# Experiment report: Data Visualization

Done by student: **Bylkova Kristina (伯汀娜), 1820249049**
Course: **Big Data Analysis Technology**
Date: **October 2024**

## Task 1

**Types of data analysis and what kind of Data Visualization tools can support which type of data analysis?**

1. Descriptive analysis: deals with what happened in the Past (Microsoft Power BI, Tableau, Matplotlib & Seaborn);

2. Diagnostic analysis: deals with why it happened in the Past (Tableau, Power BI, D3.js, Matplotlib & Seaborn);

3. Predictive analysis: deals with what will happen in the Future (Tableau, Power BI, Plotly, Matplotlib & Seaborn);

4. Perspective analysis: deals with what step to take to avoid a problem in the future (Power BI, Tableau, D3.js, Plotly).

## Task 2

**Data Visualization tools:**

1. Microsoft Power BI: provides real-time data insights through dashboards;

2. Tableau: a tool for creating interactive data visualizations;

3. D3.js (Data-Driven Documents): an open-source library written in javascript, creates dashboards in web browsers;

4. Plotly: a graphing library used for creating interactive, web-based data visualizations, including charts, graphs, and dashboards;

5. Matplotlib & Seaborn: Matplotlib is a Python library for creating visualizations and Seaborn is a Python data visualization library built on top of Matplotlib (it also makes visualizations).
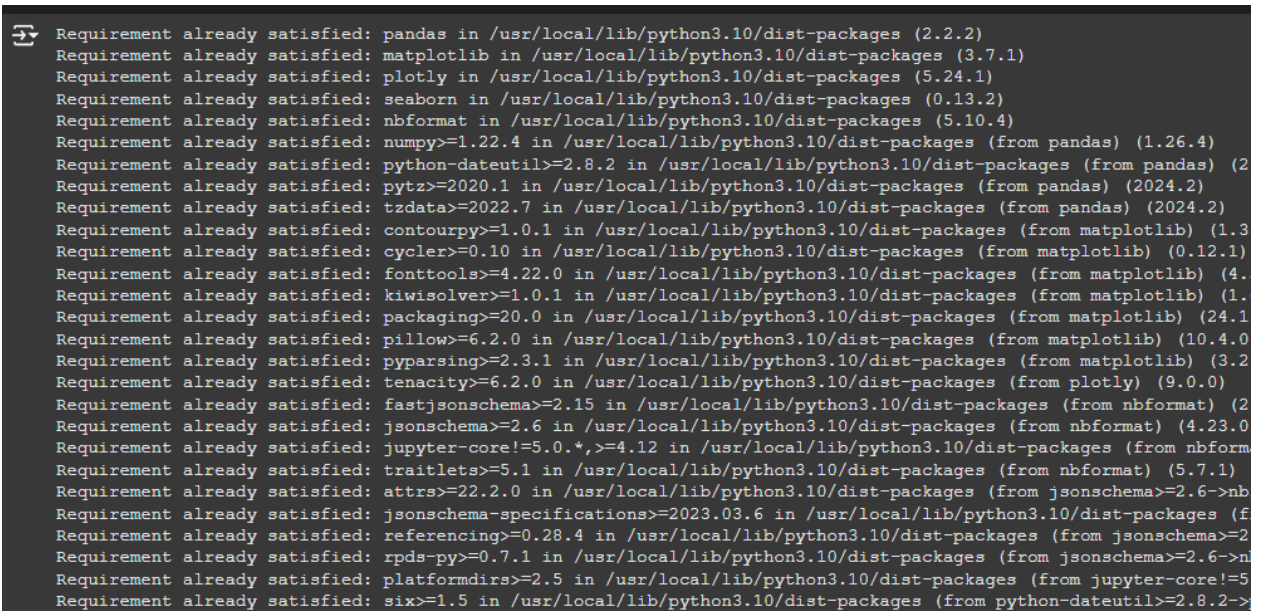
## Experiment 1: Statistical Analysis

Here we dwelve in to statistical information, where we filter data through an attribute.

**Code:**

```python
1  # Importing all necessary libraries for data manipulation, analysis and visualization
2  import sys
3  !{sys.executable} -m pip install pandas matplotlib plotly seaborn nbformat
4
5  import plotly.express as px
6  import plotly.graph_objects as go
7  import pandas as pd
8  import sqlite3
9  import matplotlib.pyplot as plt
10 import numpy as np
11 import seaborn as sns
12
13 # Creating a connection to the database and load the data
14 conn = sqlite3.connect('imdb_games.db')
15 df_cleaned = pd.read_sql_query("SELECT * FROM games", conn)
```

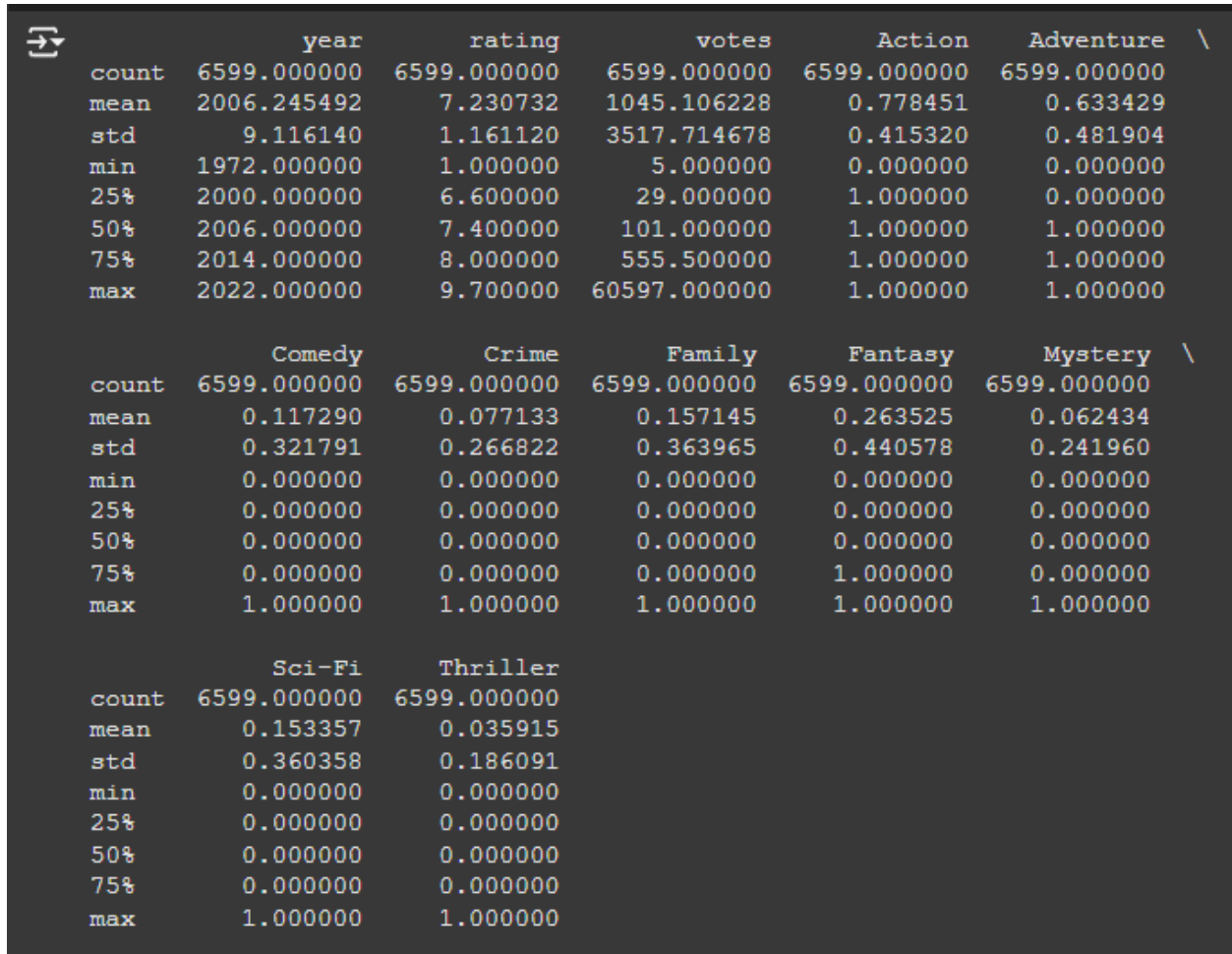Result:



Figure 1: Importing libraries

## Experiment 1.1: Descriptive Statistics

Getting the descriptive statistics, that is the Count, Mean, Standard Deviation, Percentiles (25%, 50%, 75%), minimum and maximum vaues of each column of the cleaned data from the dataset.

**Code:**

```
1  summary_stats = df_cleaned.describe()
2  print(summary_stats)
```

**Result:**

```
                   year         rating          votes         Action      Adventure  \
count      6599.000000    6599.000000    6599.000000    6599.000000    6599.000000
mean       2006.245492       7.230732    1045.106228       0.778451       0.633429
std           9.116140       1.161120    3517.714678       0.415320       0.481904
min        1972.000000       1.000000       5.000000       0.000000       0.000000
25%        2000.000000       6.600000      29.000000       1.000000       0.000000
50%        2006.000000       7.400000     101.000000       1.000000       1.000000
75%        2014.000000       8.000000     555.500000       1.000000       1.000000
max        2022.000000       9.700000   60597.000000       1.000000       1.000000


                 Comedy          Crime         Family        Fantasy        Mystery  \
count      6599.000000    6599.000000    6599.000000    6599.000000    6599.000000
mean          0.117290       0.077133       0.157145       0.263525       0.062434
std           0.321791       0.266822       0.363965       0.440578       0.241960
min           0.000000       0.000000       0.000000       0.000000       0.000000
25%           0.000000       0.000000       0.000000       0.000000       0.000000
50%           0.000000       0.000000       0.000000       0.000000       0.000000
75%           0.000000       0.000000       0.000000       1.000000       0.000000
max           1.000000       1.000000       1.000000       1.000000       1.000000


                 Sci-Fi       Thriller
count      6599.000000    6599.000000
mean          0.153357       0.035915
std           0.360358       0.186091
min           0.000000       0.000000
25%           0.000000       0.000000
50%           0.000000       0.000000
75%           0.000000       0.000000
max           1.000000       1.000000
```

Figure 2: Descriptive Statistics

## Experiment 1.2: Top 11 games by ratings

Discovering the top 11 games by ratings from the cleaned dataset.

**Code:**

```
1  conn = sqlite3.connect('imdb_games.db')
2
3  query = '''
4  SELECT name, rating
5  FROM games
6  ORDER BY rating DESC
7  LIMIT 11
8  '''
9  top_games = pd.read_sql(query, conn)
10
11 conn.close()
12
13 print(top_games.to_string(index=False))
```

**Result:**

| name | rating |
|---|---|
| Red Dead Redemption II | 9.7 |
| The Last of Us | 9.7 |
| The Witcher 3: Wild Hunt | 9.7 |
| Mass Effect: Legendary Edition | 9.7 |
| The Witcher 3: Wild Hunt – Blood and Wine | 9.7 |
| The Last of Us | 9.7 |
| The Witcher 3: Wild Hunt – Blood and Wine | 9.7 |
| Red Dead Redemption II | 9.7 |
| The Witcher 3: Wild Hunt | 9.7 |
| The Last of Us | 9.7 |
| Mass Effect: Legendary Edition | 9.7 |

Figure 3: Top 11 games by ratings

## Experiment 1.3: Total Votes for Each Game

Discovering the total votes for each game genre.

**Code:**

```python
# Defining the genre columns
genre_columns = ['Action', 'Adventure', 'Comedy', 'Crime', 'Family',
                 'Fantasy', 'Mystery', 'Sci-Fi', 'Thriller']

melted_df = df_cleaned.melt(id_vars=['votes'], value_vars=genre_columns,
                            var_name='Genre', value_name='Is_Genre')

# Filtering only applicable genres
filtered_df = melted_df[melted_df['Is_Genre'] == 1]

# Group by Genre and sum votes
total_votes_by_genre = filtered_df.groupby('Genre')['votes'].sum().reset_index()

# Sort by total votes in descending order
total_votes_by_genre = total_votes_by_genre.sort_values(by='votes', ascending=False)

print(total_votes_by_genre)
```

**Result:**



|   | Genre | votes |
|---|---|---|
| 0 | Action | 6348811 |
| 1 | Adventure | 4977777 |
| 5 | Fantasy | 1631520 |
| 3 | Crime | 1343081 |
| 7 | Sci-Fi | 982141 |
| 2 | Comedy | 457004 |
| 8 | Thriller | 440125 |
| 6 | Mystery | 429245 |
| 4 | Family | 368054 |

Figure 4: Total votes for each genre

## Experiment 2: Visualization

## Experiment 2.1: Distribution of ratings

Getting the distribution of ratings with the use of a histogram.

**Code:**

```python
mean_rating = df_cleaned['rating'].mean()

fig = px.histogram(df_cleaned, x='rating', nbins=90,
                   title='Distribution of Video Game Ratings',
                   labels={'rating': 'Rating'},
                   color_discrete_sequence=['purple'])

fig.add_vline(x=mean_rating, line_color='red', line_dash='dash',
              annotation_text=f'Mean: {mean_rating:.2f}',
              annotation_position='top right')

fig.update_layout(xaxis_title='Rating', yaxis_title='Frequency')

fig.show()
```
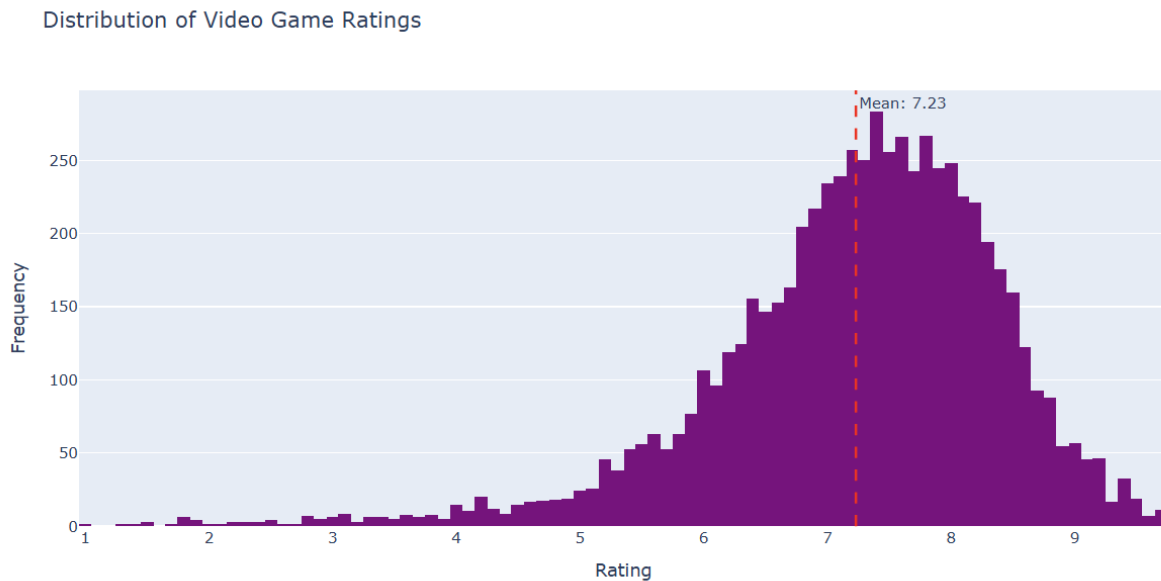
**Result:**



Figure 5: Distribution of ratings

## Experiment 2.2: Ratings over the years with trendline

Getting the ratings over the years using a line plot aided with a trendline.

**Code:**

```
1  average_ratings = df_cleaned.groupby('year')['rating'].mean().reset_index()
2
3  fig = px.line(average_ratings, x='year', y='rating',
4                title='Average Ratings of Video Games Over the Years',
5                labels={'year': 'Year', 'rating': 'Average Rating'},
6                markers=True)
7
8  trendline = go.Scatter(x=average_ratings['year'], y=average_ratings['rating'],
9                          mode='lines',
10                         name='Trendline',
11                         line=dict(color='red'))
12
13 fig.add_trace(trendline)
14
15 fig.show()
```
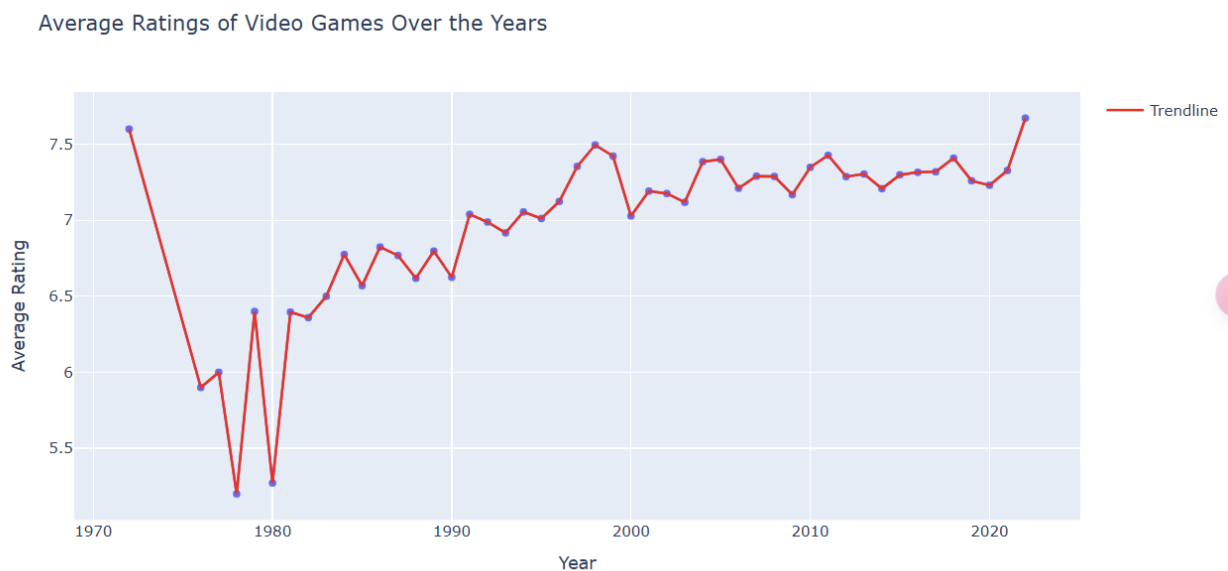
**Result:**



Figure 6: Ratings over the years with trendline

**Experiment 3: Optional Experiment**

**Experiment 3.1: How are certificates related to different genres?**

Exploring the relationship between game certificates and all the genres using a grouped bar chart to visualize.

**Code:**

```python
melted_df = df_cleaned.melt(id_vars=['certificate', 'rating'], value_vars=genre_columns,
                            var_name='Genre', value_name='Is_Genre')

filtered_df = melted_df[melted_df['Is_Genre'] == 1]

# Grouping by certificate and genre to calculate the average rating
cert_genre_relation = filtered_df.groupby(['certificate', 'Genre'])['rating'].mean().reset_index()

fig = px.bar(cert_genre_relation, x='certificate', y='rating', color='Genre',
             title='Average Genre Ratings by Certificate',
             labels={'certificate': 'Certificate', 'rating': 'Average Rating'},
             barmode='group')

fig.update_layout(xaxis_title='Certificate',
                  yaxis_title='Average Rating',
                  legend_title='Genre')

fig.show()
```
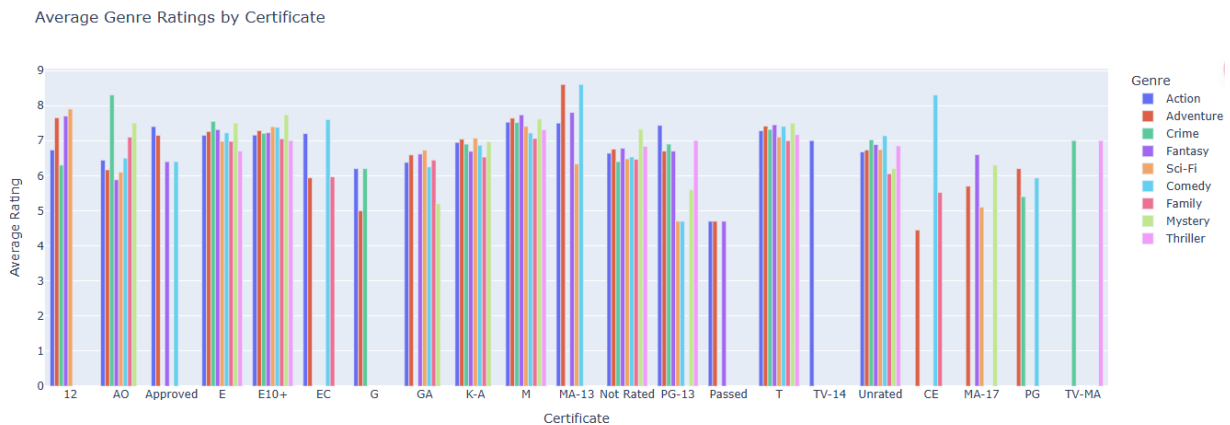
**Result:**



Figure 7: Relation of certificates

**Experiment 3.2: What is the relationship between the Genres, ratings, release year and votes?**

Making a correlation heatmap.

**Code:**

```
1  numeric_df = df_cleaned.select_dtypes(include=[np.number])
2
3  plt.figure(figsize=(10, 8))
4  correlation_matrix = numeric_df.corr()   # Calculating correlation on numeric DataFrame
5  sns.heatmap(correlation_matrix, annot=True, fmt=".2f", cmap='coolwarm')
6  plt.title('Correlation Heatmap')
7  plt.show()
```
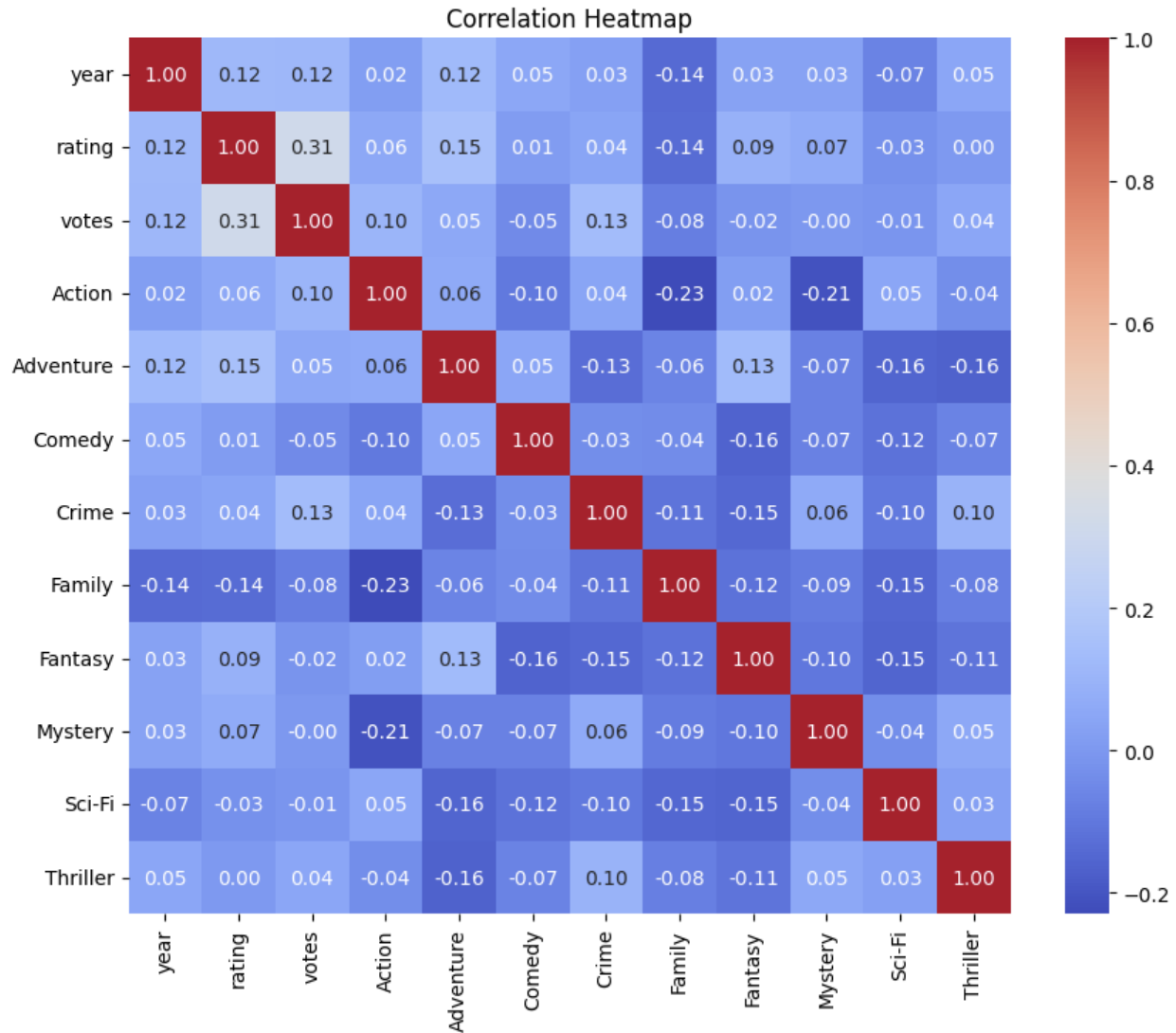
**Result:**



Figure 8: Correlation Heatmap