



Netflix Movies Data Analysis

This project analyzes a movie dataset to extract meaningful insights about popularity, genres, languages, and ratings of Netflix-like content.



Importing Required Libraries

Here we import essential libraries such as NumPy, Pandas for data manipulation, and Matplotlib, Seaborn for visualization.

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



Loading the Dataset

The dataset is loaded from a CSV file named `mymoviedb.csv`.

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



Displaying First Few Records

We display the first 5 rows to understand the structure and contents of the dataset.

```
In [3]: df.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	Ad
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	An

	Release_Date	Title	Overview	Popularity	Vote_Count	Vote_Average	Original_Language	
4	2021-12-22	The King's Man	As a collection of history's worst tyrants and...	1895.511	1793	7.0	en	Ad



Statistical Summary

Summary statistics of numerical columns to understand distribution.

In [4]: `df.describe()`

```
Out[4]:
```

	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



Dataset Info

Basic information such as column data types and non-null counts.

In [5]: `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
```



Genre Analysis

Breaking down or analyzing movie genres to find popular or frequent genres.

```
In [8]: df['Genre'].head()
```

```
Out[8]: 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 [10]: df.duplicated().sum()
```

```
Out[10]: 0
```

Date Conversion

We convert the 'Release_Date' column to proper datetime format.

```
In [12]: df['Release_Date'] = pd.to_datetime(df['Release_Date'])
print(df['Release_Date'].dtypes)
```

```
datetime64[ns]
```

Date Conversion

We convert the 'Release_Date' column to proper datetime format.

```
In [13]: df['Release_Date'] = df['Release_Date'].dt.year
df['Release_Date'].dtypes
```

```
Out[13]: dtype('int64')
```

Displaying First Few Records

We display the first 5 rows to understand the structure and contents of the dataset.

```
In [14]: df.head()
```

```
Out[14]:
```

	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	en	Ad
1	2022	The Batman	In his second year of fighting crime, Batman u...	3827.658	1151	8.1	en	

	Release_Date	Title	Overview	Popularity	Vote_Count	Vote_Average	Original_Language
2	2022	No Exit	Stranded at a rest stop in the mountains durin...	2618.087	122	6.3	en
3	2021	Encanto	The tale of an extraordinary family, the Madri...	2402.201	5076	7.7	en
4	2021	The King's Man	As a collection of history's worst tyrants and...	1895.511	1793	7.0	en

```
In [15]: #dropping the columns
cols = ['Overview', 'Original_Language', 'Poster_Url']
```

```
In [17]: df.drop(cols, axis = 1, inplace = True)
df.columns
```

```
Out[17]: Index(['Release_Date', 'Title', 'Popularity', 'Vote_Count', 'Vote_Average',
        'Genre'],
        dtype='object')
```

👁️ Displaying First Few Records

We display the first 5 rows to understand the structure and contents of the dataset.

```
In [18]: df.head()
```

```
Out[18]:
```

	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 [19]: #categorizing the vote_averager column
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
```

★ Ratings and Popularity

Analyzing which movies received higher ratings and popularity scores.

```
In [23]: labels = ['not_popular', 'below_avg', 'average', 'popular']

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

```
Out[23]: ['popular', 'below_avg', 'average', 'not_popular', NaN]
Categories (4, object): ['not_popular' < 'below_avg' < 'average' < 'popular']
```

👁 Displaying First Few Records

We display the first 5 rows to understand the structure and contents of the dataset.

```
In [24]: df.head()
```

```
Out[24]:
```

	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

★ Ratings and Popularity

Analyzing which movies received higher ratings and popularity scores.

```
In [25]: df['Vote_Average'].value_counts()
```

```
Out[25]: not_popular    2467
popular      2450
average      2412
below_avg    2398
Name: Vote_Average, dtype: int64
```

```
In [26]: df.dropna(inplace = True)
df.isna().sum()
```

```
Out[26]: Release_Date    0
Title                  0
Popularity             0
Vote_Count             0
Vote_Average          0
```

Genre 0
dtype: int64

Displaying First Few Records

We display the first 5 rows to understand the structure and contents of the dataset.

```
In [27]: #splitting Genre into a list and then exploding the dataframe

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

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

Genre Analysis

Breaking down or analyzing movie genres to find popular or frequent genres.

```
In [29]: #casting column into category

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

df['Genre'].dtypes

Out[29]: 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)
```

Dataset Info

Basic information such as column data types and non-null counts.

```
In [30]: 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 [31]: df.nunique()
```

```
Out[31]: Release_Date    100
Title          9415
Popularity     8088
Vote_Count     3265
Vote_Average    4
Genre          19
dtype: int64
```

Data Visualization

Visual representation of data trends using Seaborn and Matplotlib.

```
In [32]: #data visualization
sns.set_style('whitegrid')
```

what is the most frequent genre of movies released on netflix

Genre Analysis

Breaking down or analyzing movie genres to find popular or frequent genres.

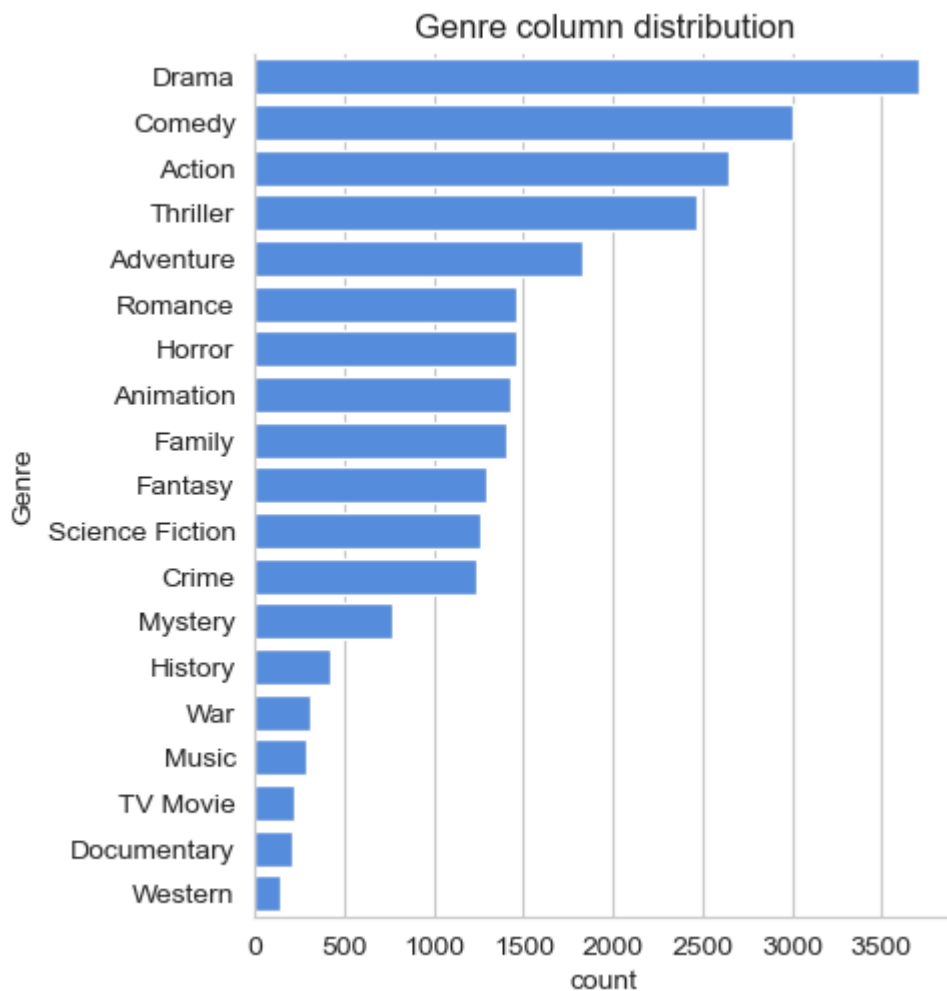
```
In [34]: df['Genre'].describe()
```

```
Out[34]: count      25552
unique         19
top           Drama
freq          3715
Name: Genre, dtype: object
```

Data Visualization

Visual representation of data trends using Seaborn and Matplotlib.

```
In [39]: sns.catplot(y = 'Genre', data = df, kind = 'count',
                    order = df['Genre'].value_counts().index,
                    color = '#4287f5')
plt.title('Genre column distribution')
plt.show()
```

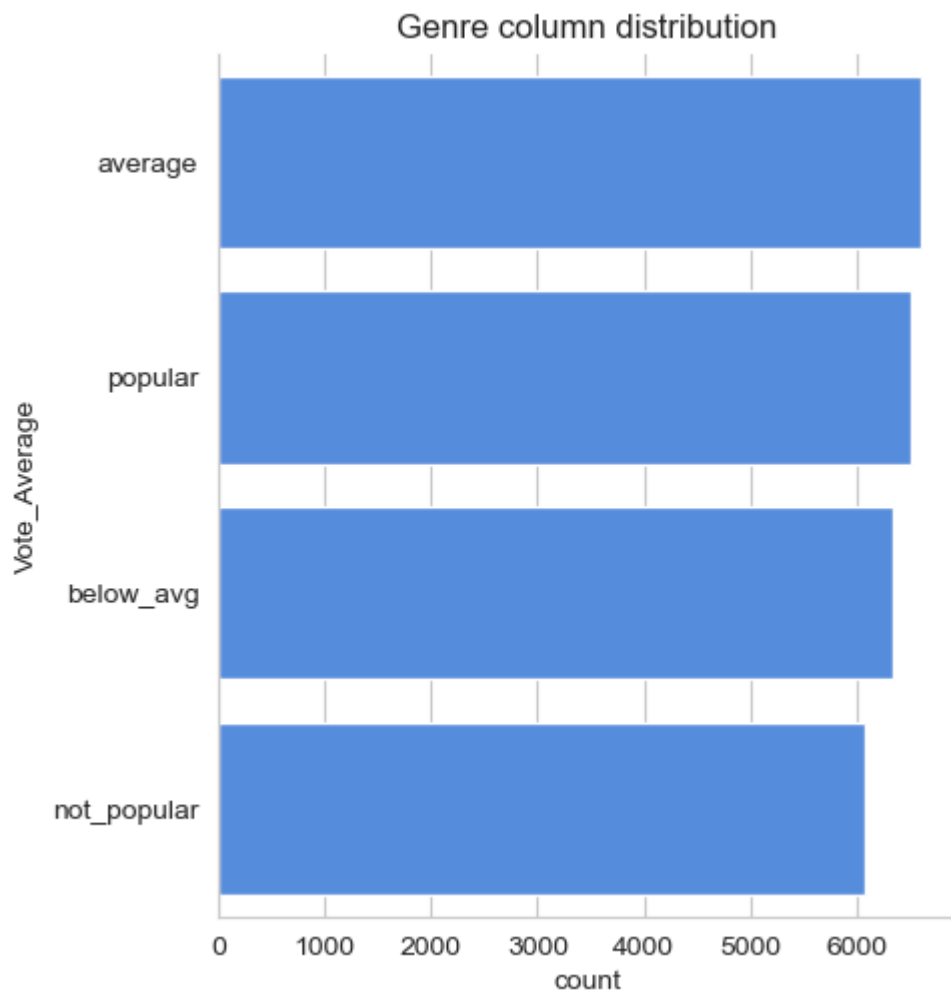


what has heighest votes in vote avg column

Data Visualization

Visual representation of data trends using Seaborn and Matplotlib.

```
In [38]: sns.catplot(y = 'Vote_Average', data = df, kind = 'count',  
                  order = df['Vote_Average'].value_counts().index,  
                  color = '#4287f5')  
plt.title('Genre column distribution')  
plt.show()
```

Which movie got the lowest popularity

✨ Ratings and Popularity

Analyzing which movies received higher ratings and popularity scores.

```
In [41]: df[df['Popularity'] == df['Popularity'].min()]
```

```
Out[41]:
```

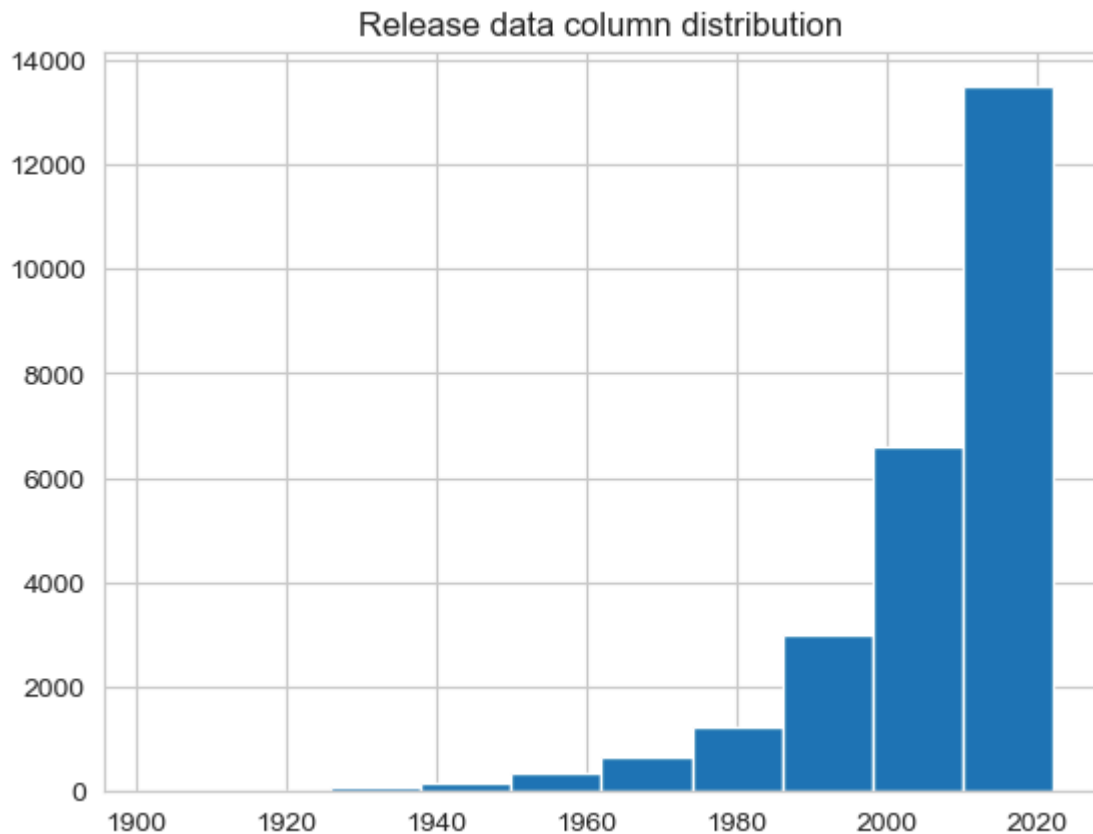
	Release_Date	Title	Popularity	Vote_Count	Vote_Average	Genre
25546	2021	The United States vs. Billie Holiday	13.354	152	average	Music
25547	2021	The United States vs. Billie Holiday	13.354	152	average	Drama
25548	2021	The United States vs. Billie Holiday	13.354	152	average	History
25549	1984	Threads	13.354	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

Data Visualization

Visual representation of data trends using Seaborn and Matplotlib.

```
In [42]: df['Release_Date'].hist()  
plt.title('Release data column distribution')  
plt.show()
```



```
In [ ]:
```

Conclusion

In this project, I explored and analyzed a dataset containing Netflix-like movie data. I performed preprocessing, visualizations, and extracted insights about:

- Most frequent genres
- Trends in movie production over the years
- Distribution of languages
- Ratings and popularity of movies

The dataset allowed me to understand key drivers behind movie success and trends in user preferences.

Problems Solved

- **Data Cleaning:** Converted date formats and handled missing values.
- **Exploratory Data Analysis (EDA):** Used visualization to uncover patterns and trends.
- **Genre Analysis:** Identified the most common genres.
- **Language & Year Distribution:** Analyzed movie counts per language and per year.
- **Popularity & Rating Insights:** Determined which movies were most liked and most viewed.

This project provides a foundation for recommendation systems or business strategy planning based on content performance.

In []: