2022	No Exit	Stranded at a rest stop in the mountains durin...	2618.087	122	6.3	en	Thriller	<a href="https://image.tmdb.org">https://image.tmdb.org</a>
3	2021	Encanto	The tale of an extraordinary family, the Madri...	2402.201	5076	7.7	en	Animation, Comedy, Family, Fantasy	<a href="https://image.tmdb.org">https://image.tmdb.org</a>
4	2021	The King's Man	As a collection of history's worst tyrants and...	1895.511	1793	7.0	en	Action, Adventure, Thriller, War	<a href="https://image.tmdb.org">https://image.tmdb.org</a>

```
In [11]: # dropping the columns

cols= ['Overview', 'Original_Language', 'Poster_Url']

df.drop(cols, axis=1 , inplace= True)
```

```
In [12]: df.columns
```

```
Out[12]: Index(['Release_Date', 'Title', 'Popularity', 'Vote_Count', 'Vote_Average',
              'Genre'],
              dtype='object')
```

```
In [13]: df.head()
```

```
Out[13]:
```

	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 [14]: # categorizing vote-Average column
# we would cut the vote_average values and make 4 categories: popular, average ,below-avg, not_popular

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 [15]: labels=['not_popular', 'below_avg', 'average', 'popular']

# calling the fuction categorize_col

categorize_col(df, 'Vote_Average', labels)

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

```
Out[15]: ['popular', 'below_avg', 'average', 'not_popular', NaN]
Categories (4, object): ['not_popular' < 'below_avg' < 'average' < 'popular']
```

```
In [16]: df.head()
```

```
Out[16]:
```

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

```
Out[17]: not_popular    2467
popular      2450
average      2412
below_avg    2398
Name: Vote_Average, dtype: int64
```

```
In [18]: df.dropna(inplace= True)
```

```
In [19]: df.isna().sum()
```

```
Out[19]: Release_Date    0
Title                  0
Popularity             0
Vote_Count             0
Vote_Average           0
Genre                  0
dtype: int64
```

```
In [20]: df.head()
```

```
Out[20]:
```

	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 [21]: # we'd split genere into a list then explod our dataframe to have only one genre per row for each mov

df['Genre']= df['Genre'].str.split(', ')

df= df.explode('Genre').reset_index(drop= True)

df.head()
```

```
Out[21]:
```

	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 [22]: # casting column into category

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

df['Genre'].dtype
```

```
Out[22]: 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)
```

In [23]: 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  int64
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), int64(2), object(1)
memory usage: 849.4+ KB
```

In [24]: df.nunique()

```
Out[24]: Release_Date    100
Title                9415
Popularity           8088
Vote_Count           3265
Vote_Average         4
Genre                19
dtype: int64
```

In [25]: df.head()

```
Out[25]:
```

	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 [26]: # Data Visualization
sns.set_style('whitegrid')
```

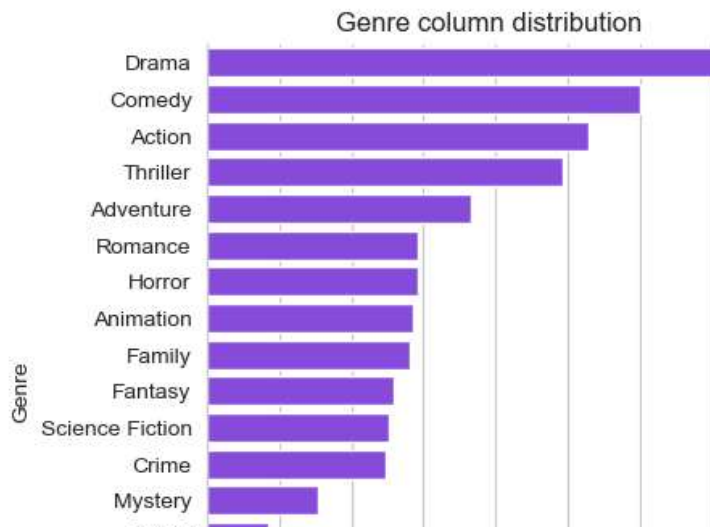
```
In [27]: # 1.What is the most frequent genre of movies released on Netflix?

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

```
Out[27]: count      25552
unique         19
top           Drama
freq          3715
Name: Genre, dtype: object
```

In [31]: *# visualizing genere column*

```
sns.catplot(y= 'Genre', data= df, kind= 'count', order= df['Genre'].value_counts().index, color= '#8532f')
plt.title('Genre column distribution')
plt.show()
```

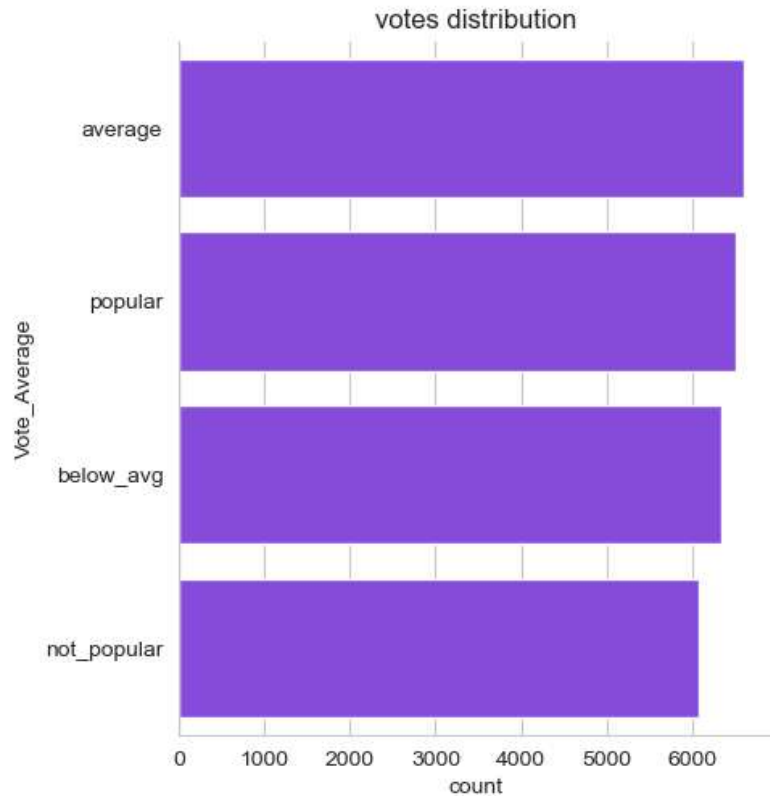


In [32]: *# 2. Which has highest votes in vote avg column?*  
df.head()

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 [34]: sns.catplot(y= 'Vote_Average', data= df, kind= 'count', order= df['Vote_Average'].value_counts().index,
plt.title('votes distribution')
plt.show()
```



```
In [36]: # 3.What movie got the highest popularity? what's its genre?
```

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

Out[36]:

	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 [37]: # 4.What movie got the lowest popularity? what's its genre?
```

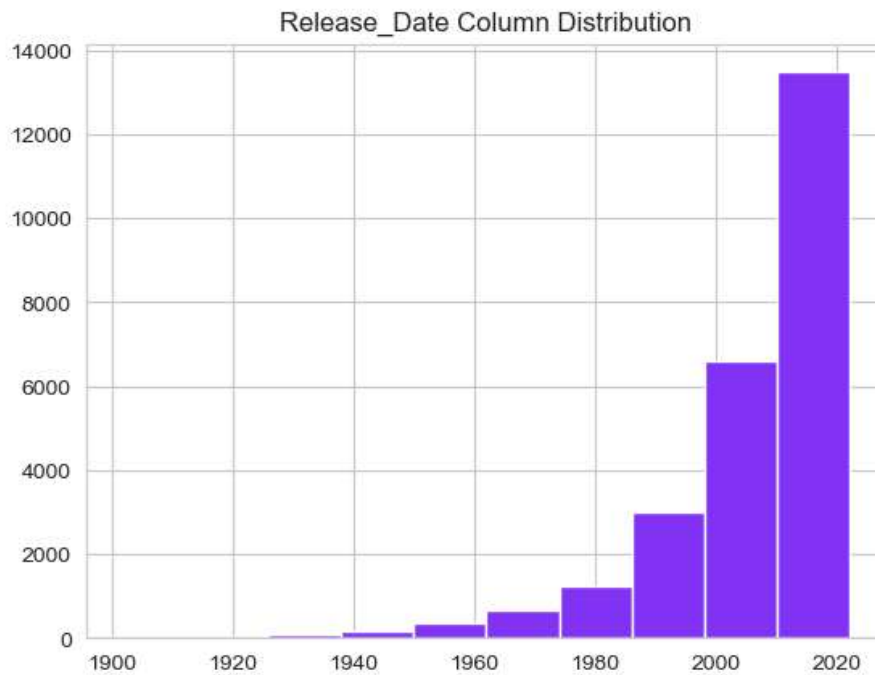
```
df[df['Popularity']== df['Popularity'].min()]
```

Out[37]:

	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

In [39]: *# 5.Which year has the most filmed movies?*

```
df['Release_Date'].hist(color='#8532f9')  
plt.title('Release_Date Column Distribution')  
plt.show()
```



In [ ]: