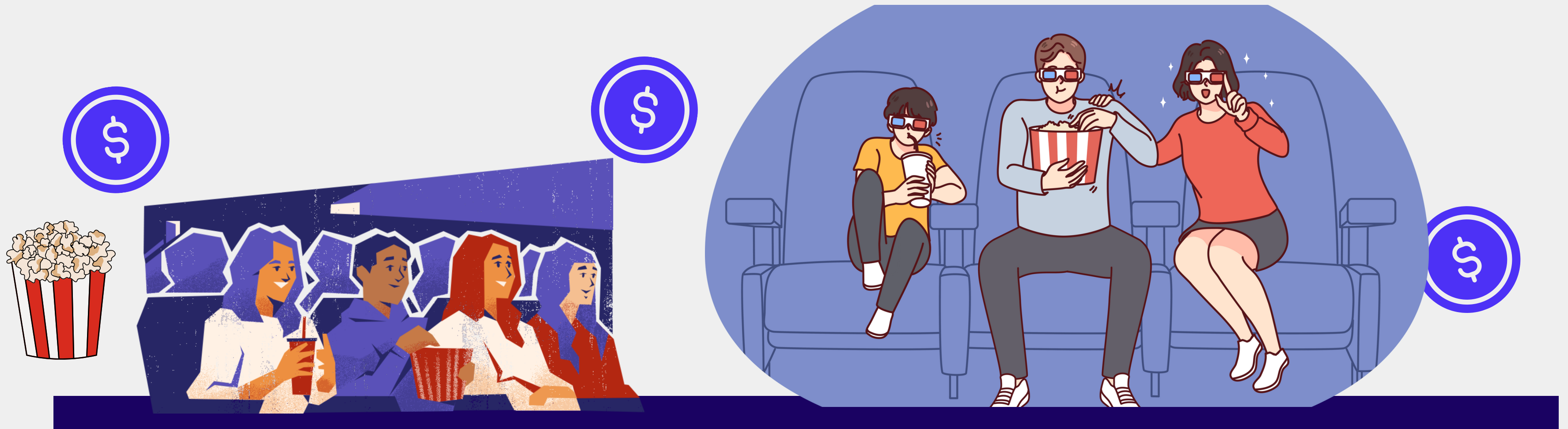# Movie Correlation

# Data Used
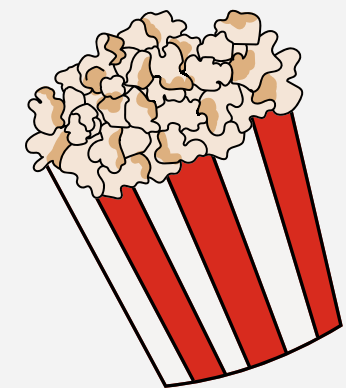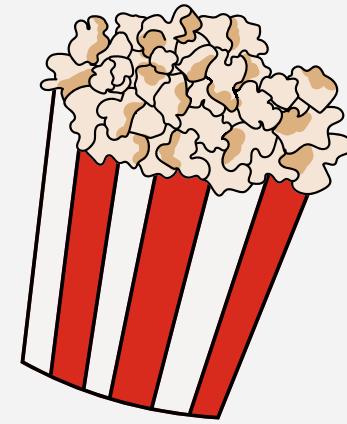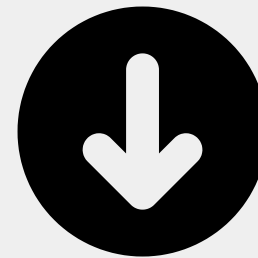
- Movies.csv

# Language

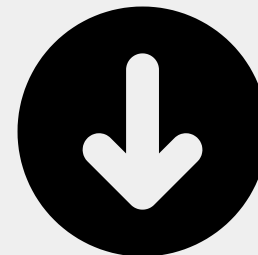- Python

# Tool

- Jupyter Notebook

# Process Used In Project

Getting Data

↓

Cleaning Data (ETL Process)

↓

Data Analysis

↓

Graphs & Conclusion

# Geting Data

Importing all the modules and reading the data

```python
# You can do this all now or as you need them
import pandas as pd
import seaborn as sns

import matplotlib
import matplotlib.pyplot as plt
plt.style.use('ggplot')
from matplotlib.pyplot import figure

%matplotlib inline
matplotlib.rcParams['figure.figsize'] = (12,8) # Adjust yhe configuration of the plot we will create
```

```python
[3]:  # Read in the data
      df = pd.read_csv(r"movies.csv")
```

```python
[3]:  df
```

# Data Cleaning (ETL Process)

Checking the Total Number of rows and columns & Null Value

```
[4]:  # We need to see if we have any missing data
      df.isnull().sum()
```

```
[4]:  name             0
      rating          77
      genre            0
      year             0
      released         2
      score            3
      votes            3
      director         0
      writer           3
      star             1
      country          3
      budget        2171
      gross          189
      company         17
      runtime          4
      dtype: int64
```

```
[5]:  # Volume of data
      df.shape
```

```
[5]:  (7668, 15)
```

# Data Cleaning (ETL Process)

Filling the Null Value of the columns

```python
# Fill The Missing Data in Rating Coulmn using Fill down method
df['rating'] = df['rating'].fillna(method='ffill')
```

```python
# Filling the budget column using mean of the column
df['budget'] = df['budget'].fillna(df['budget'].mean())
```

```python
# Filling the gross column using mean of the column
df['gross'] = df['gross'].fillna(df['gross'].mean())
```

# Data Analysis

```python
# Data Visualization for better understanding
plt.scatter(x=df['budget'],y=df['gross'])
plt.title('Budget vs Gross Earning')
plt.xlabel('Budget of the film')
plt.ylabel('Gross Earning')
plt.show()
```

Finding Relation between budget of the movie and the gross earning of the movie

Correlation between the numeric feature of the data

```python
# Correlation of Numberic feature of movies
correlation_matrix = df.corr(method = 'pearson', numeric_only = True)
sns.heatmap(correlation_matrix, annot=True)
plt.title('Correlation Matrix for Numeric Feature')
plt.xlabel('Movie Feature')
plt.ylabel('Movie Feature')
plt.show()
```

# Data Analysis

```python
# Changing the non numeric data type value of the data to numeric data type
df_numerized = df

for col_name in df_numerized.columns:
    if(df_numerized[col_name].dtype == 'object'):
        df_numerized[col_name] = df_numerized[col_name].astype('category')
        df_numerized[col_name] = df_numerized[col_name].cat.codes

df_numerized
```
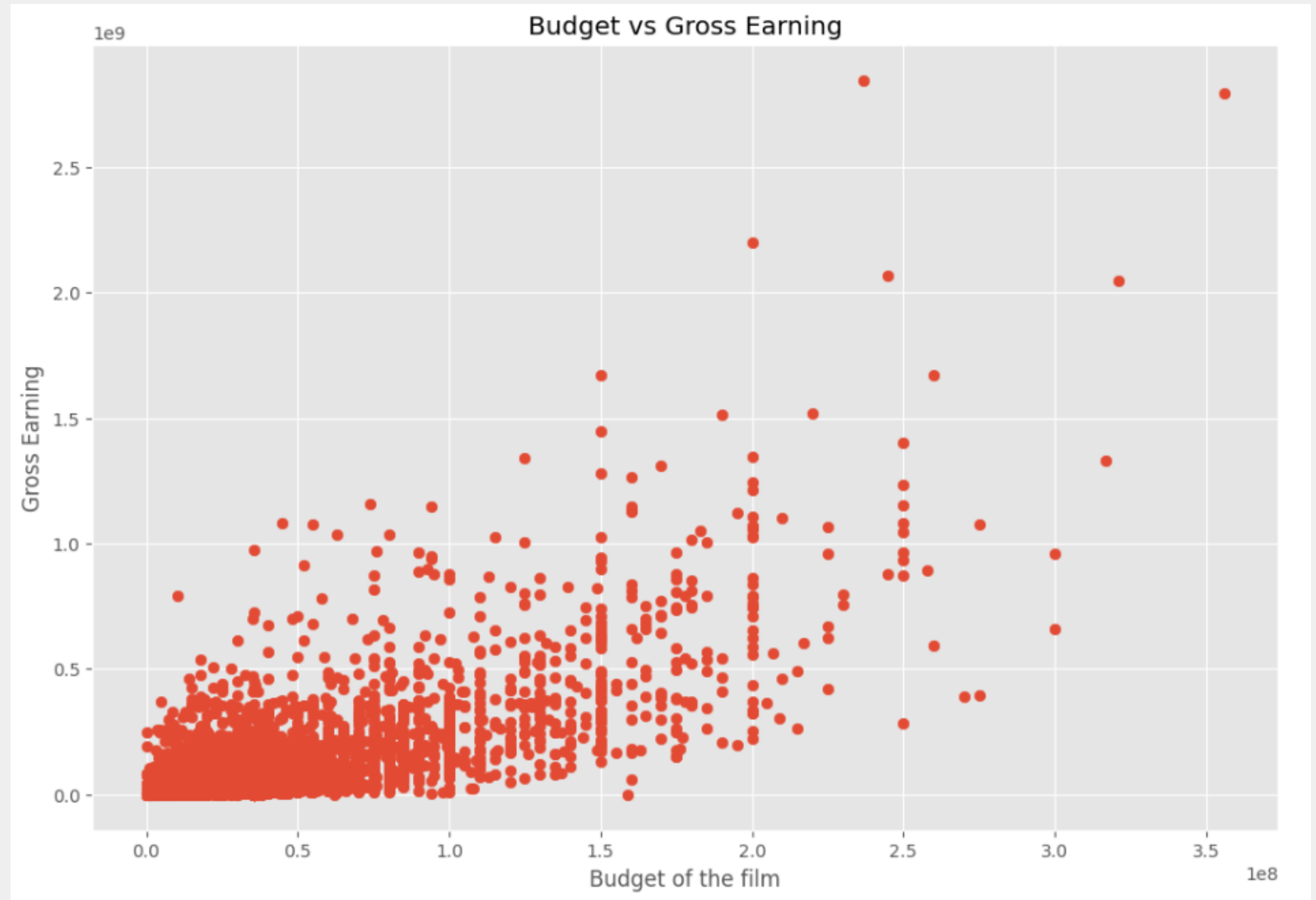
Changing the non numeric value data type of the data to the numeric value data type of the data

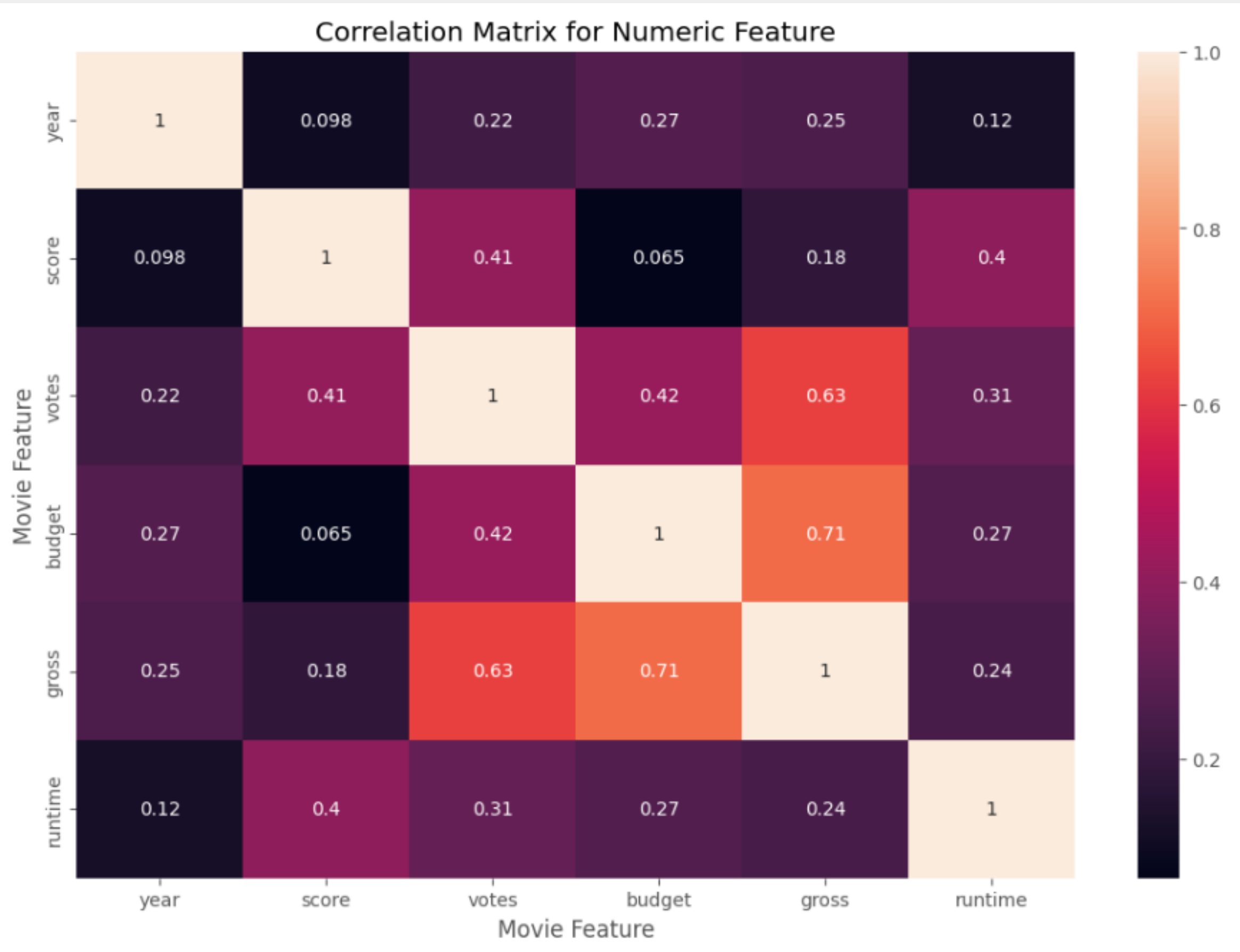Correlation between the non numeric feature of the data

```python
# Correlation of Non Numeric feature
correlation_matrix = df_numerized.corr(method = 'pearson')
sns.heatmap(correlation_matrix, annot=True)
plt.title('Correlation Matrix for Non - Numeric Feature')
plt.xlabel('Movie Feature')
plt.ylabel('Movie Feature')
plt.show()
```

# Graphs

Scatter plot shows relationship between Budget of the movie and the Gross earning of the movie
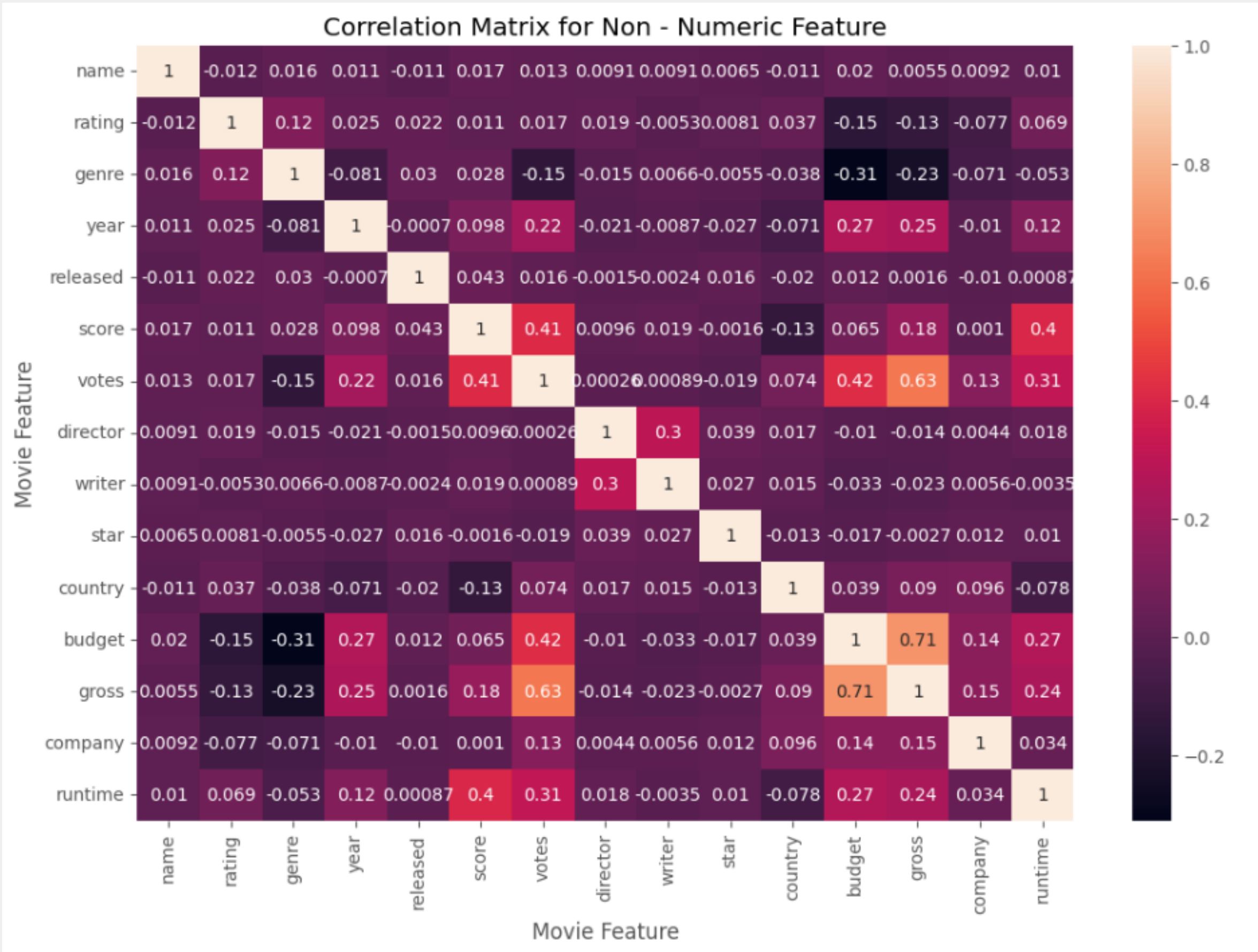


Budget vs Gross Earning

# Graphs



Correlation Matrix for Numeric Feature

The Heatmap here shows the correlation between the numeric feature of the movie

# Graphs

The Heatmap here shows the correlation between the non numeric feature of the movie

# Conclusion

The Conclusion of the project is that the Votes and Gross Earning of the movies has the highest correlation.

# Thank you for your time.

## CONTACT

Feel free to reach out:

Email: tiwariv617@gmail.com
LinkedIn: www.linkedin.com/in/vishal-tiwari-6a380a2aa
Portifolio: https://vishaltiwari1.github.io/Portfolio/