

# Movie Data Analysis



# Table of contents

01

Members

02

Project Overview

03

Data Overview

04

Data Preprocessing

05

Questions

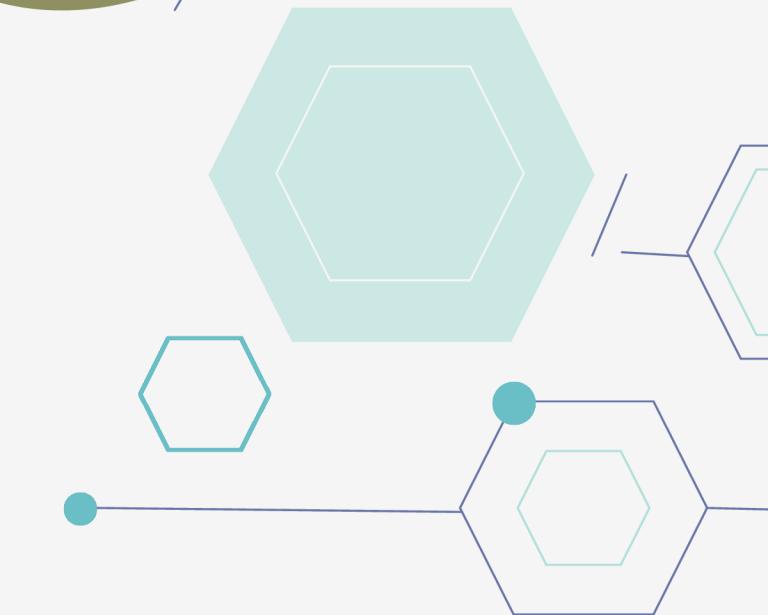
06

Modelling

# 01

# Members

# of Group



# Members of group

22127148

Dương Nhật Huy

22127224

Trương Thuận Kiệt

22127257

Phạm Minh Mẫn

22127492

Hồ Đăng Phúc

# 02

# Project

# Overview



# Reasons for Project



To explore and analyze movie data to uncover insights into factors affecting movie performance, including ratings, revenue, and critical reception.



To address missing data challenges and demonstrate effective preprocessing and modeling techniques.



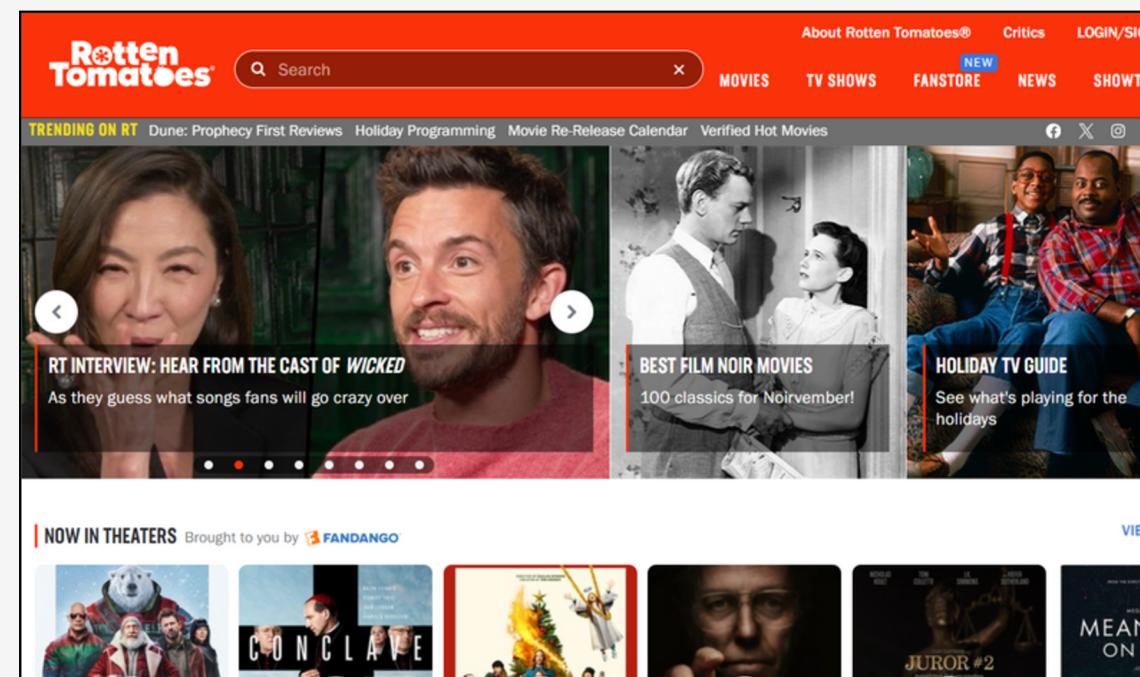
To enhance understanding of data cleaning, visualization, and modeling in a real-world dataset.



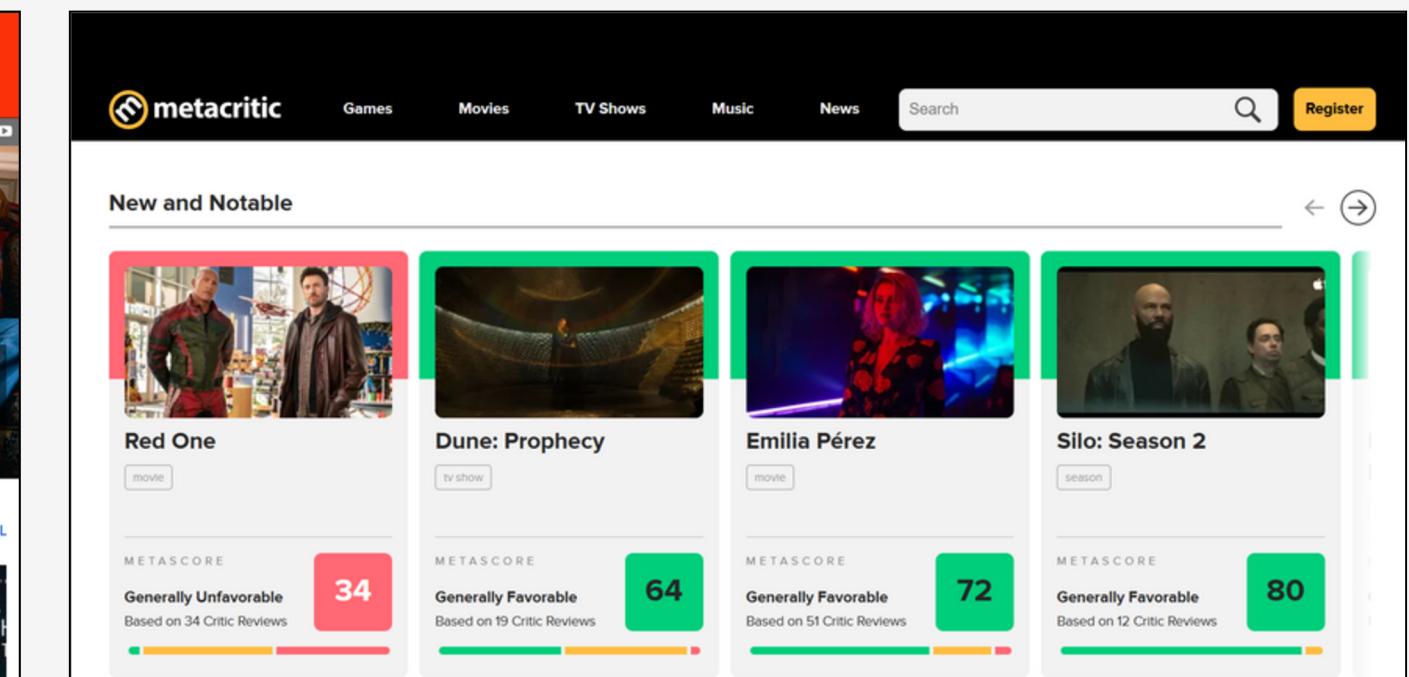
# Data Sources

## Data Rating Web

### 01 Rotten Tomatoes



### 02 Meta Critic



# Data Sources

Budget/ Revenue  
and Inflation rate

## 01 The Numbers

**INSIGHTFUL.  
INTELLIGENT.  
INVALUABLE.** **THE NUMBERS BUSINESS REPORT**

Heretic weekly (November 8, 2024)

Daily Domestic Chart for Thursday November 14, 2024

Movie Title	Gross	%YD	%LW	Theaters	Per Theater	Total Gross	Days In Release
- P Red One	\$3,700,000					\$3,700,000	
1 (1) Venom: The Last Dance	\$800,104	-14%	-51%	3,905	\$205	\$120,244,619	21
2 (2) Heretic	\$789,652	-11%		3,221	\$245	\$15,269,794	7

See the full daily chart

Weekend predictions: Red One set to kick off the Holidays with moderate opening

November 15, 2024



Quick Links

- DEG Watched at Home Top 20
- Netflix Daily Top 10
- Weekly DVD+Blu-ray Chart
- News
- Release Schedule
- Daily Box Office
- Weekend Box Office
- Weekly Box Office
- Annual Box Office
- Box Office Records
- International Box Office
- Distributors
- People Records
- People Index
- Genre Tracking
- Keyword Tracking
- Franchises
- Research Tools
- Bankability Index

Most Anticipated Movies

- Trap
- Avengers: Doomsday
- Borderlands
- Never Let Go
- My Penguin Friend
- The Strangers: Chapter 2
- Venom: The Last Dance
- Piece by Piece
- Azrael
- The Ark and the Aardvark

Investopedia INVESTING SIMULATOR BANKING PERSONAL FINANCE ECONOMY NEWS REVIEWS TRADE 

(-1.28%) QQQ - 496.57  -12.12 (-2.38%) \$ DIA - 434.51  -3.19 (-0.73%) VALUG - 624.9500  -6.7100 (-1.06%) EURUSD - 1.05394  +0.00100 (+0.09%) SPY -

SẴN SÀNG CẮT CÁNH:  
BAY THẮNG TỚI  
SINGAPORE!  
Vé một chiều tới Singapore từ  
**1.102.692đ\***  
Đặt ngay

VÌ HÀNH TRÌNH  
ĐÁNG GIÁ  
flyscoot.com



Find the Best Financial Products

5.50%  
Best Savings or Money Market Account Rate  
[Compare Rates →](#)

5.50%  
Best 6-Month CD Rate  
[Compare Rates →](#)

Stocks End Lower as U.S. Indexes



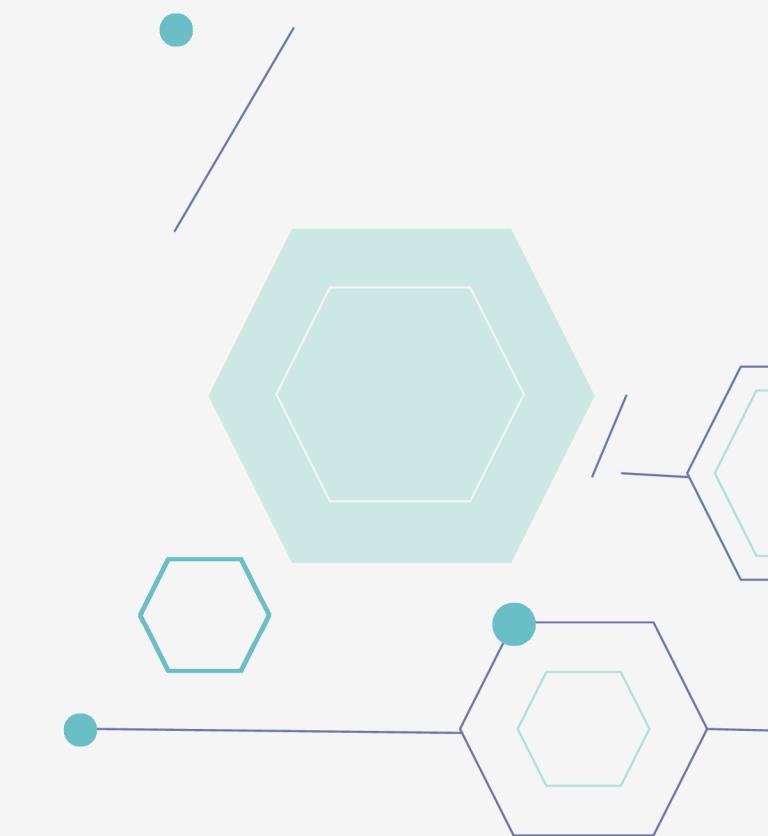
What Wall Street Analysts Think of Walmart's Stock Ahead of Earnings  
Retail giant Walmart is set to report earnings before the market opens Tuesday.



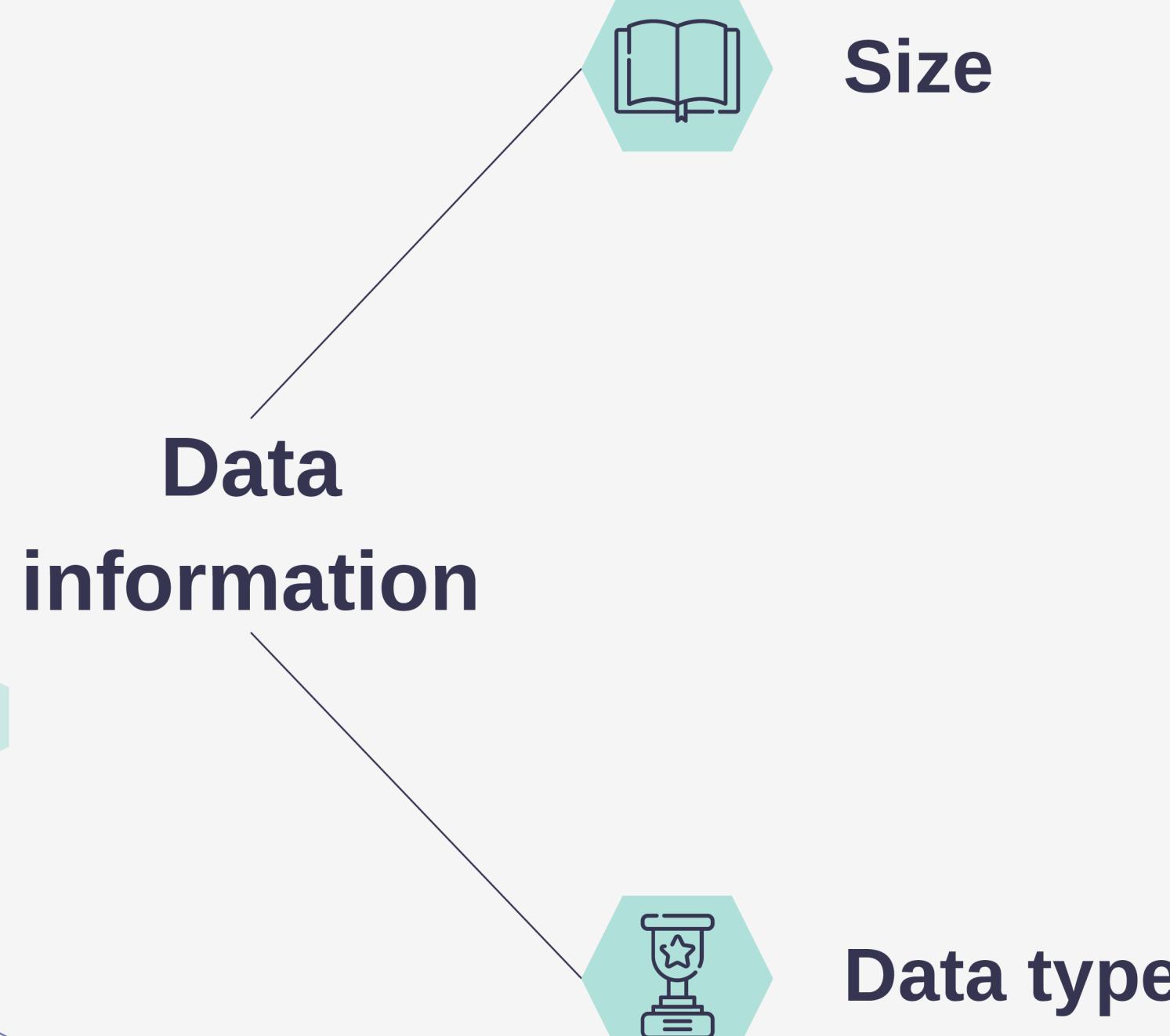
# 03

# Data

# overview



# Data Overview



The dataset contains approximately 5,106 movies with 17 features.

#	Column	Non-Null Count	Dtype
0	Title	5106 non-null	object
1	Tomatoes CriticScore	4129 non-null	float64
2	Tomatoes UserScore	4523 non-null	float64
3	Link	5106 non-null	object
4	PlatformReleased	5106 non-null	object
5	Cast	5074 non-null	object
6	Director	5071 non-null	object
7	Genre	5060 non-null	object
8	Rating	4015 non-null	category
9	Runtime	5036 non-null	object
10	Studio	5059 non-null	object
11	Release Date	5031 non-null	datetime64[ns]
12	Production Budget	5106 non-null	float64
13	Domestic Gross	5106 non-null	float64
14	Worldwide Gross	5106 non-null	float64
15	Metascore	4353 non-null	float64
16	Meta UserScore	4186 non-null	float64
dtypes: category(1), datetime64[ns](1), float64(7), object(8)			

The dataset fully contains the general features of each movie

# Key features



## Title

Movie name



## Critic/User Scores

Ratings from critics and users



## Financials

Production budgets, domestic and worldwide gross revenue



## Categoricals

Genre, platform, studio, cast, director

## Time Information

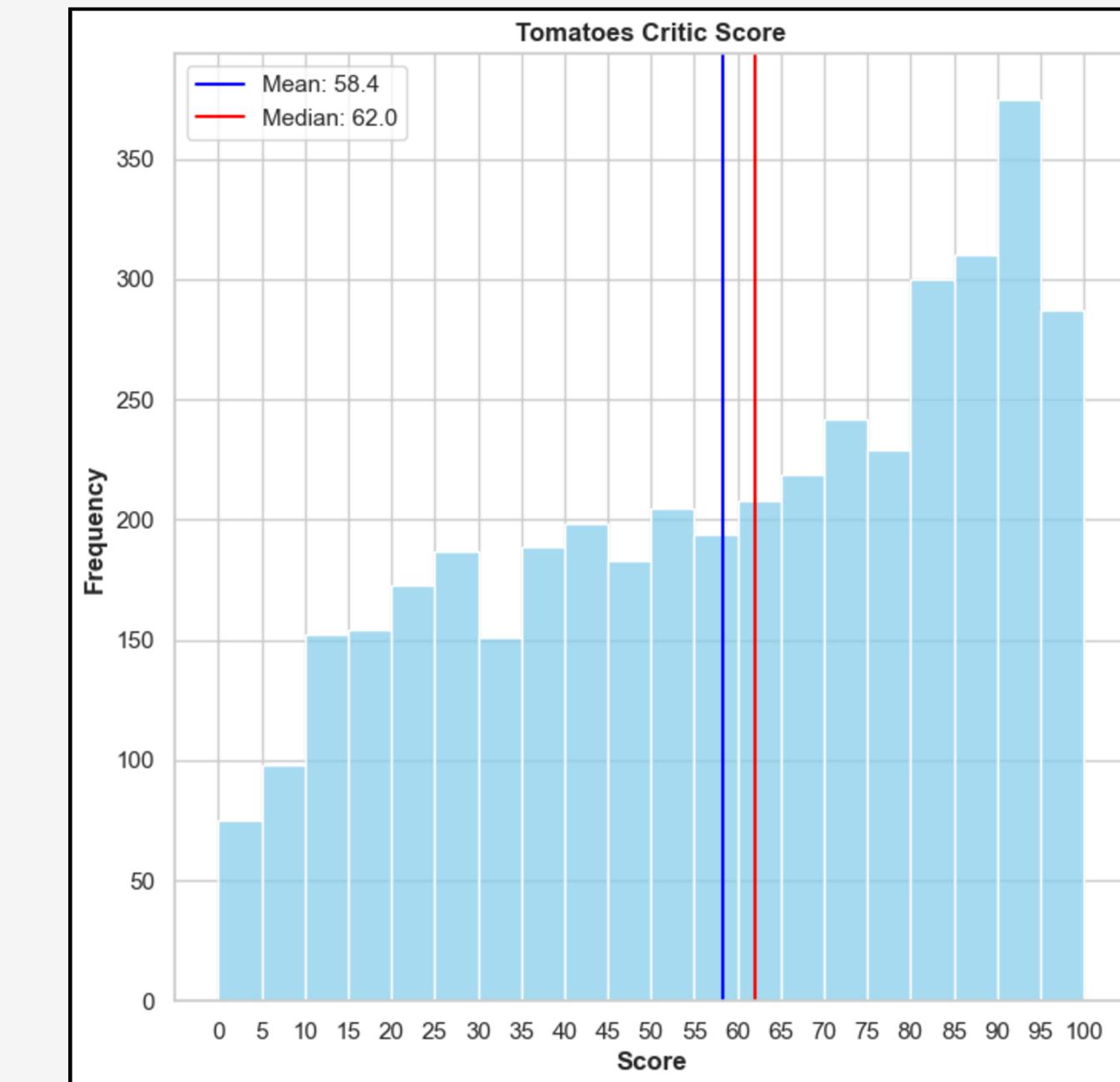
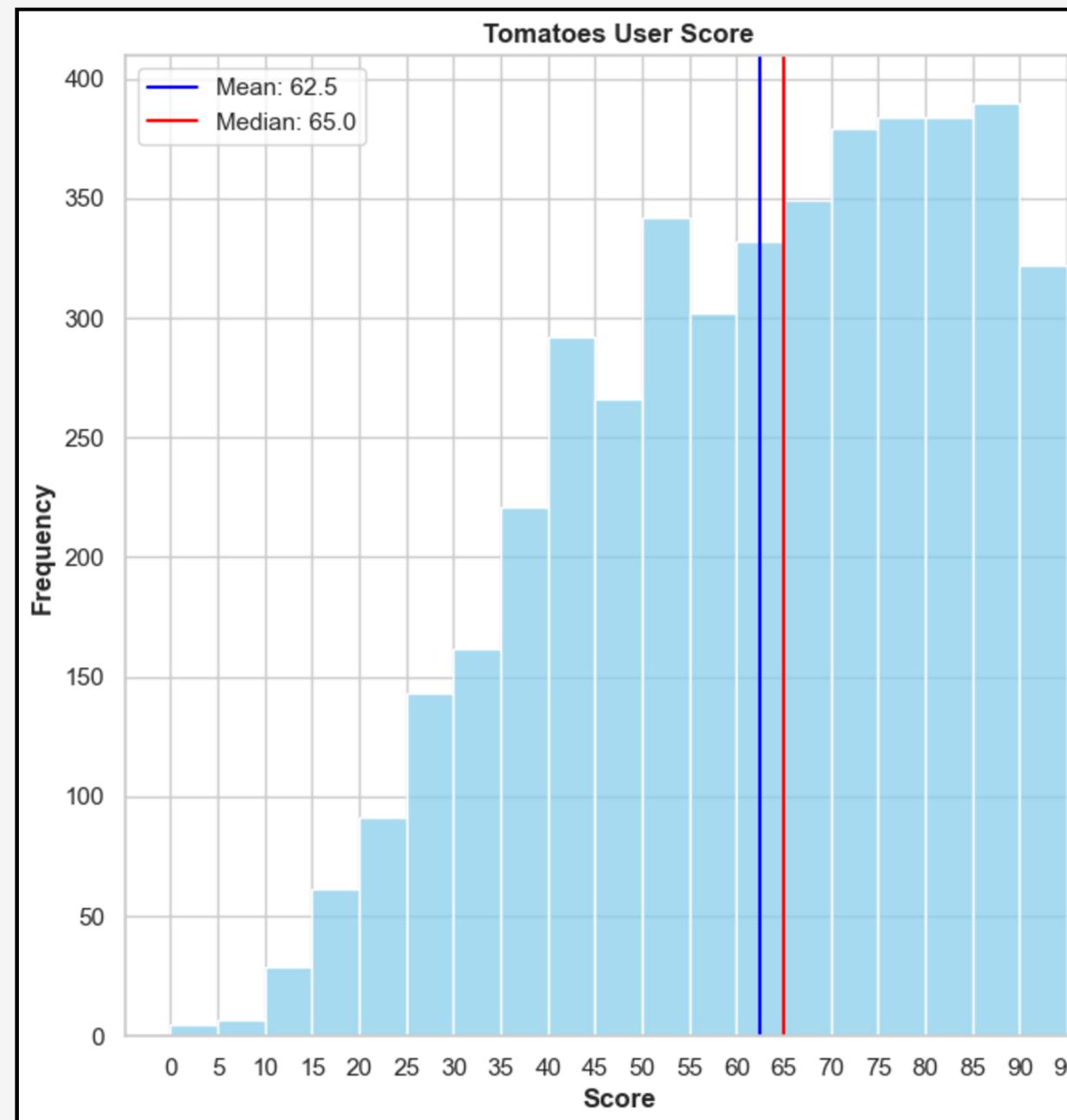
Release date, runtime

# Numerical Columns

## Rotten Tomatoes Scores

All the scores are out of 100

- Tomatoes UserScore: Scores from users/ viewers from Rotten Tomatoes
- Tomatoes CriticScore: Scores from critics from Rotten Tomatoes

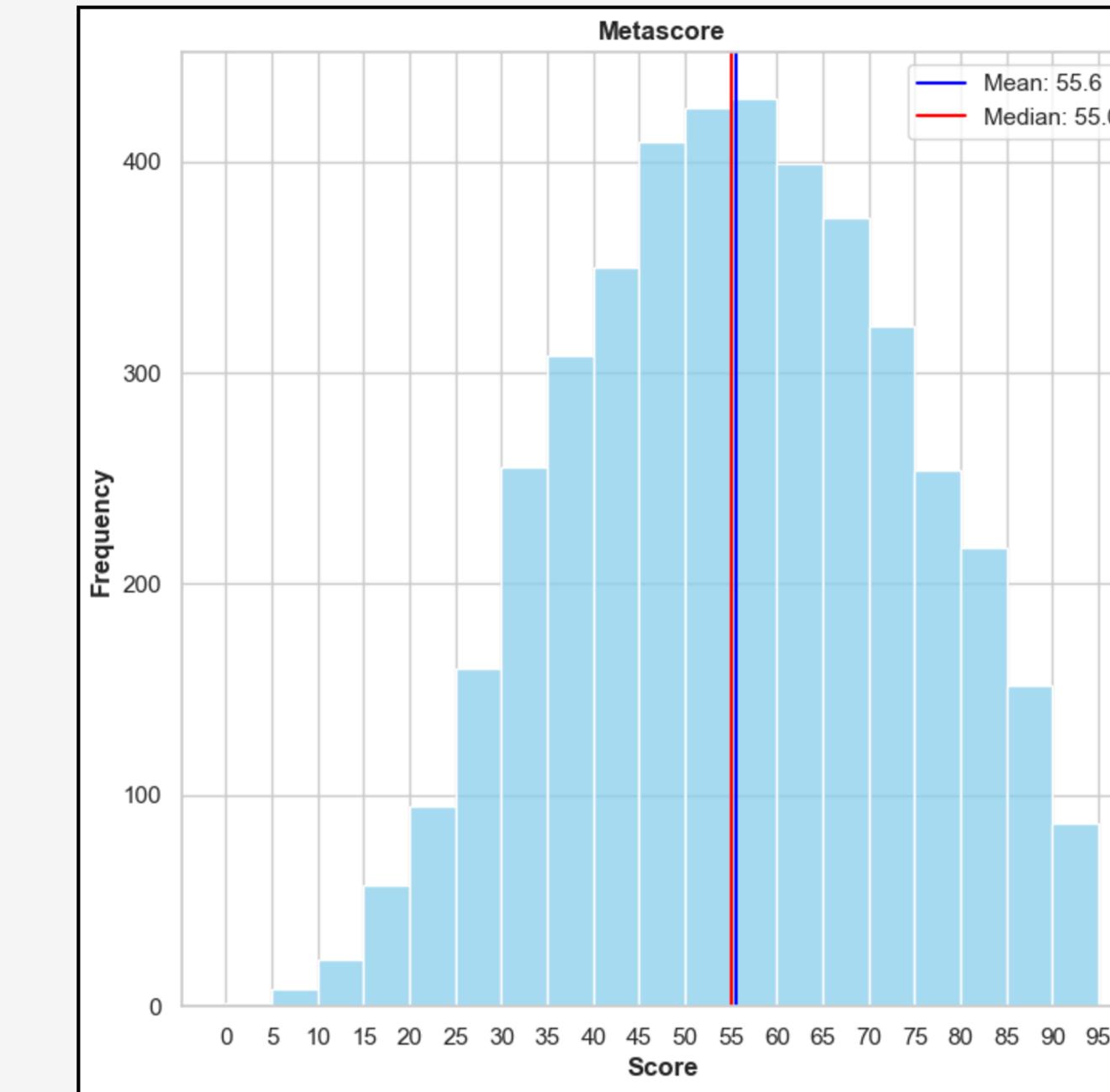
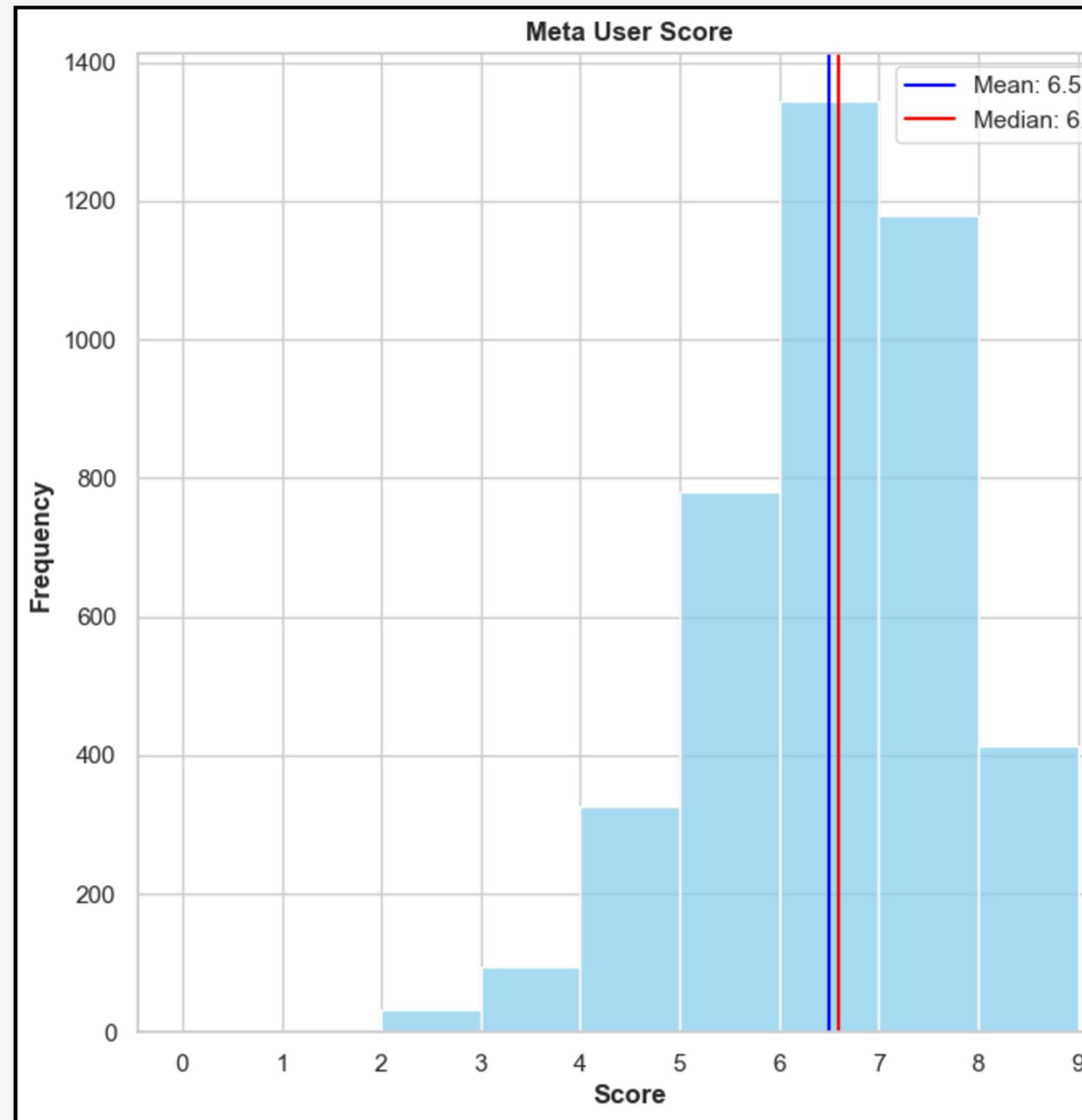


# Numerical Columns

## • MetaCritic Scores

MetaScore is out of 100, except for Meta User Score is out of 10

- Meta User Score: Scores from users/ viewers from MetaCritic
- MetaScore: Scores from critics from MetaCritic



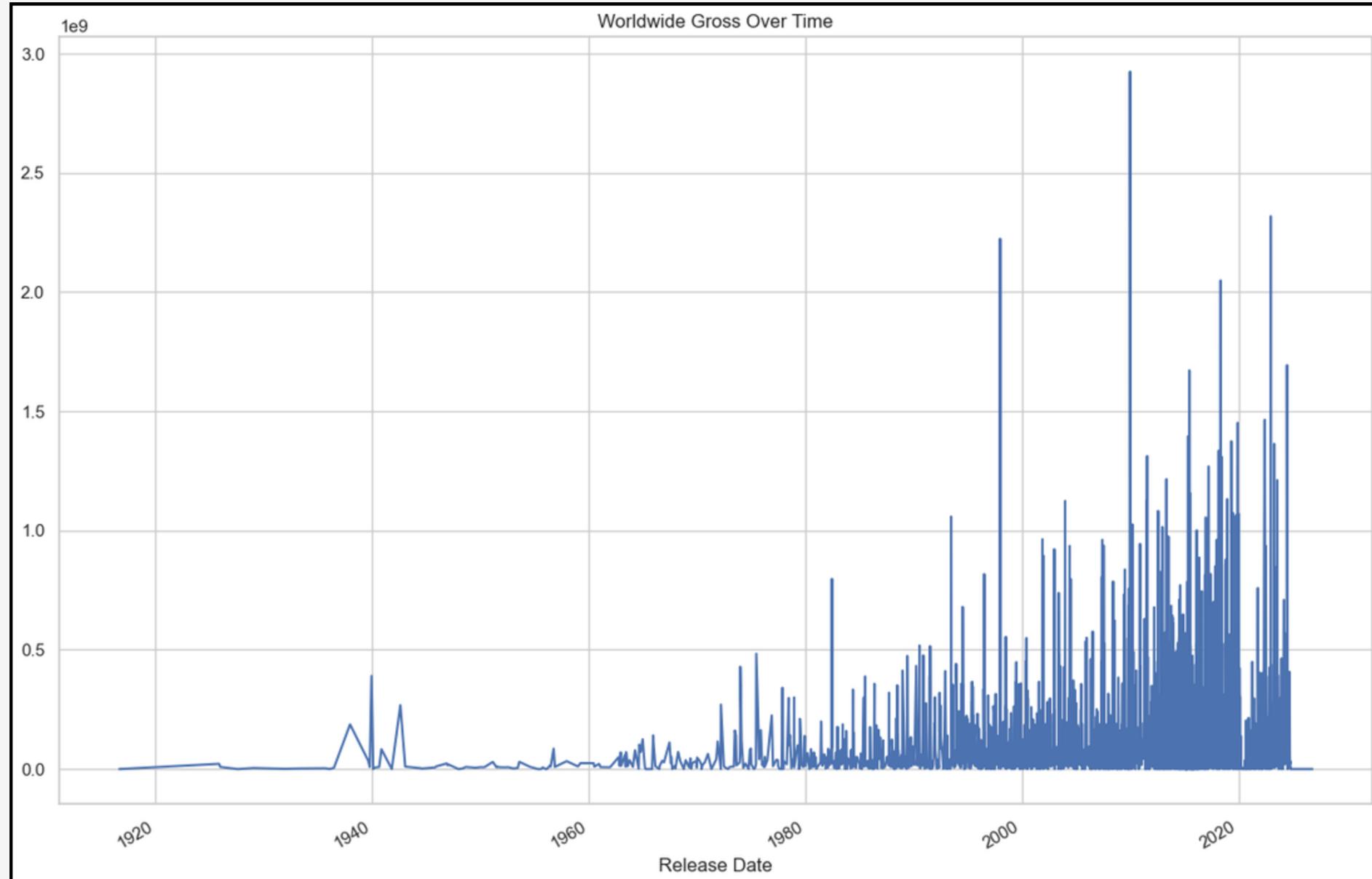
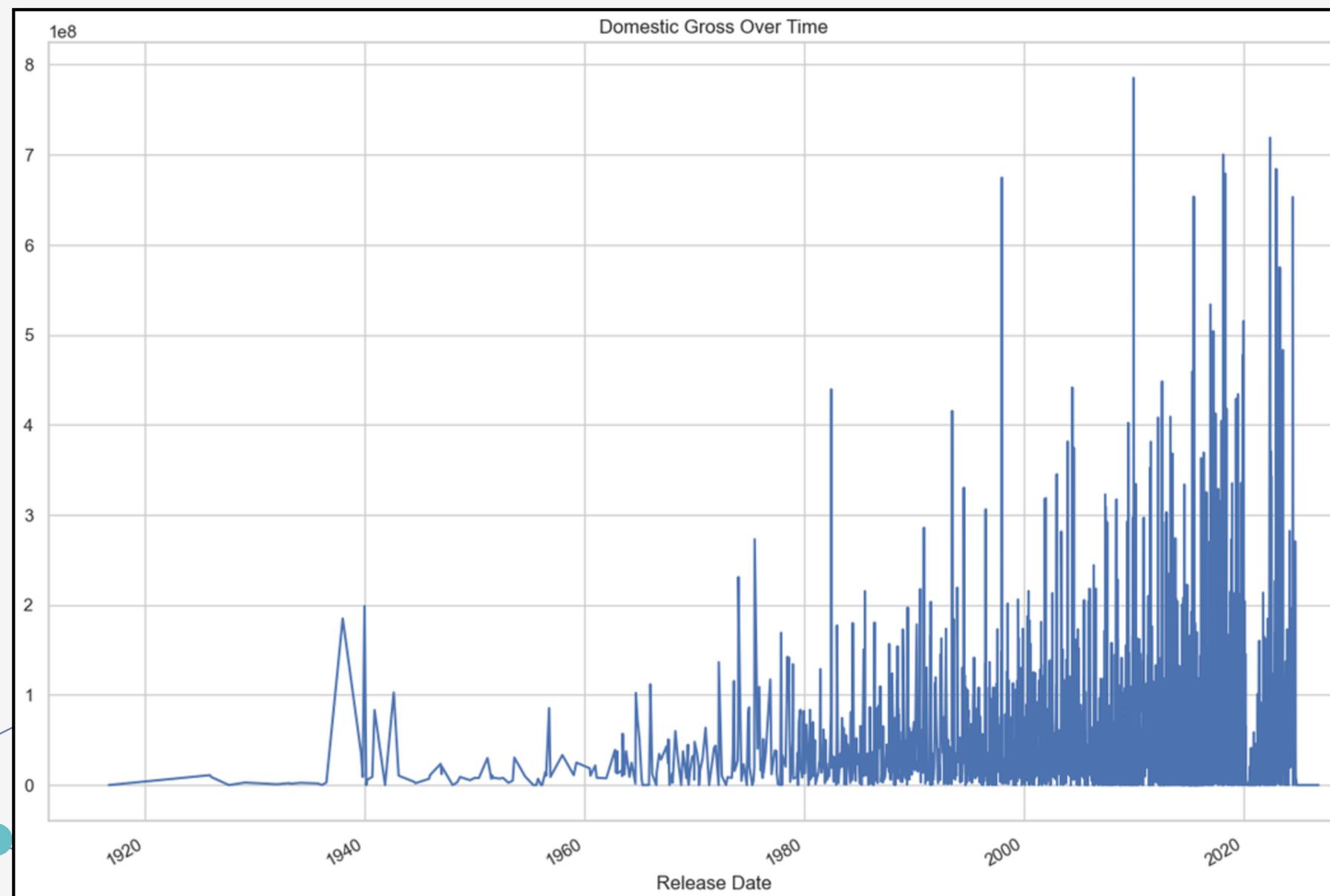


# Numerical Columns

## Economic Features

All the factors are calculated in billions

- Domestic/Worldwide Gross: Revenues from national/international market
- Production Budget: The expenditure spent to produce movies

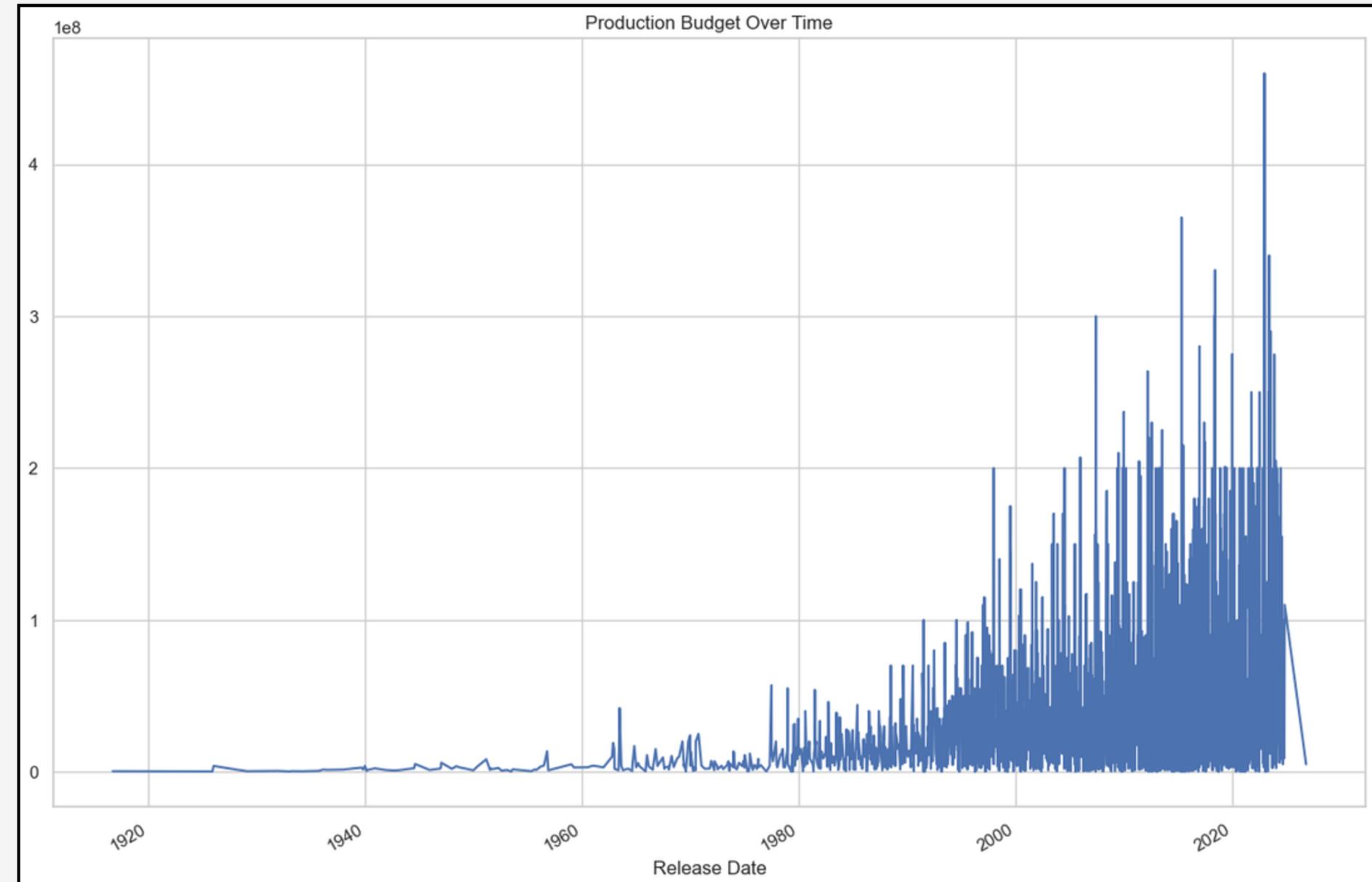


# Numerical Columns

## Economic Features

All the factors are calculated in billions

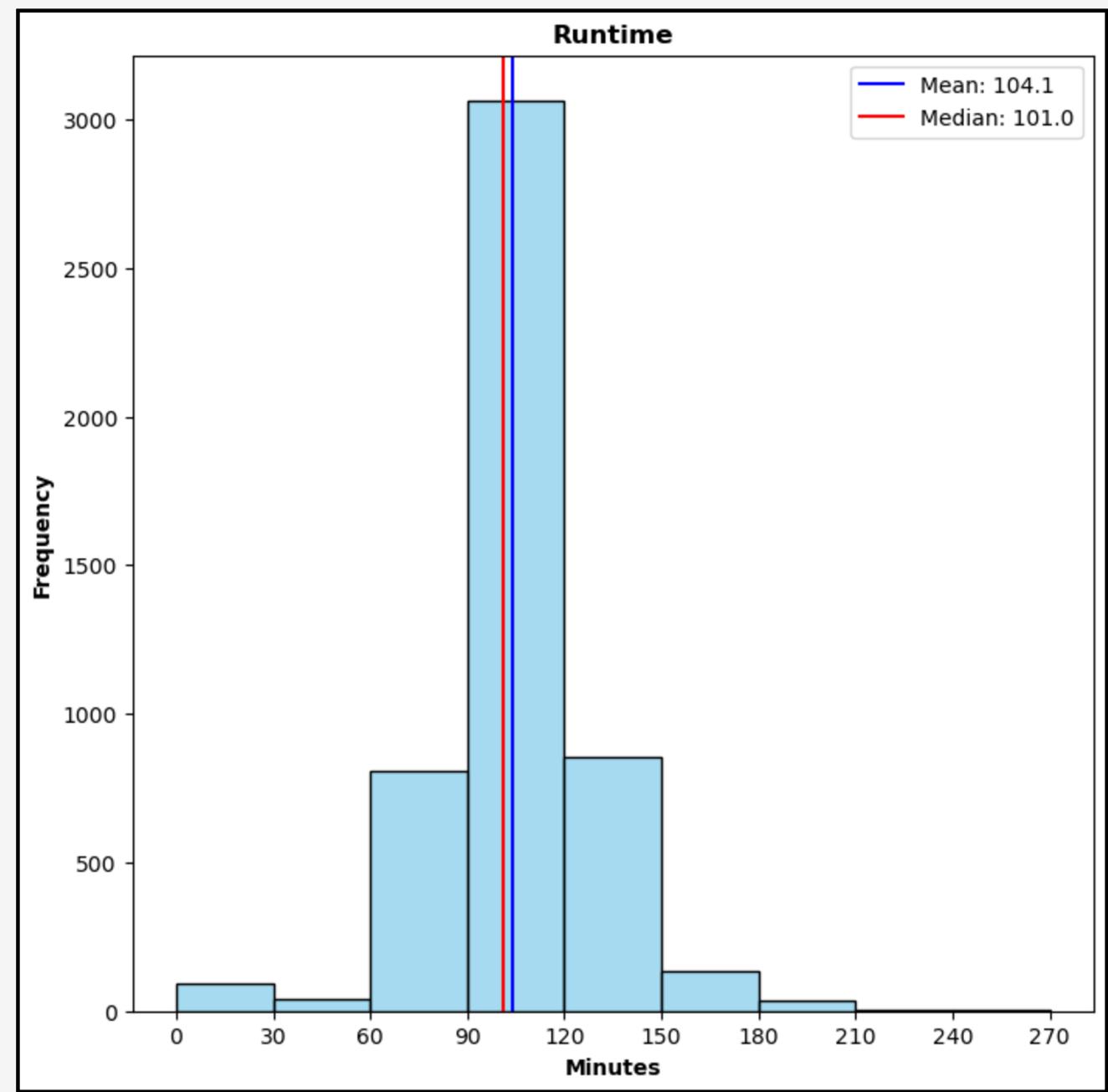
- Domestic/Worldwide Gross: Revenues from national/international market
- Production Budget: The expenditure spent to produce movies



# Numerical Columns

## Runtime

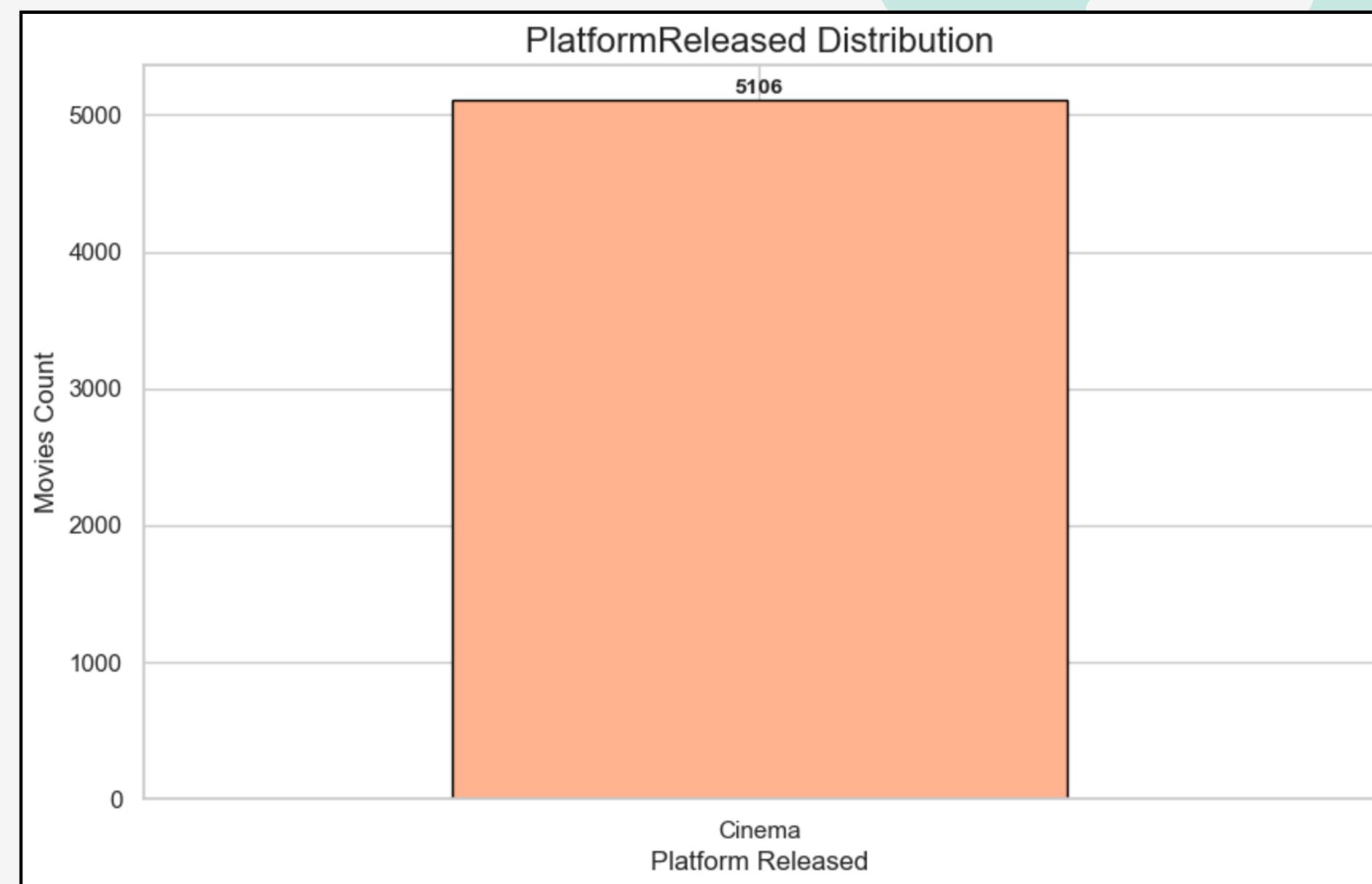
Time duration of each movie



# Categorical Columns

- Platform Released

Places that movies are released like Netflix, Cinema, etc

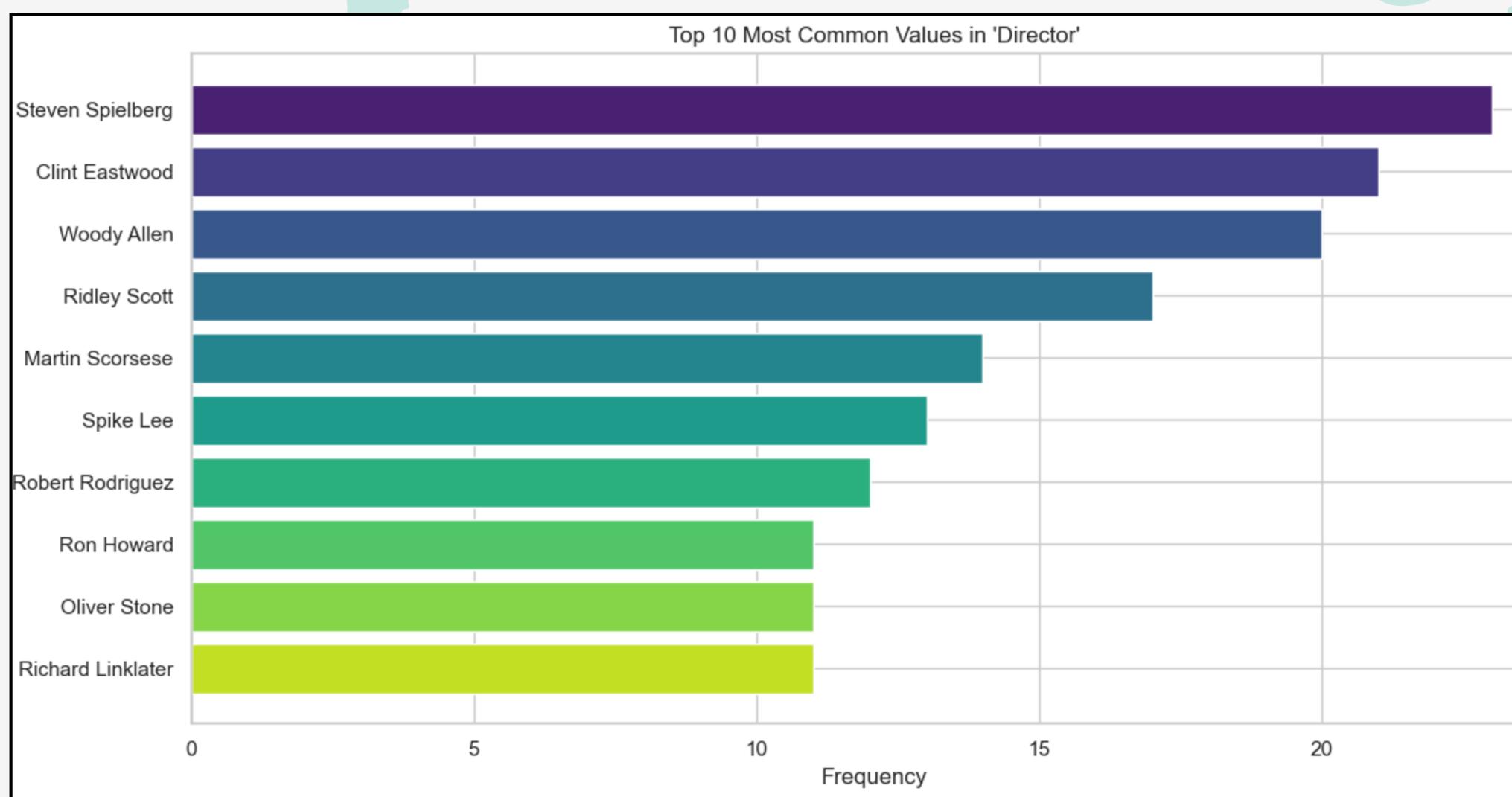
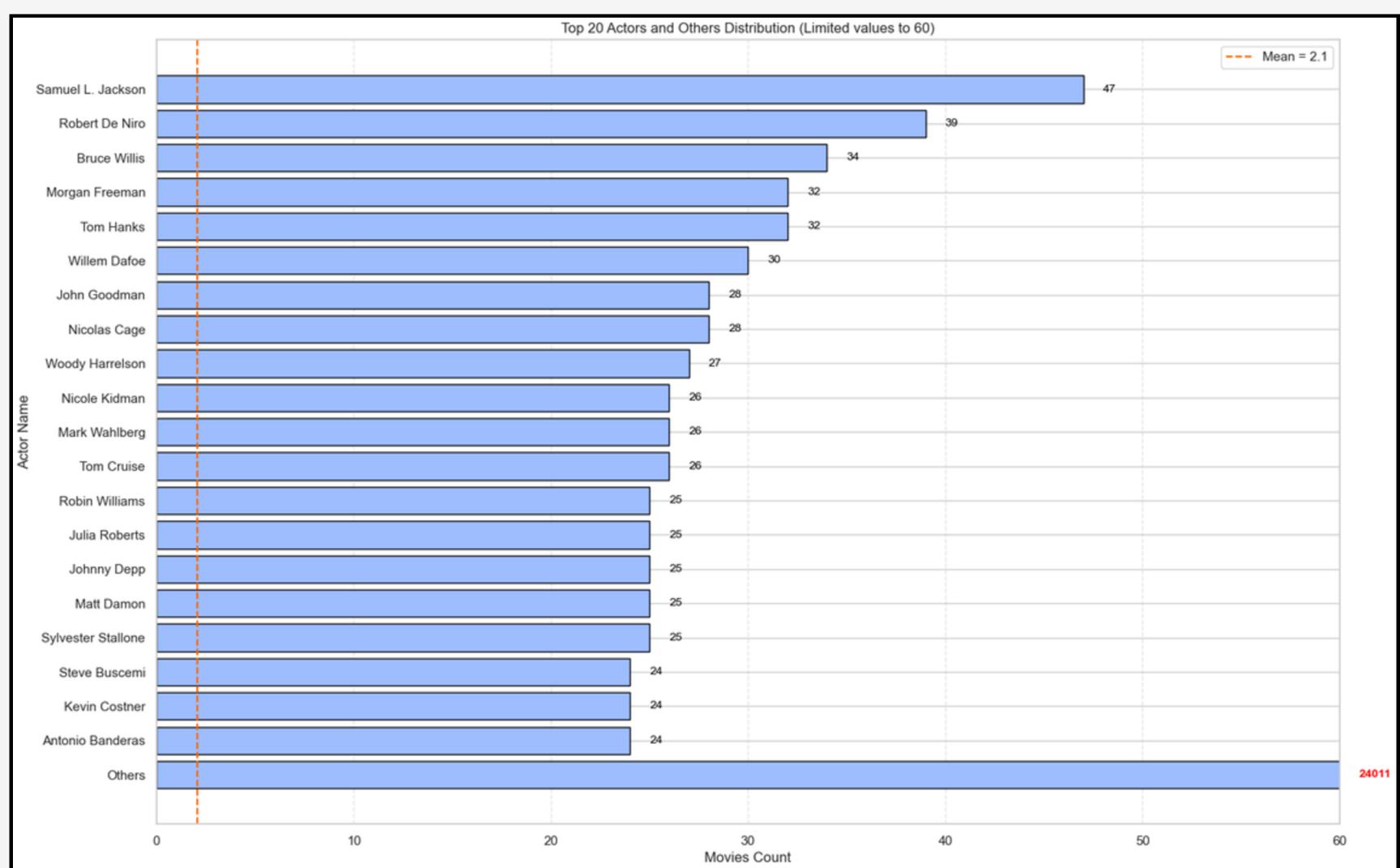


# Categorical Columns

## Casts, Directors and Studios

All factors that contribute to creating a movie

- Casts: Actors/Actresses that played in a movie
- Directors: Who have control over movie production
- Studios: Companies that are responsible for releasing movies

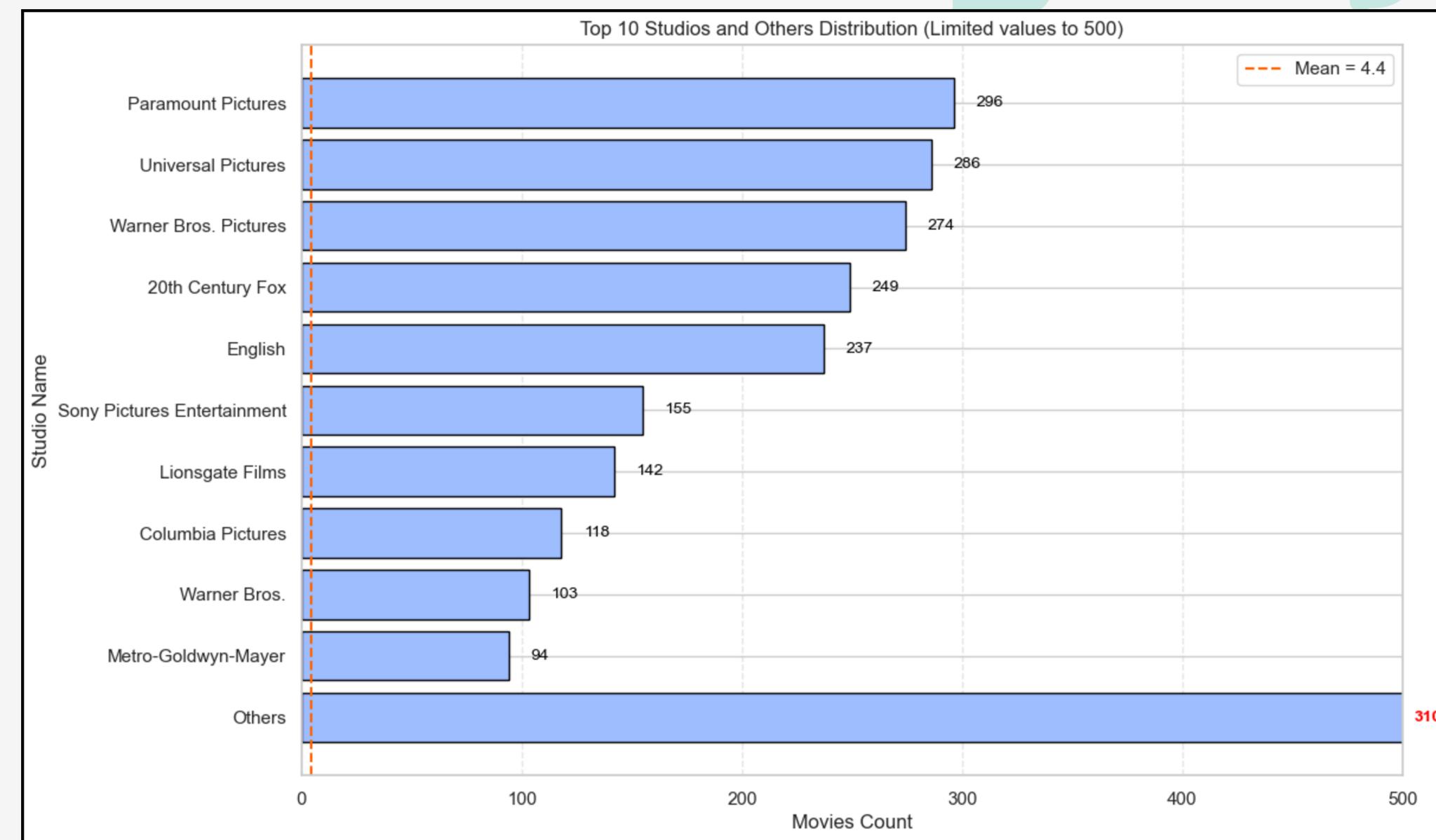


# Categorical Columns

## • Casts, Directors and Studios

All factors that contribute to creating a movie

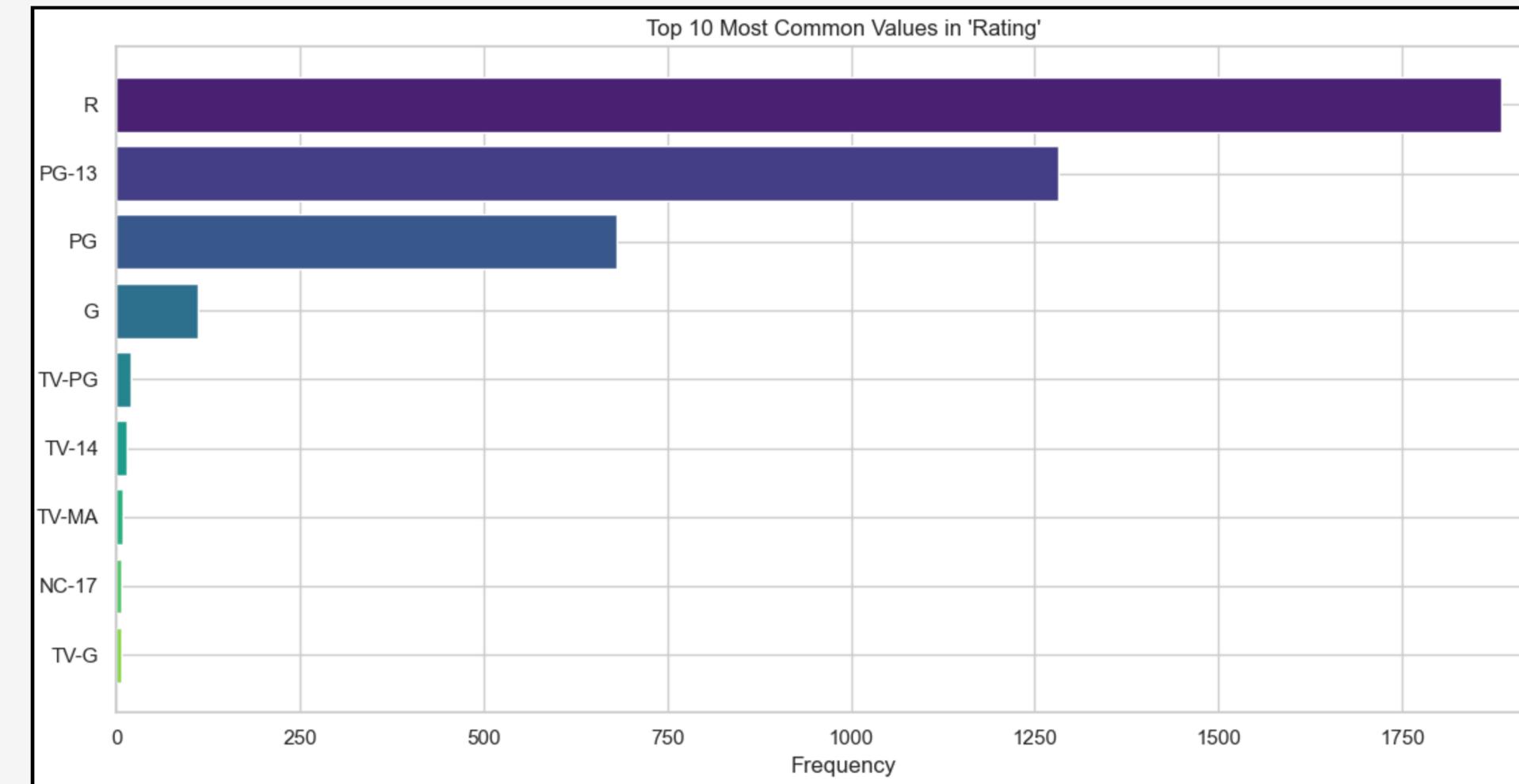
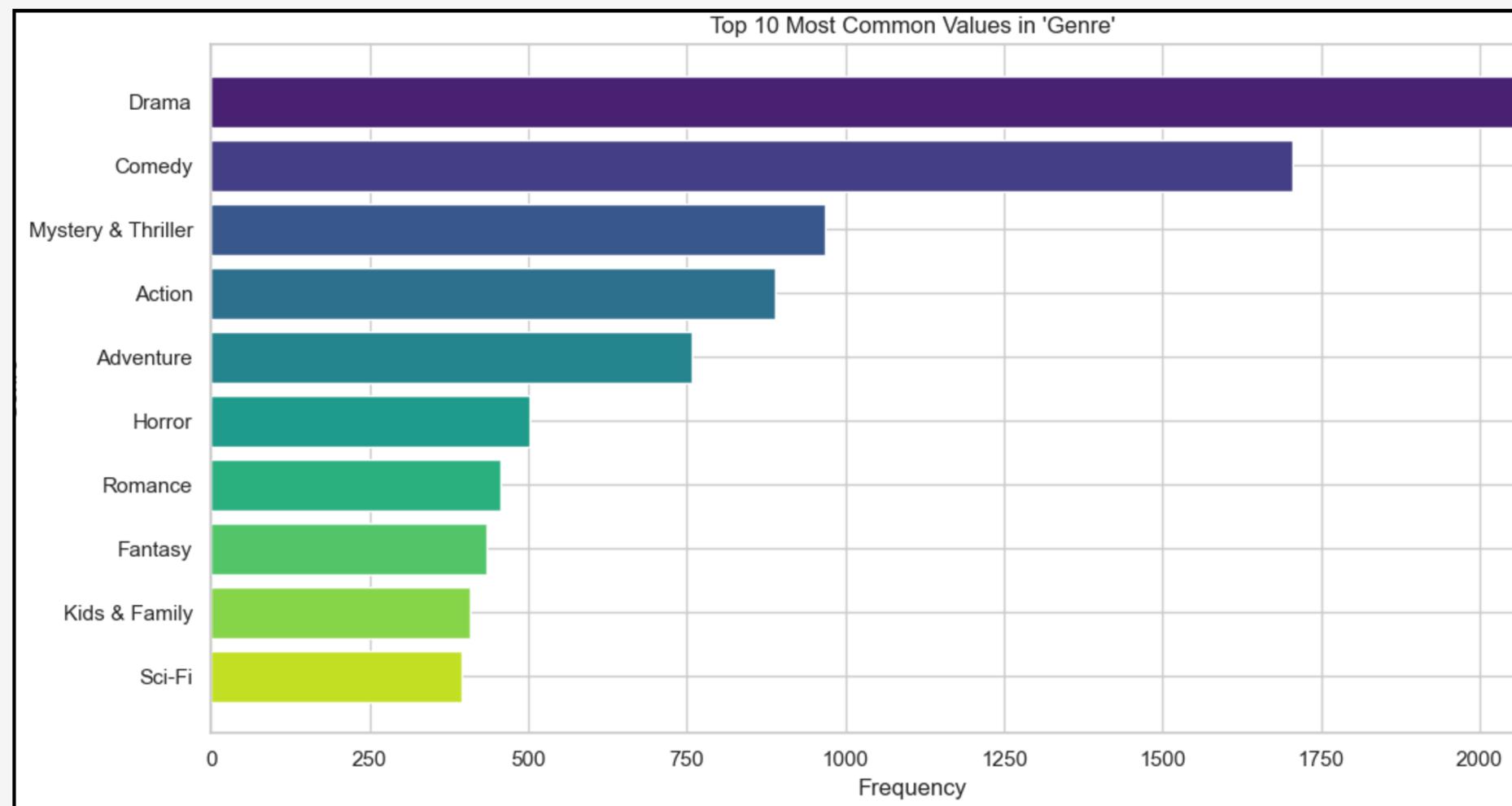
- Casts: Actors/Actresses that played in a movie
- Directors: Who have control over movie production
- Studios: Companies that are responsible for releasing movies



# Categorical Columns

## Genres and Ratings

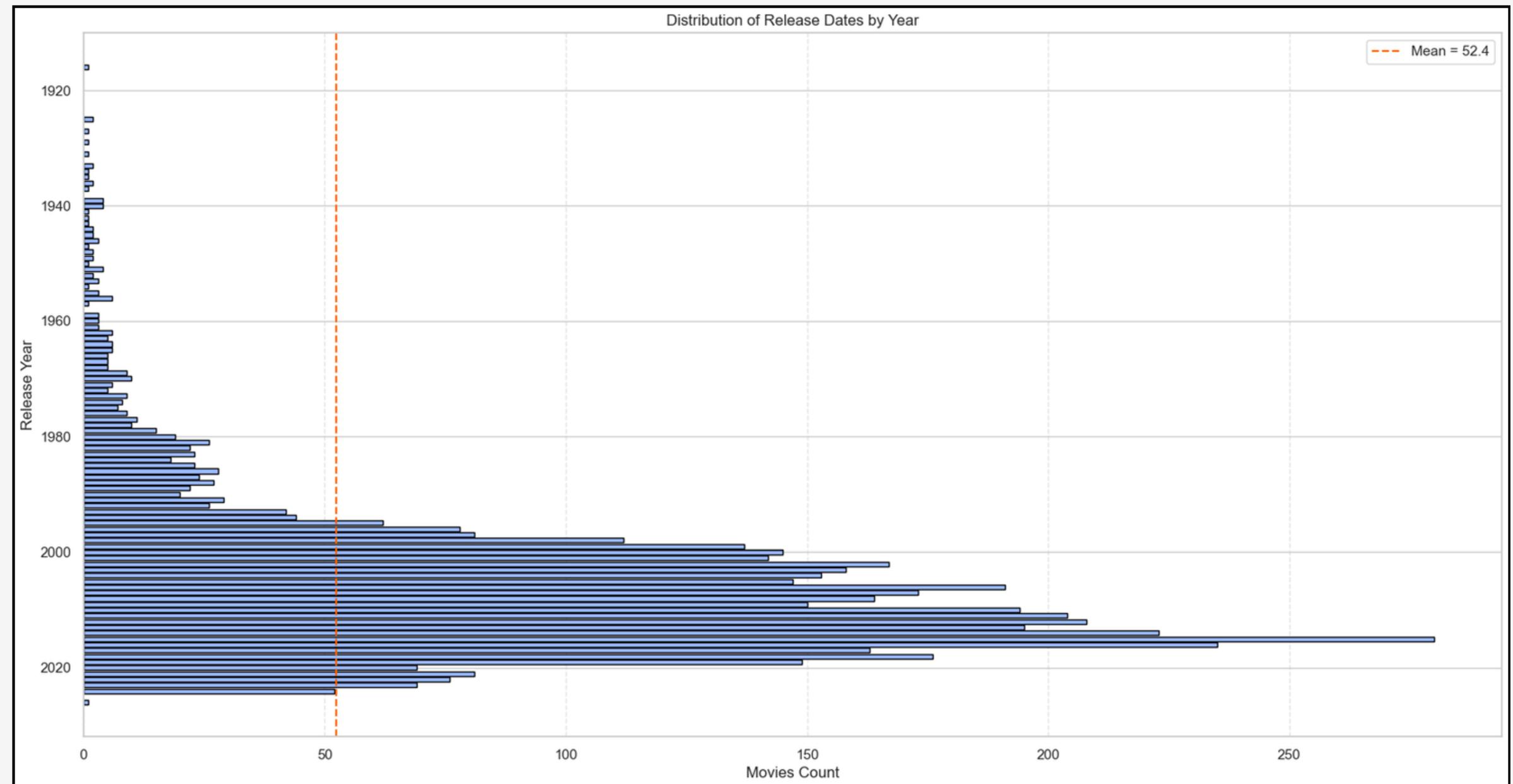
- Genres: A category of artistic composition that creates a movie
- Ratings: Classifies films based on their suitability for audiences



# Categorical Columns

## Release Dates

Date that a movie is released



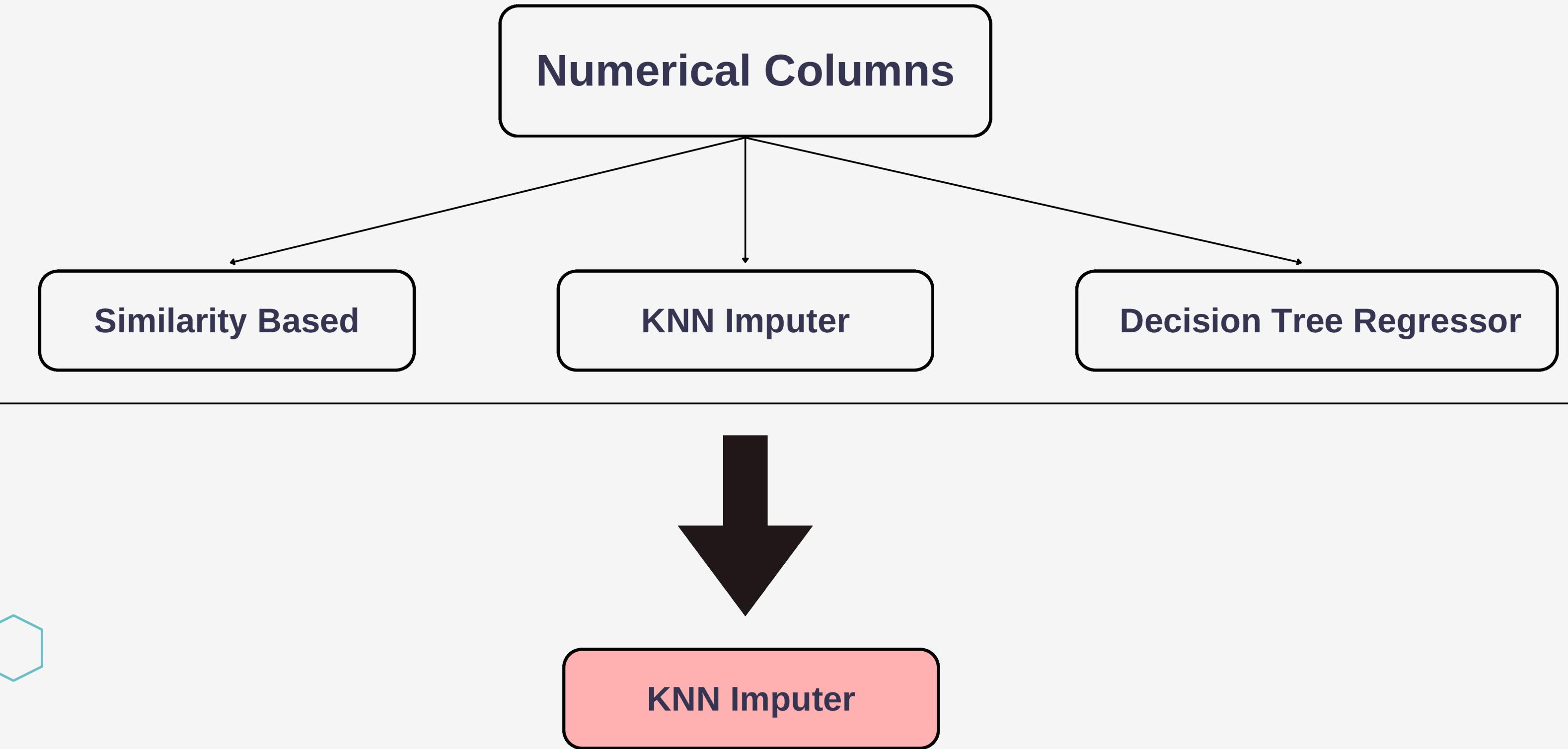
# 04

# Data

# Preprocessing

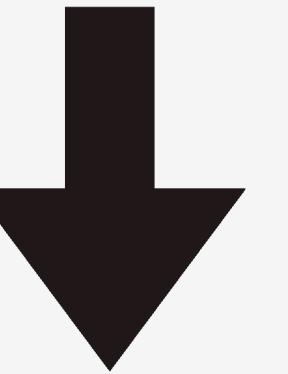


# Handle Missing Values



# Handle Missing Values

Categorical Columns



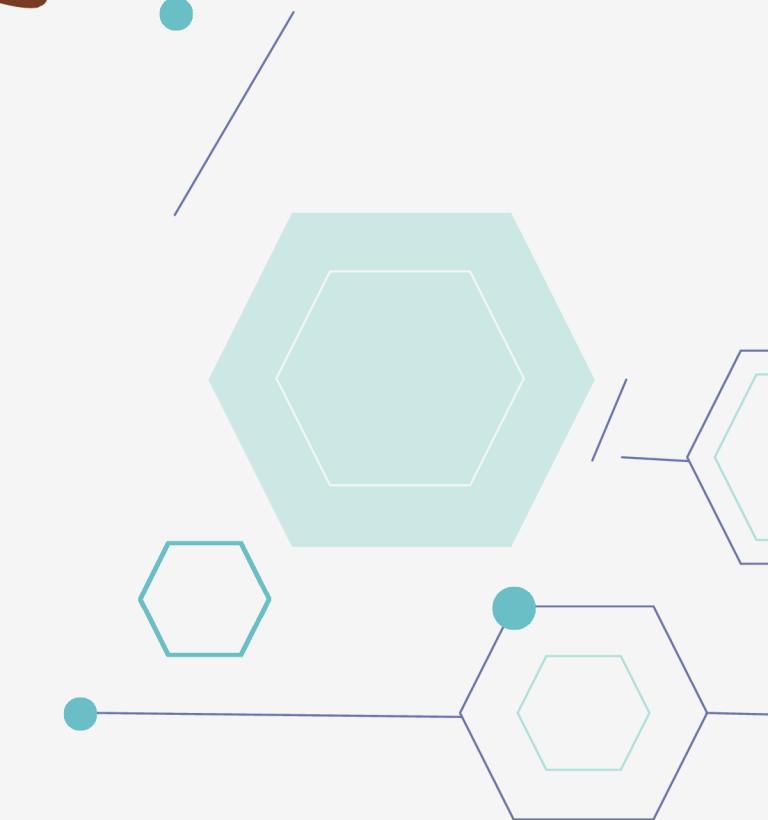
Similarity Based



Decision Tree Classifier

# 05

# Questions



# Question 1:

What are the primary factors influencing critic and user scores?



# Some critical factors (hypothesis)

Genre

Production Budget

Release Time

Marketing & Promotion



# Preprocessing



One-hot encode  
the Genre column



Split release date into  
year, month and day



Calculate the mean  
score for each movie

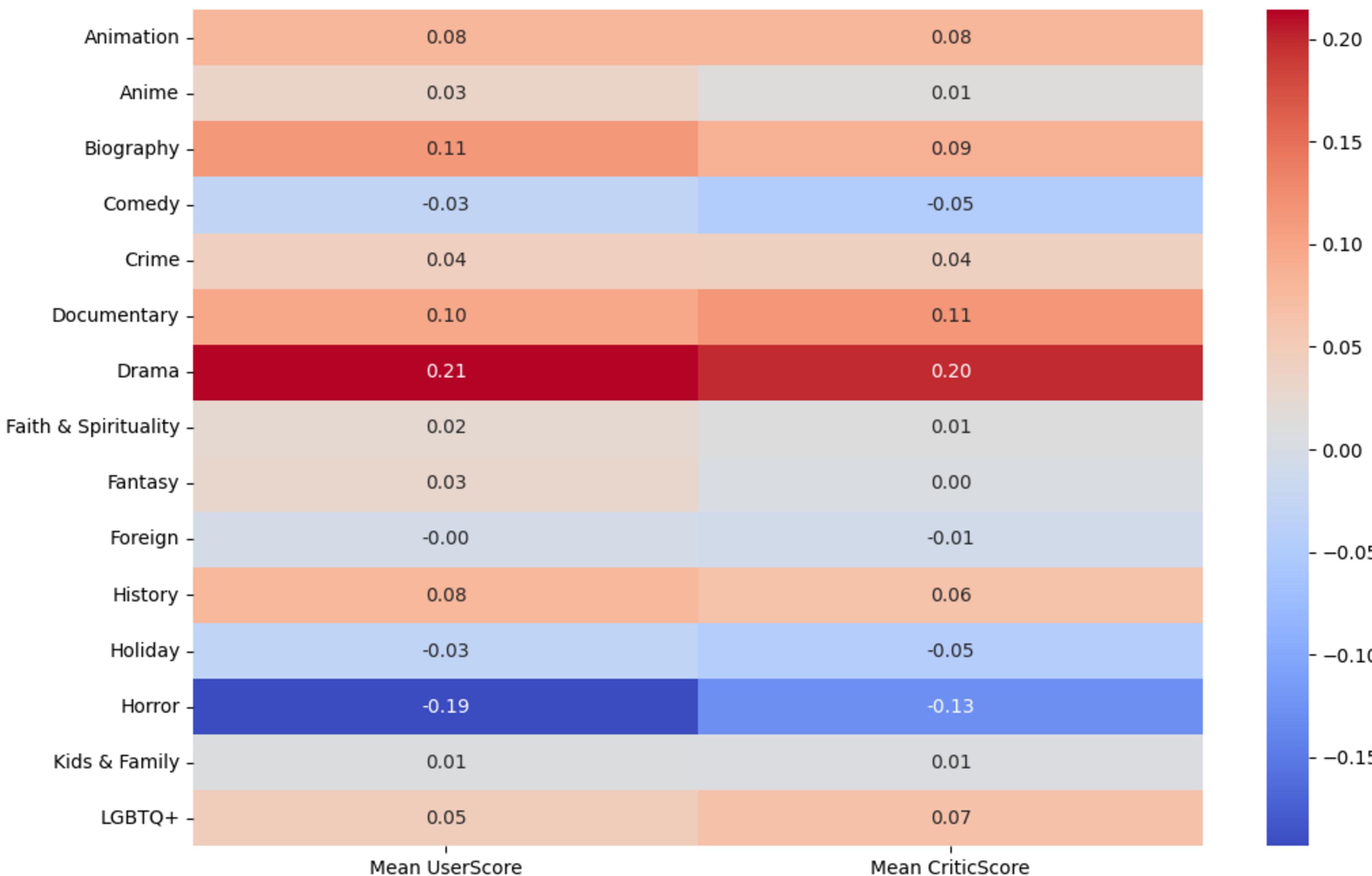


One-hot encode  
the Rating column



Drop columns that  
are not needed

Correlation Between Genres and Scores

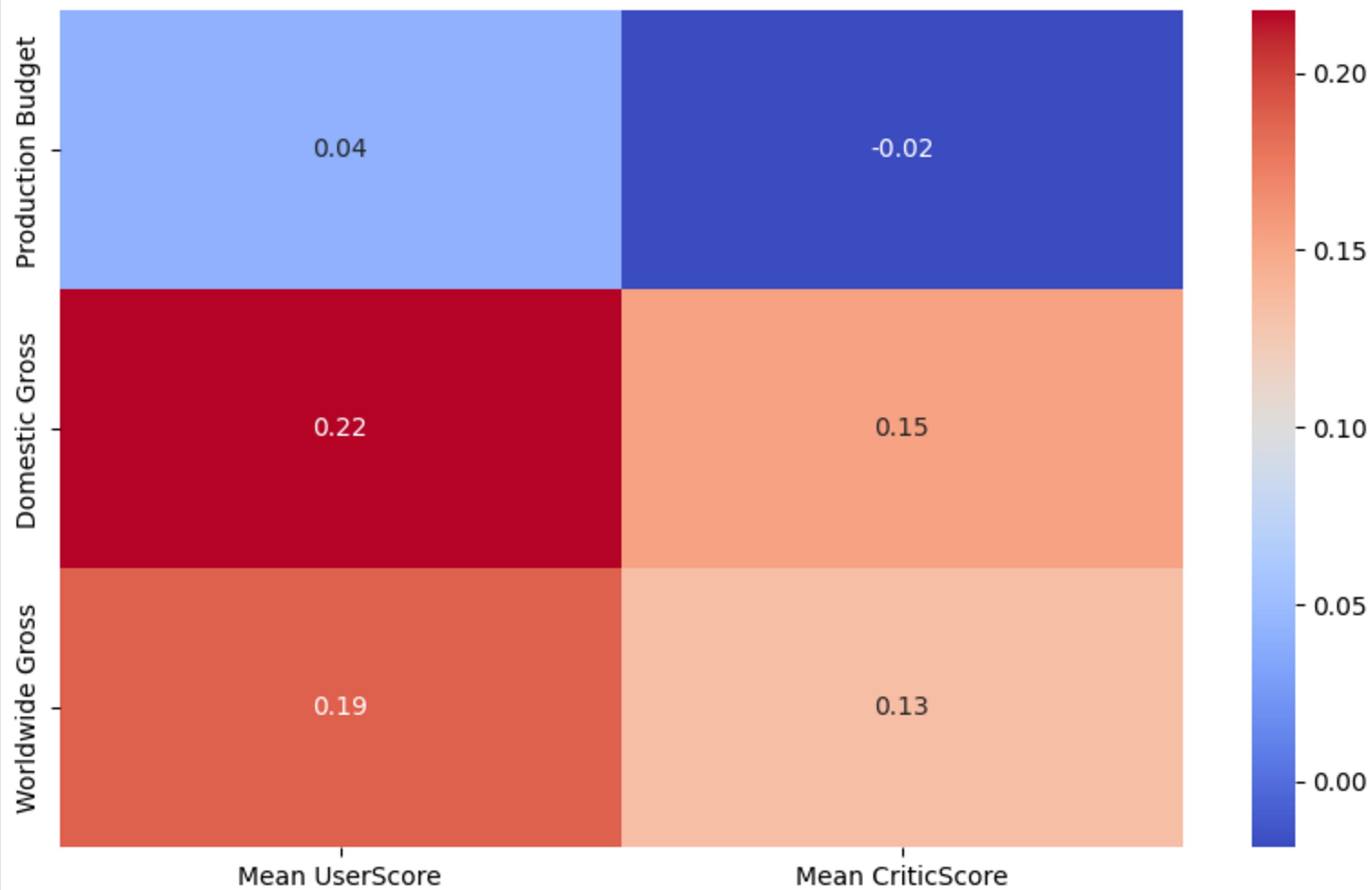


# Key Observations

- **Strong Positive Correlations:**
  - Drama
  - Documentary
- **Weak or Negative Correlations:**
  - Horror
  - Comedy
- **Mixed Correlations:**
  - Animation
  - Biography

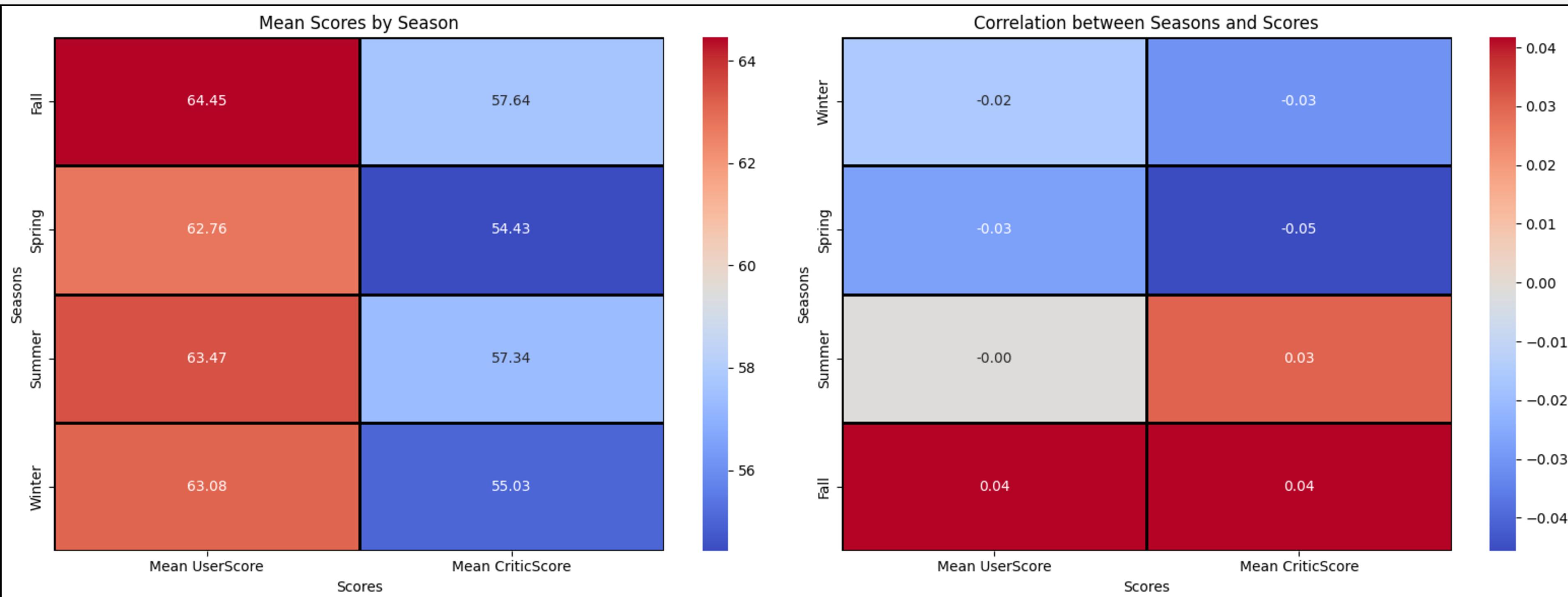


### Correlation Between Economic and Score Columns



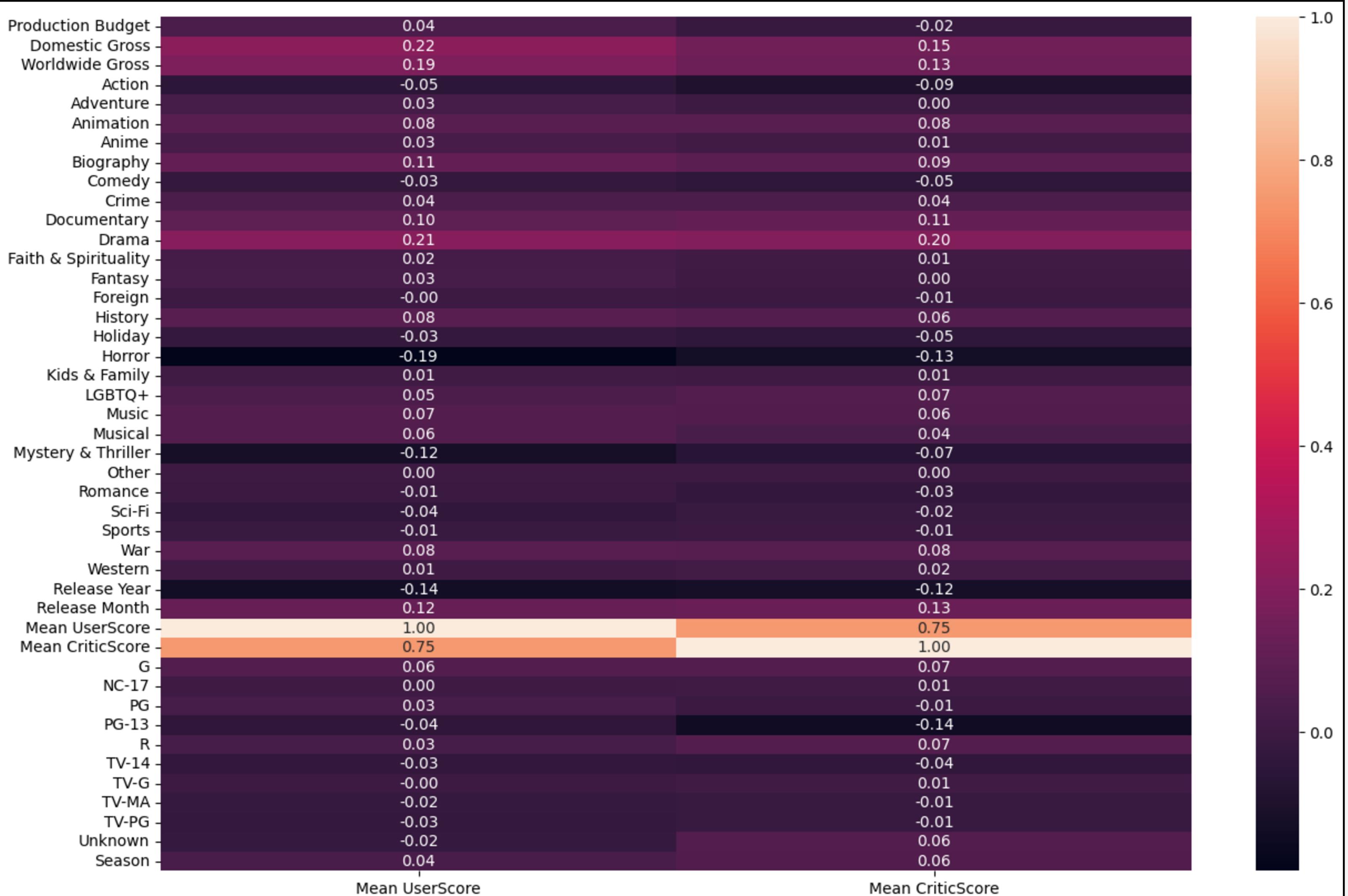
# Key Observations

- **Summary:**
  - **Production Budget:** Weak correlation with ratings.
  - **Domestic and Worldwide Gross:** Strong positive correlation with ratings
- **Inferences:**
  - **Box Office and Acclaim:** Commercial success often aligns with critical and audience approval.
  - **Budget vs. Quality:** A large budget alone is not a reliable predictor of a film's quality or success.



# Key Observations

- **Mean Scores by Season:**
  - **Winter:** Generally sees higher scores from both users and critics.
  - **Spring:** Lower scores compared to other seasons.
  - **Summer and Fall:** Moderate level of scores.
- **Correlation Between Seasons and Scores**
  - **Weak Correlations:** The correlations are generally weak.
  - **Slight Variations:** There are subtle variations in the correlations for user and critic scores across seasons.



# Key Observations

- **Summary:**

- **Consistent Perception:** Movies that perform well with critics tend to be appreciated by users too.
- **Mutual Influences:** High critic scores may influence user expectations, or user enthusiasm might validate critic assessments.

# Conclusion

Domestic and Worldwide Gross

Genre Influences

Production Budget

Strategic release timing

Seasonality

Ratings

Marketing Indicators

Effective marketing

## Question 2:

How do financial metrics like  
budget and gross revenue  
correlate with ratings?



# Preprocessing



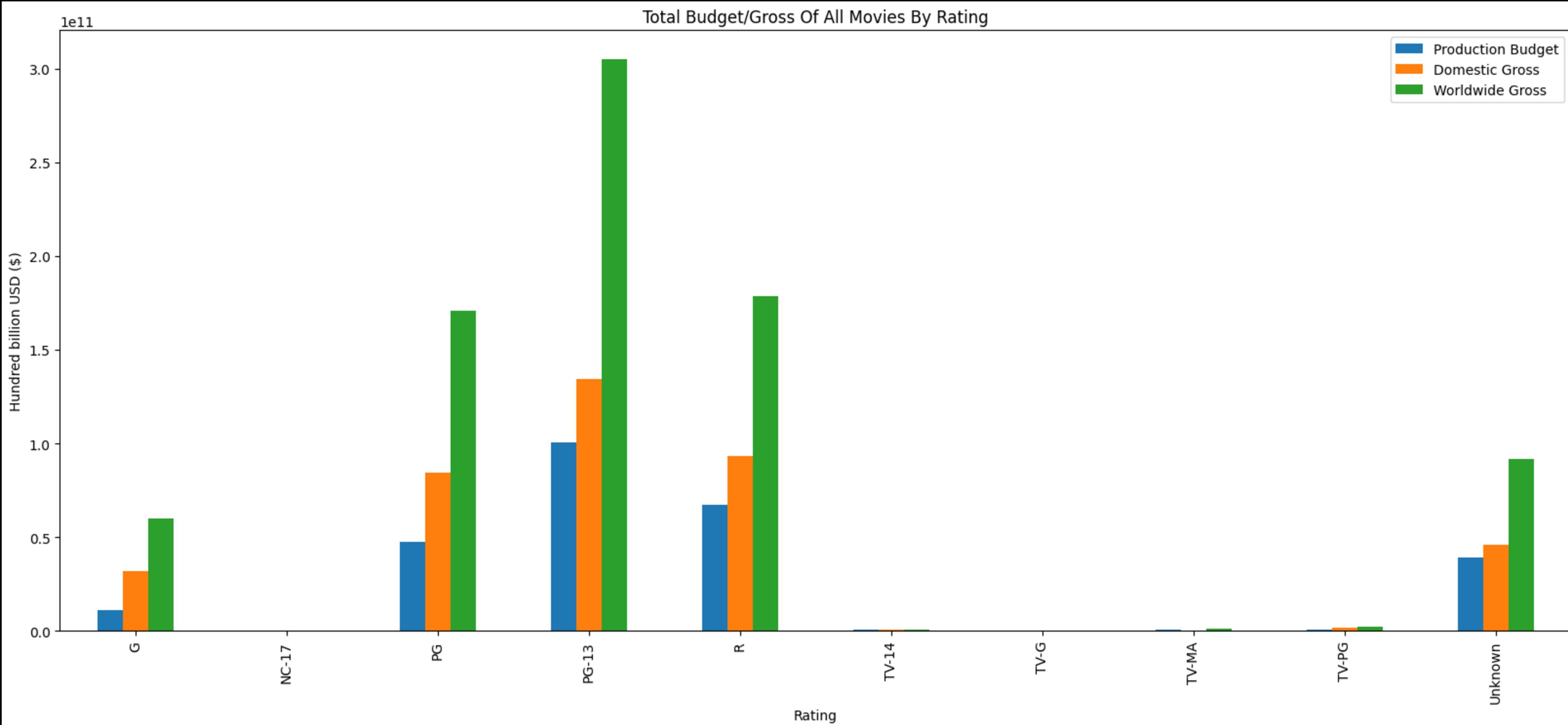
Remove unnecessary  
columns

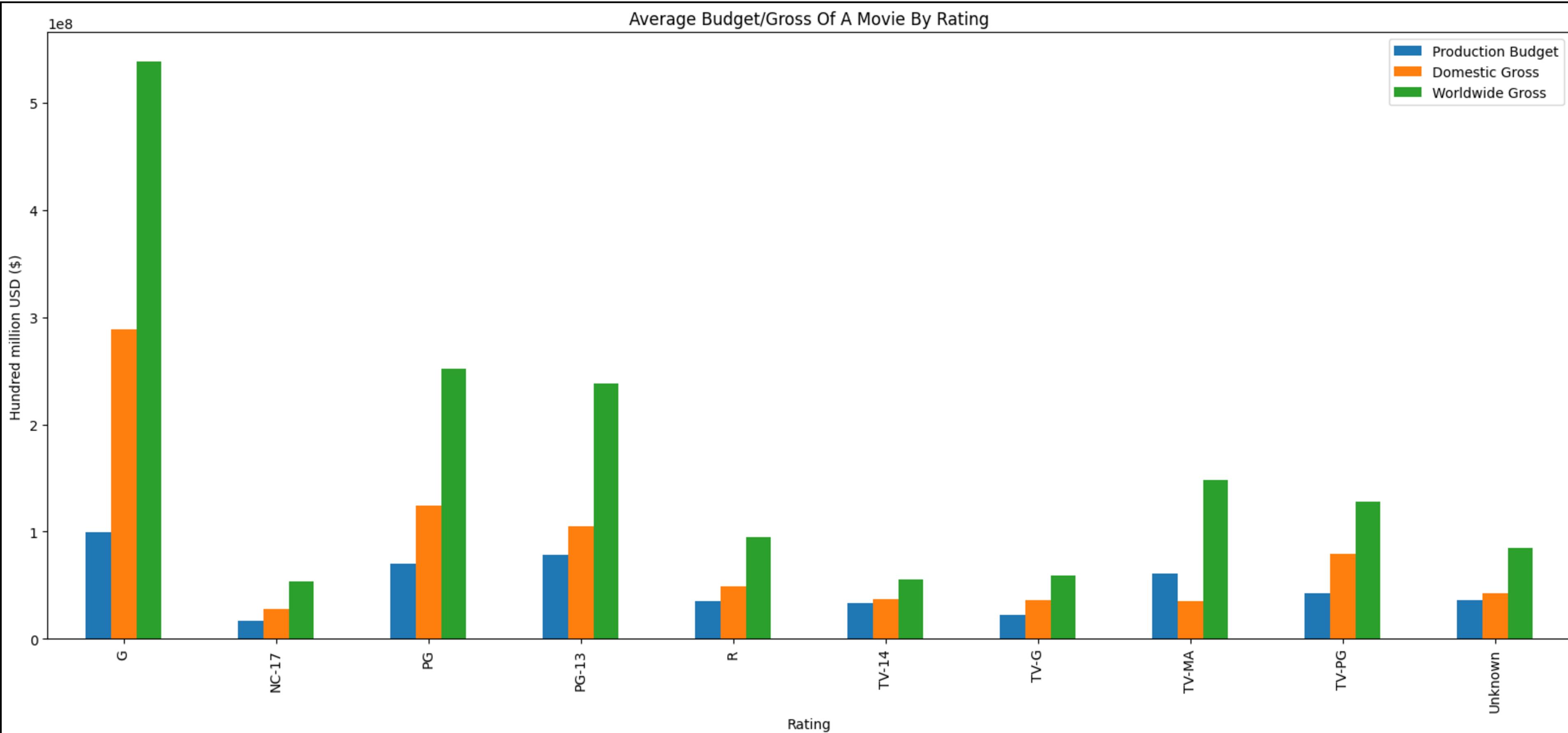


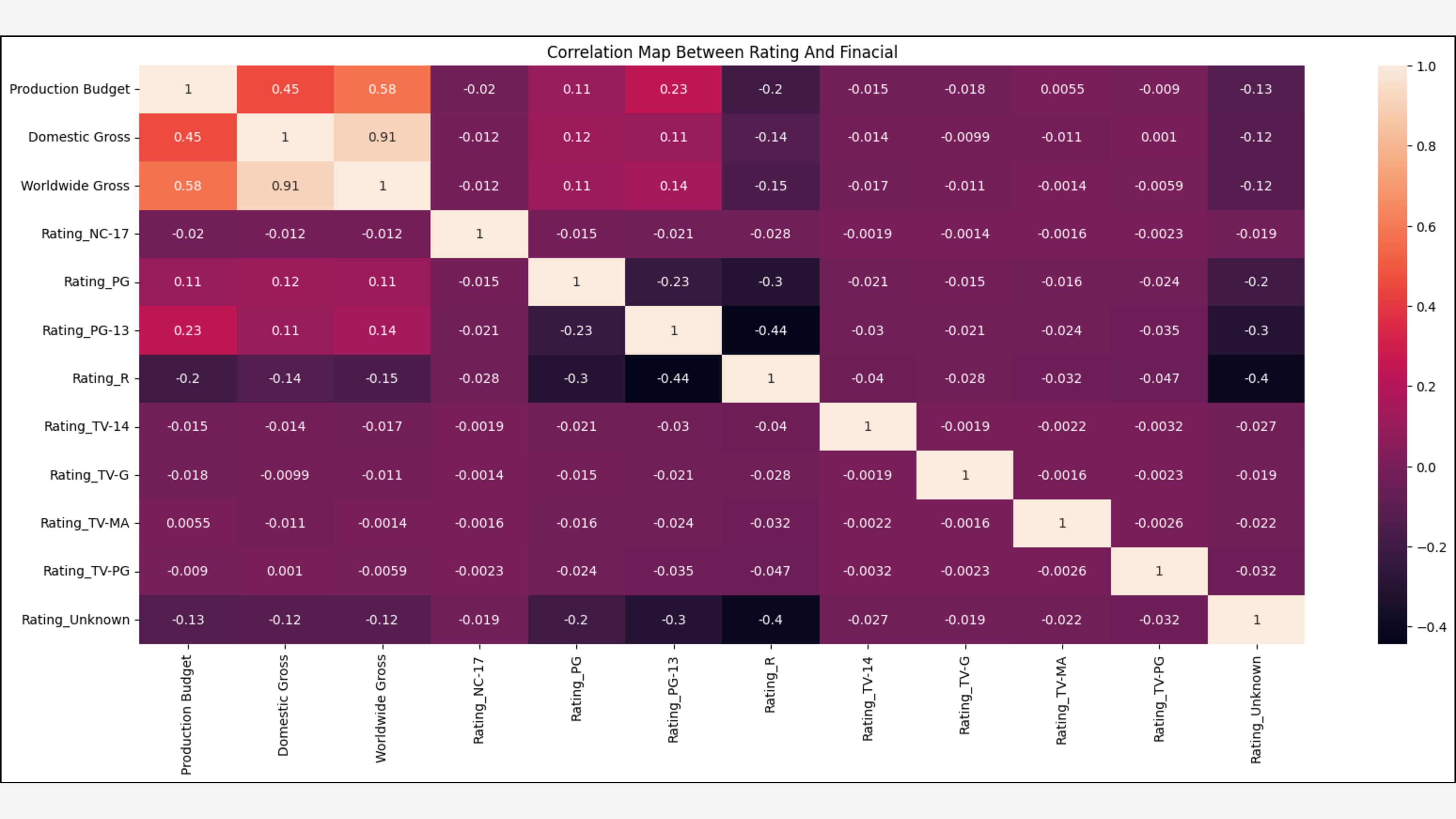
Get release year and  
unit for each year



Apply inflation







# Key Observations

- **PG-13 Movies:**
  - Highest budget and gross.
  - Moderate budget-gross correlation.
  - Broad appeal and global success.
- **R & PG Movies:**
  - Strong budget and gross.
  - R: Mature content, limited reach.
  - PG: Kids-focused, less broad appeal.
- **G-Rated Movies:**
  - Best gross-to-budget ratio.
  - Low costs, high global appeal.
  - Dominated by Disney and Pixar.



# Conclusions

PG-13

Highest Production Budgets  
and box office Grosses

R and PG

Perform well but high  
Production Budget

G

Best return on investment

# Question 3:

What genres tend to perform better in terms of critic and user scores?



# Preprocessing



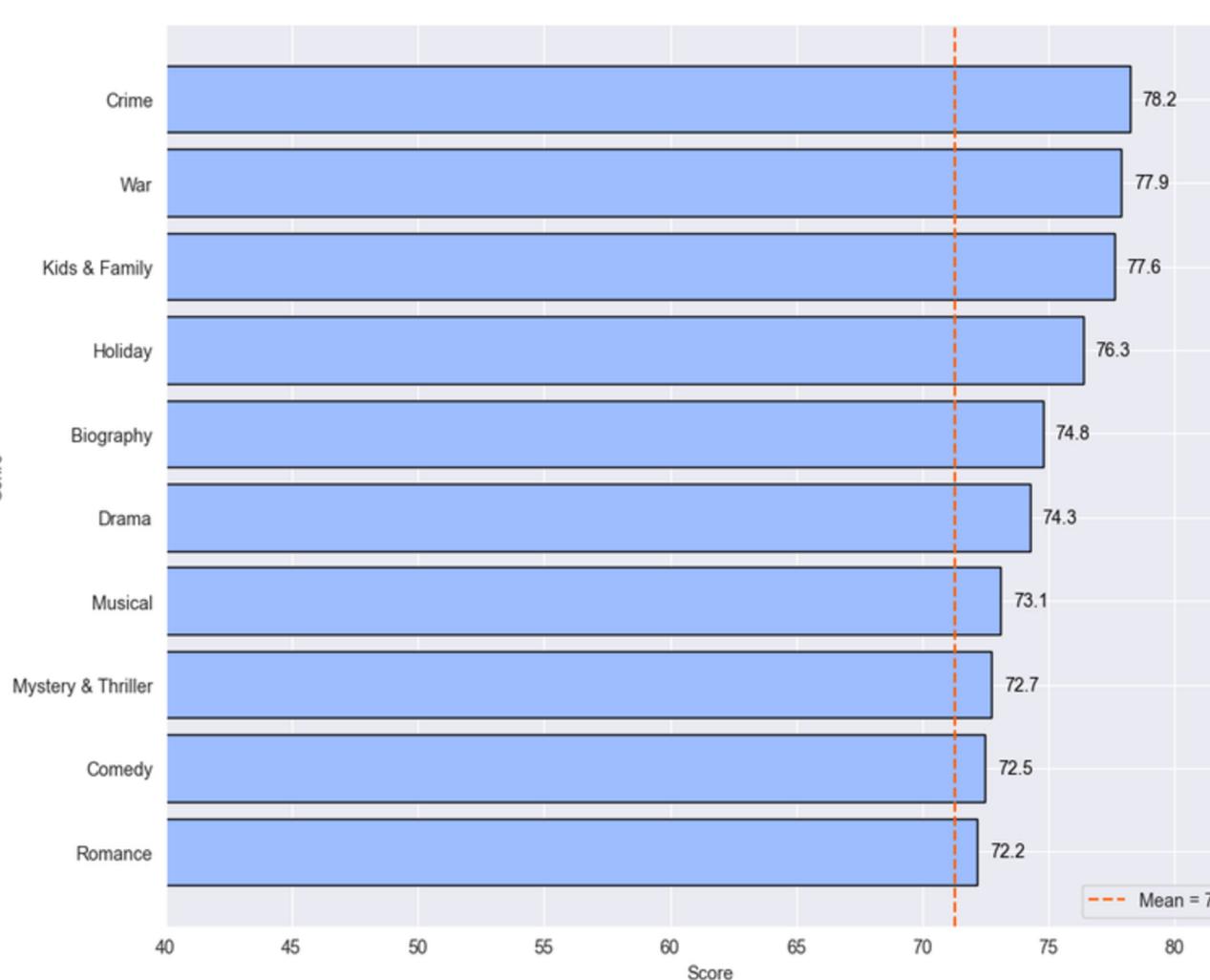
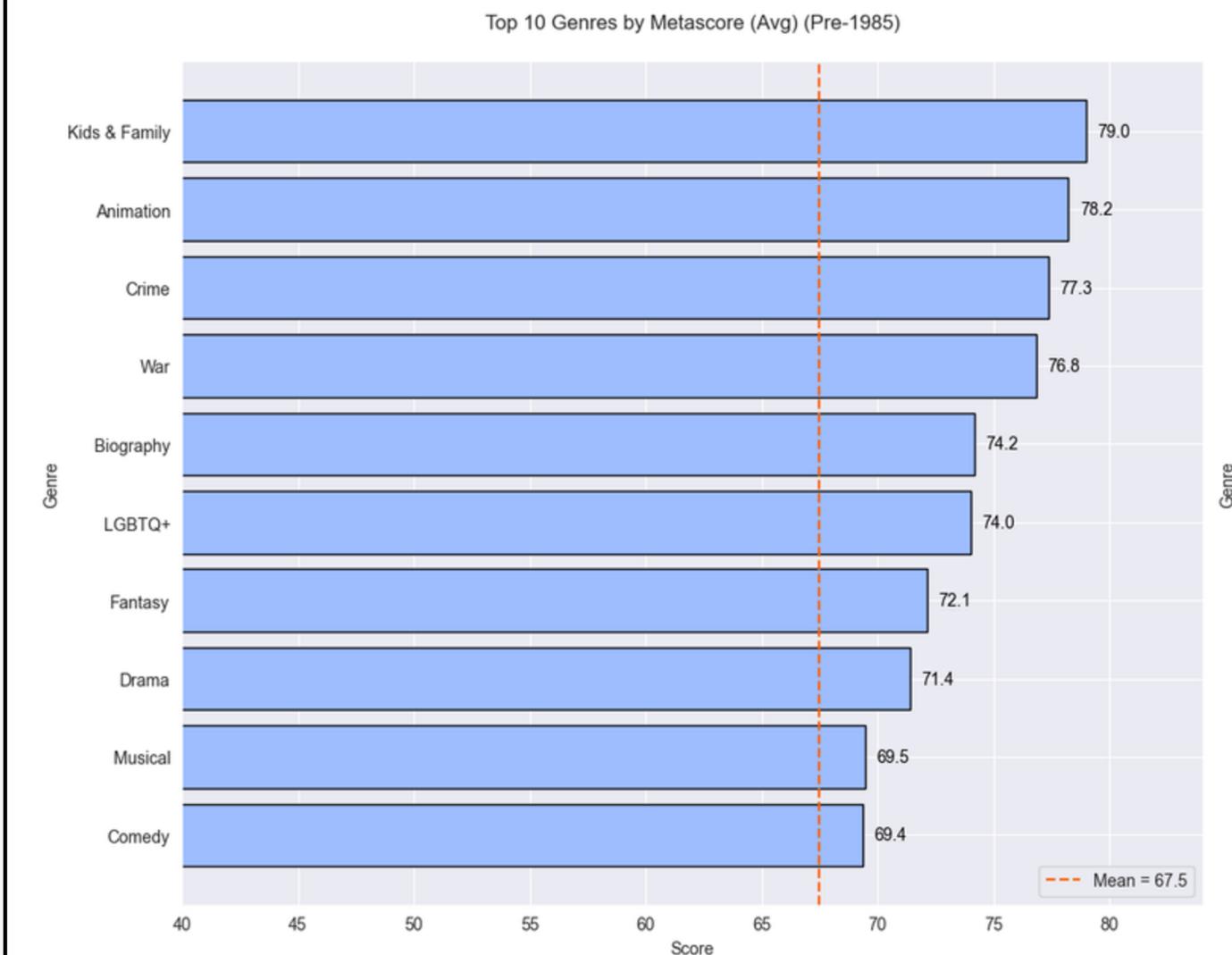
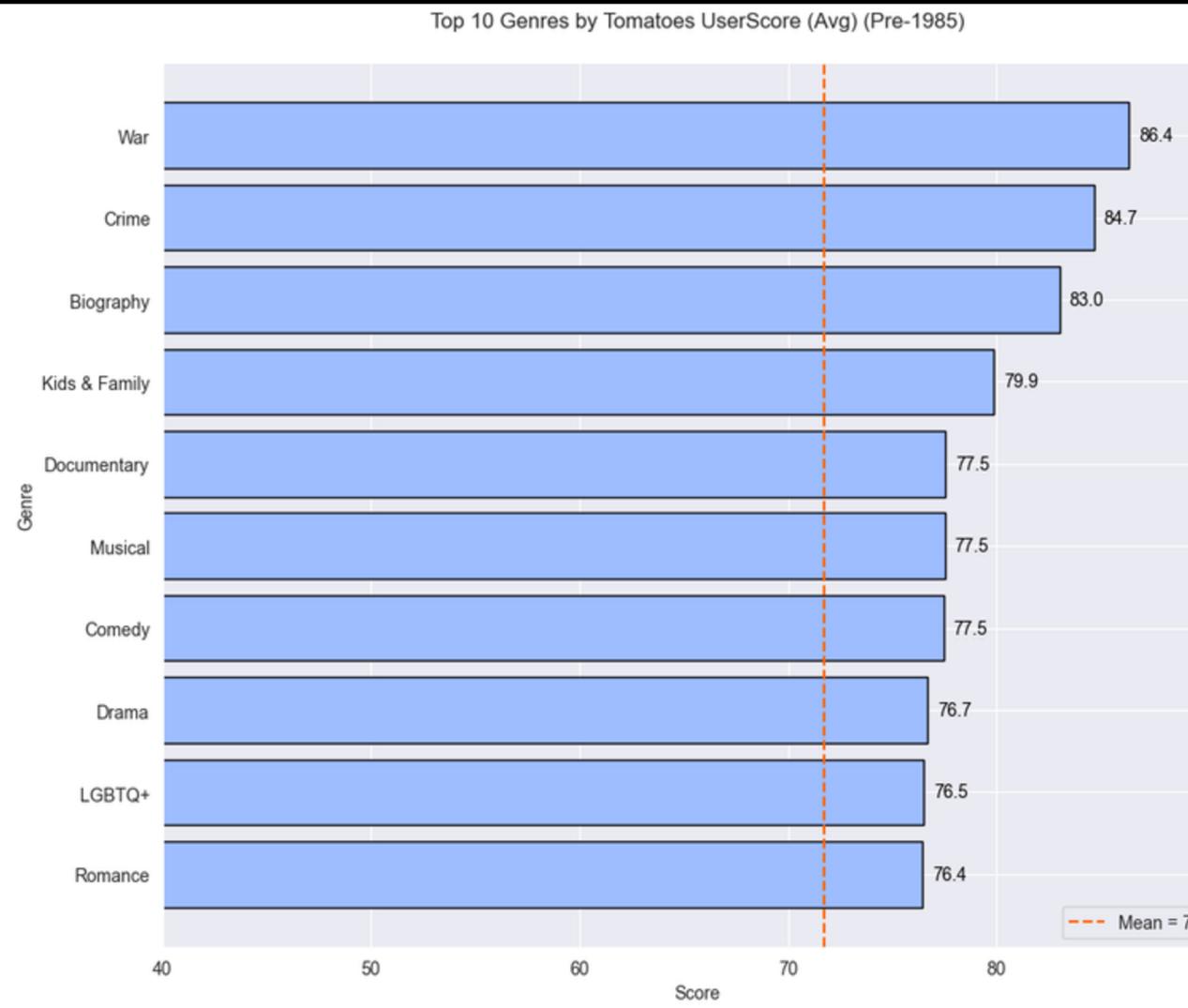
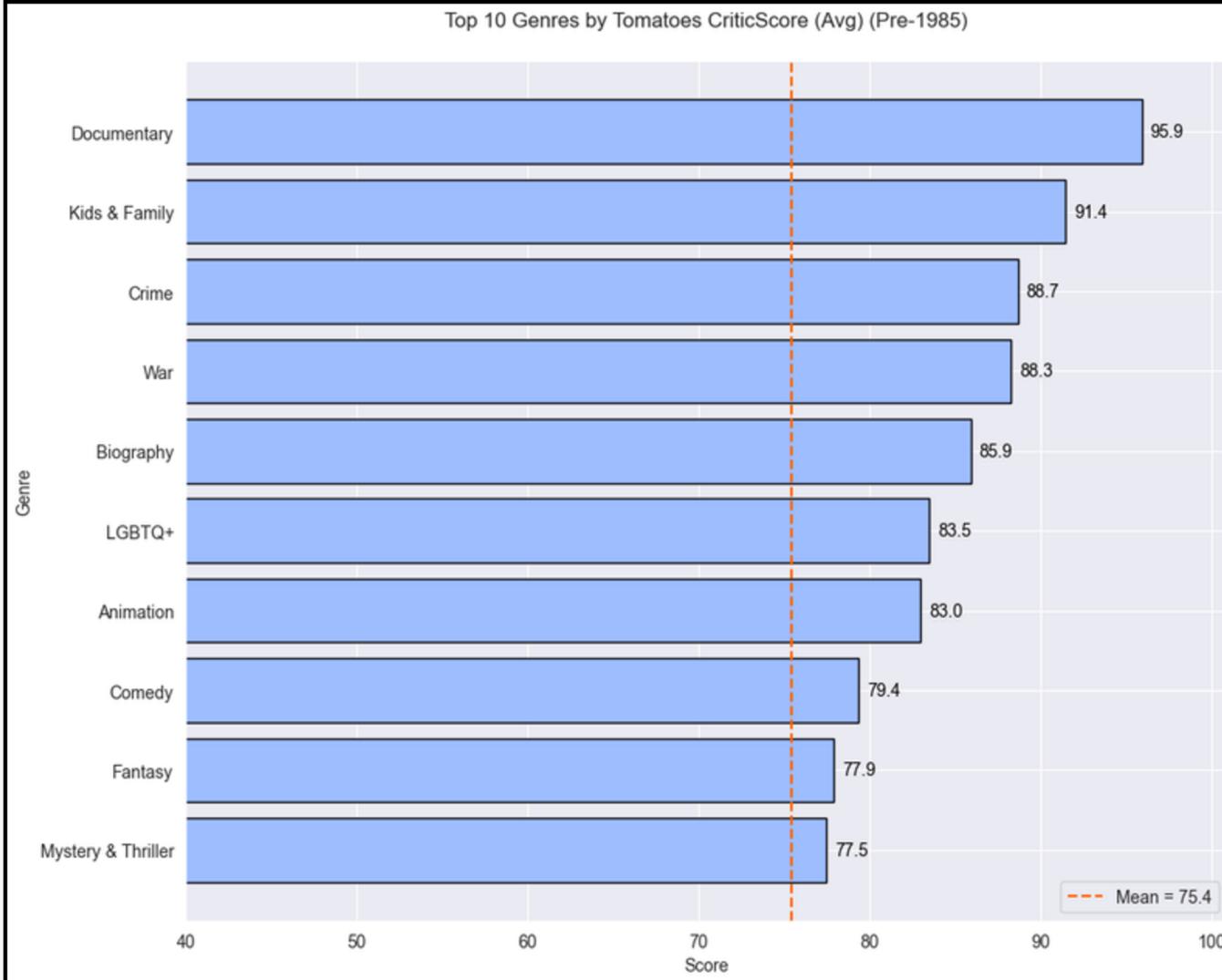
Separate rows with  
multiple genres

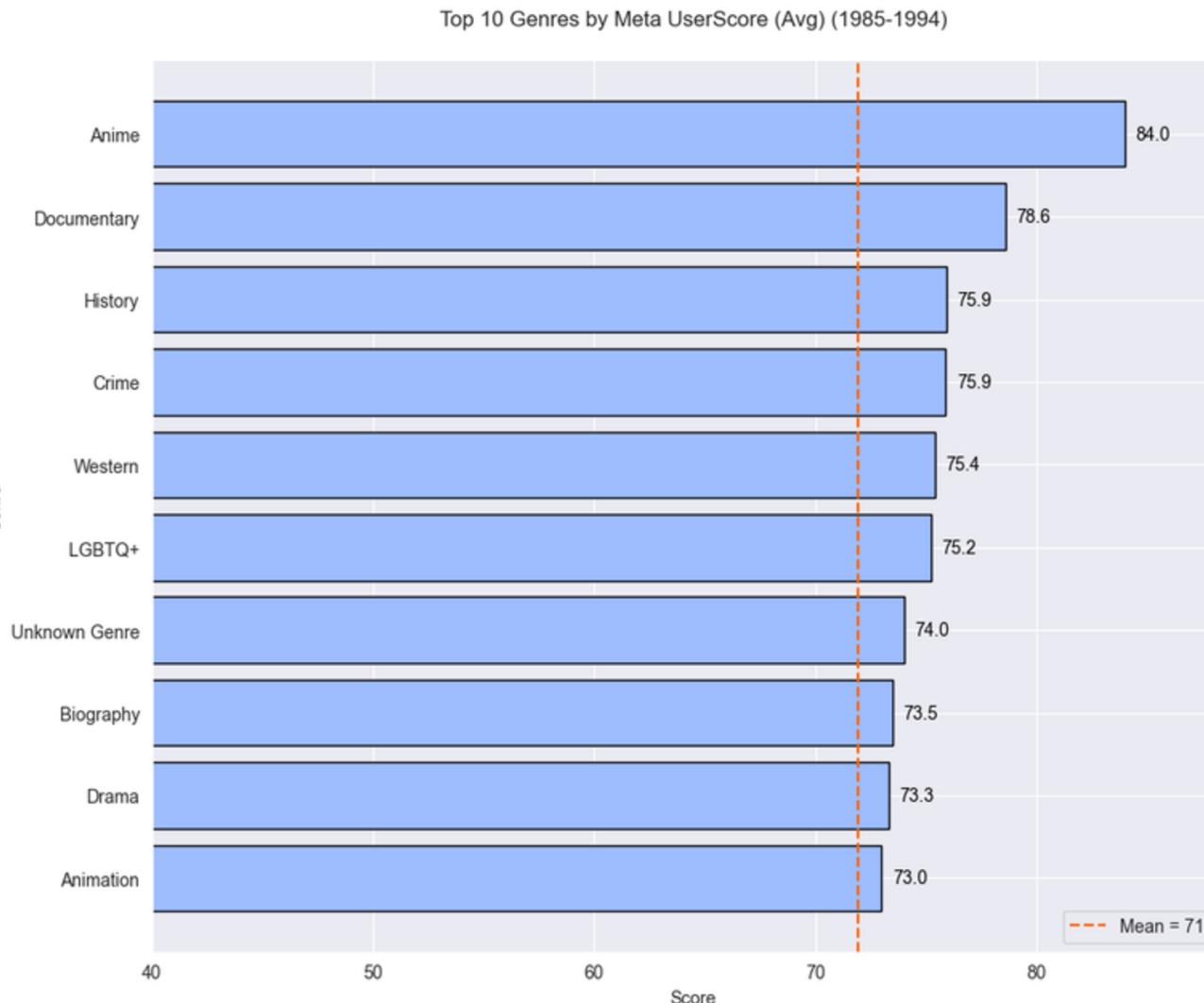
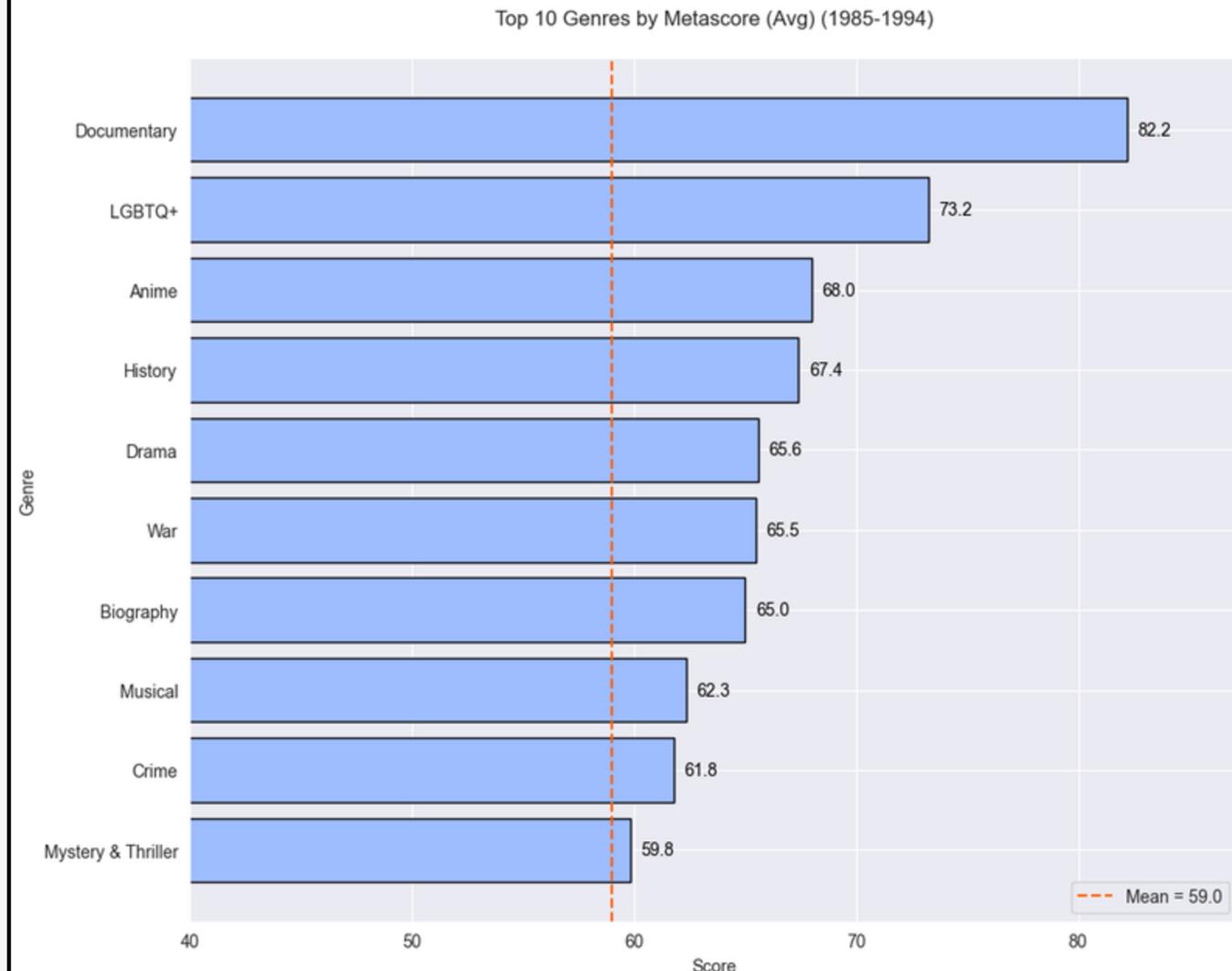
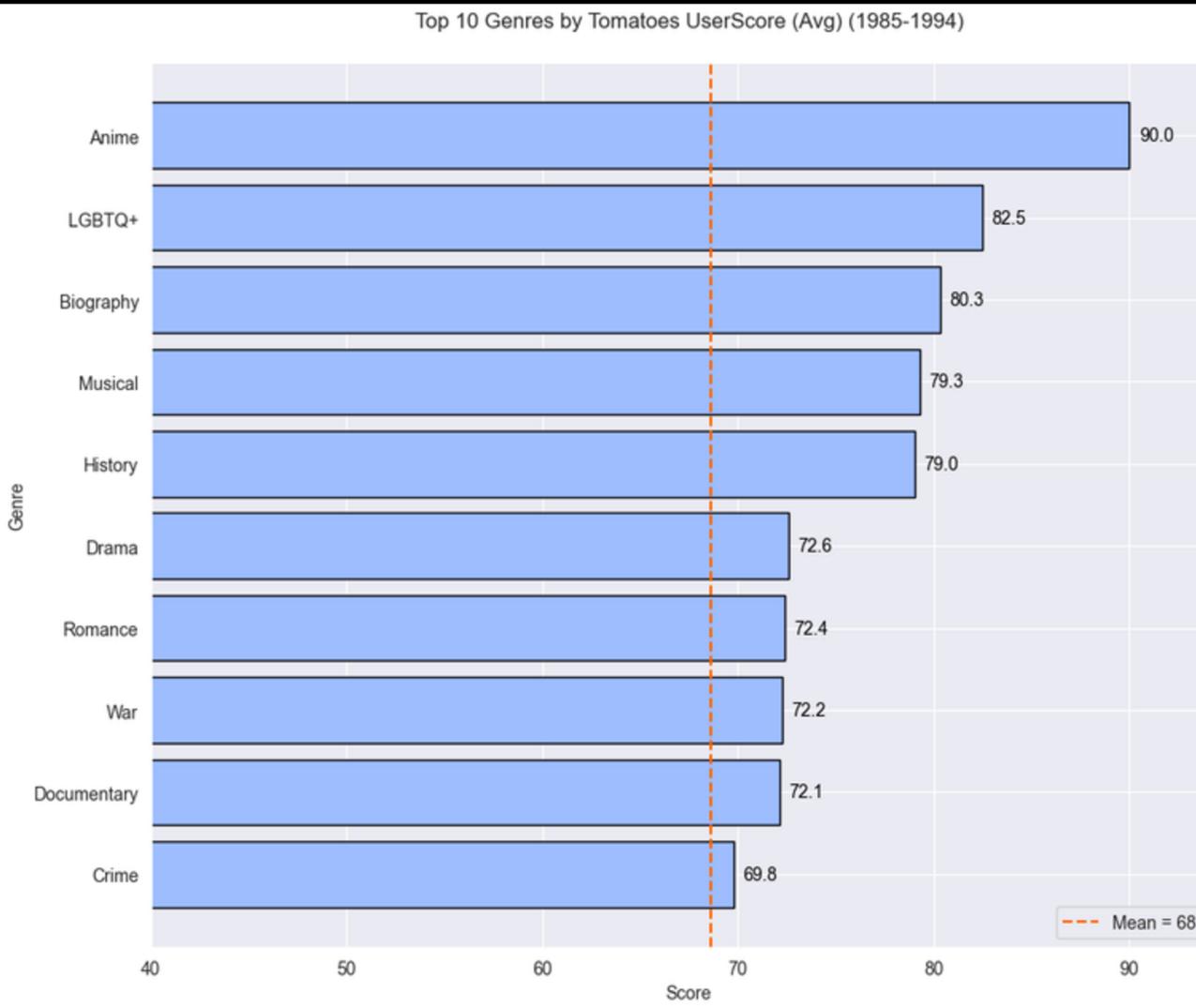
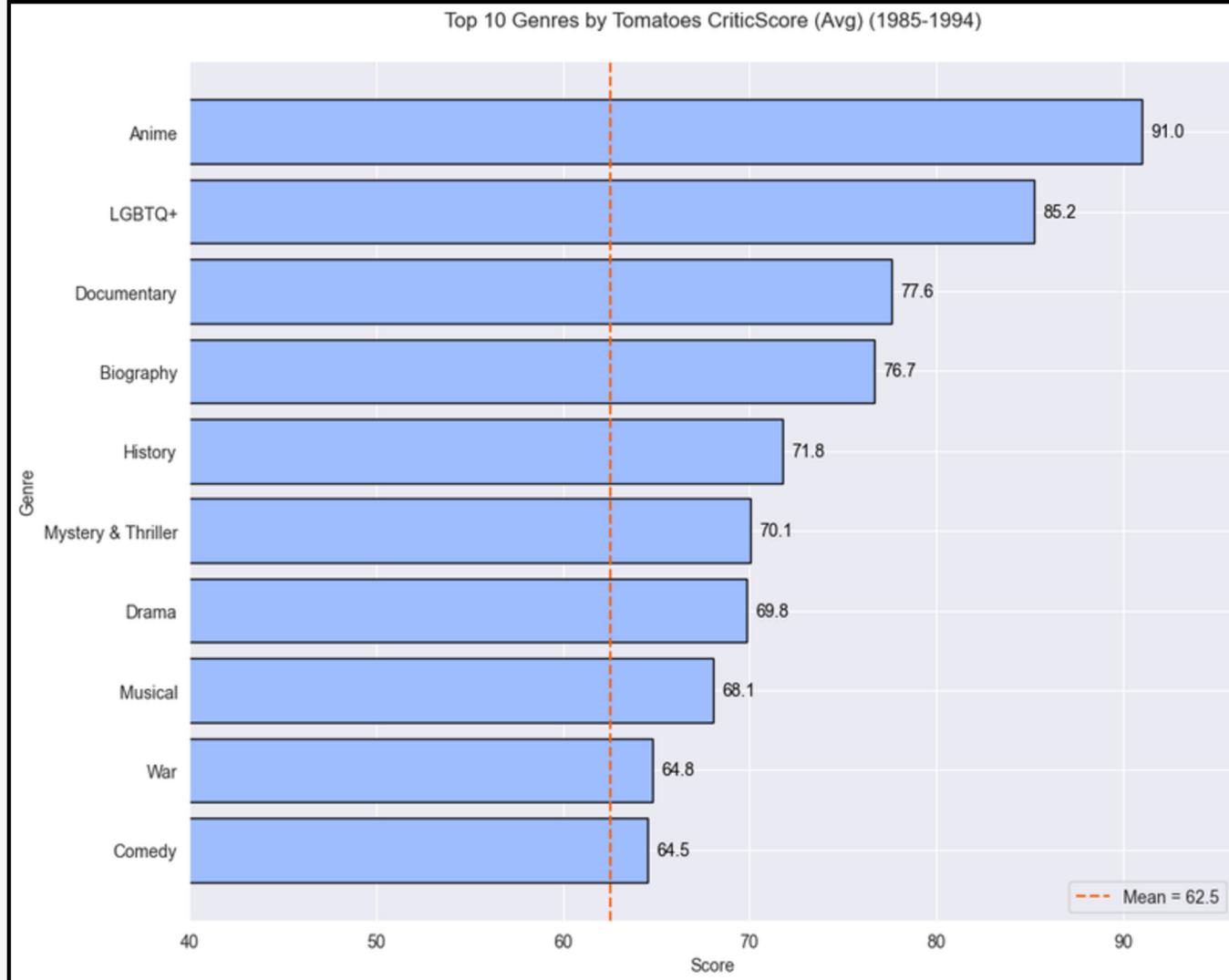


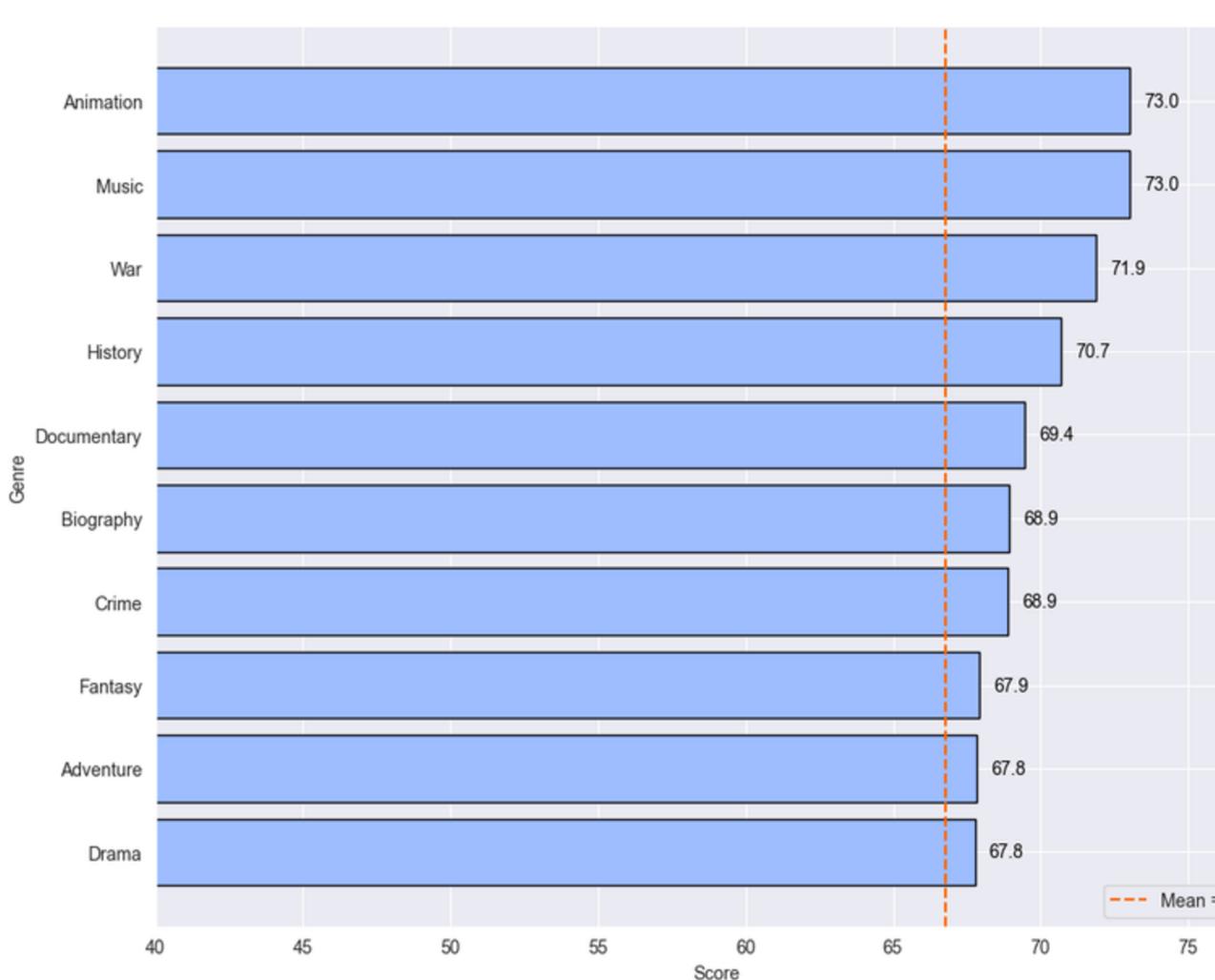
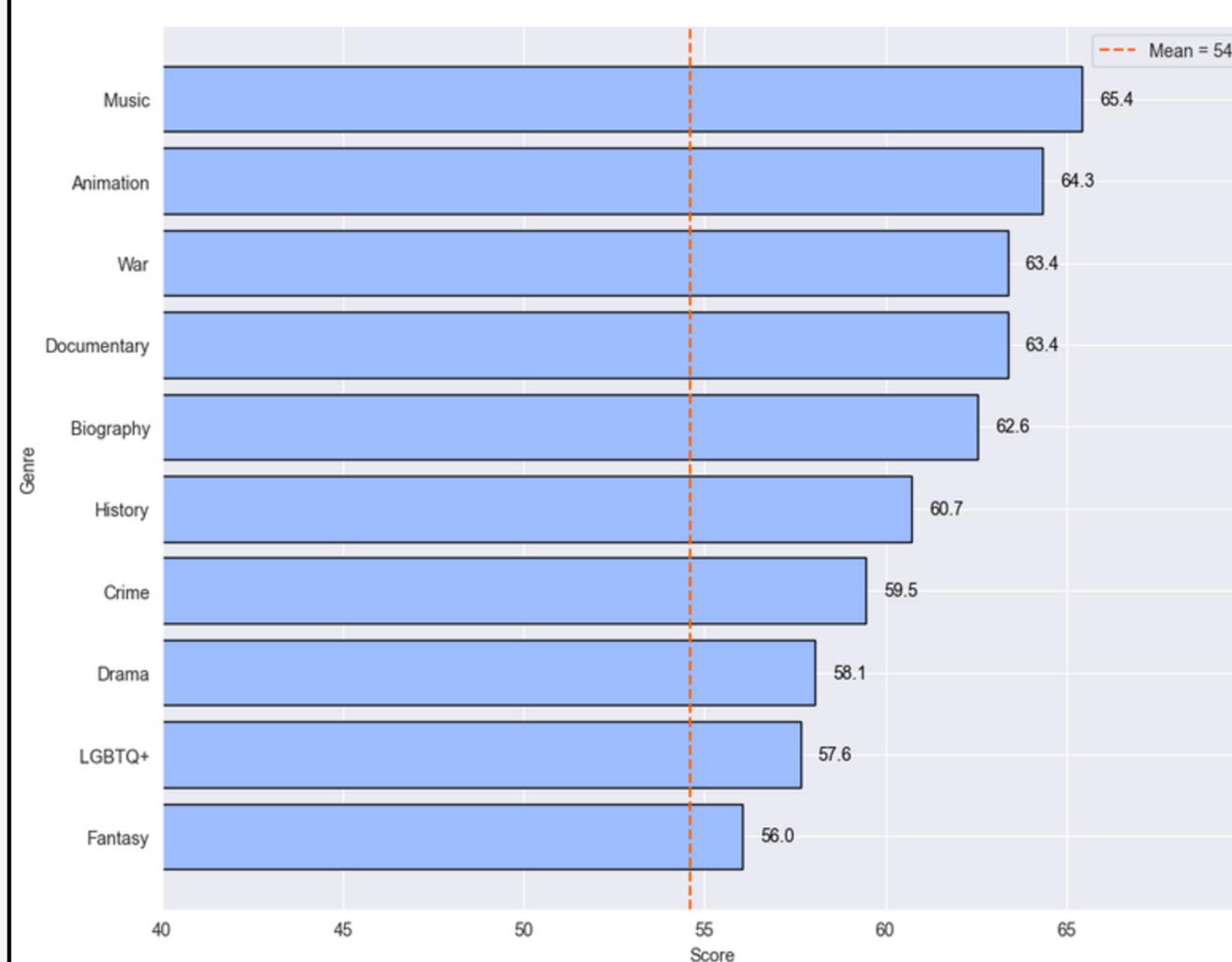
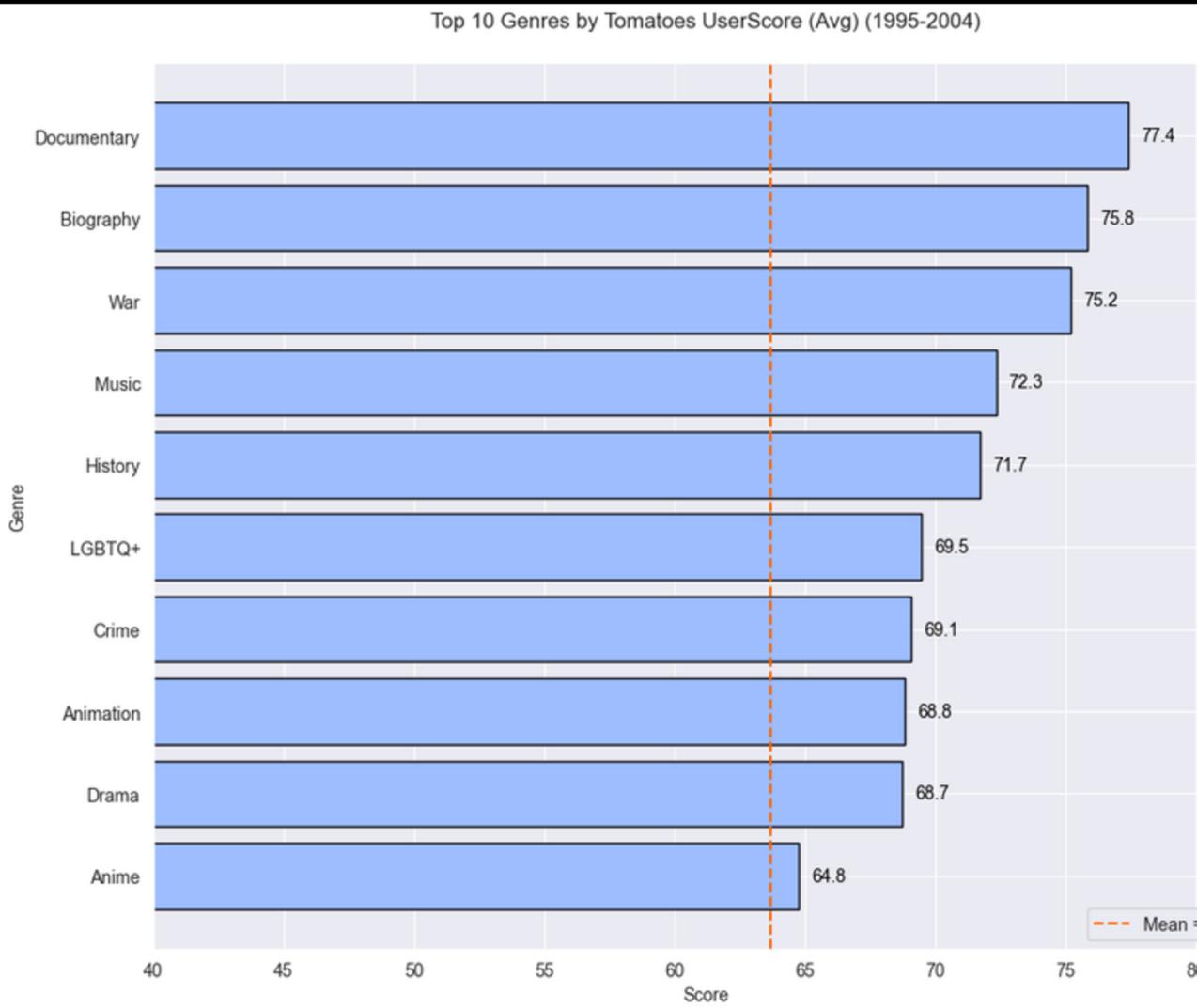
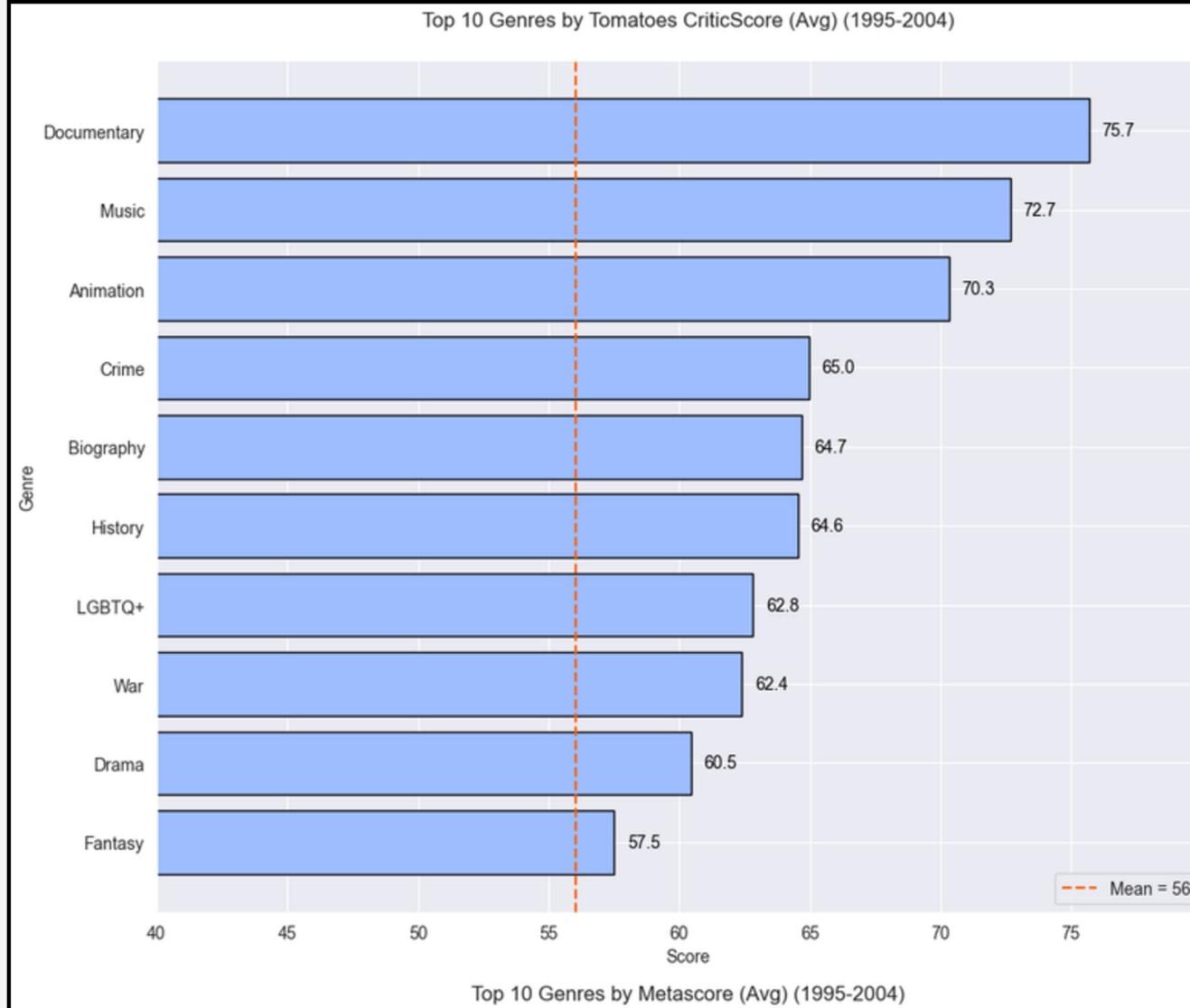
Divide data into different  
"time period" bins

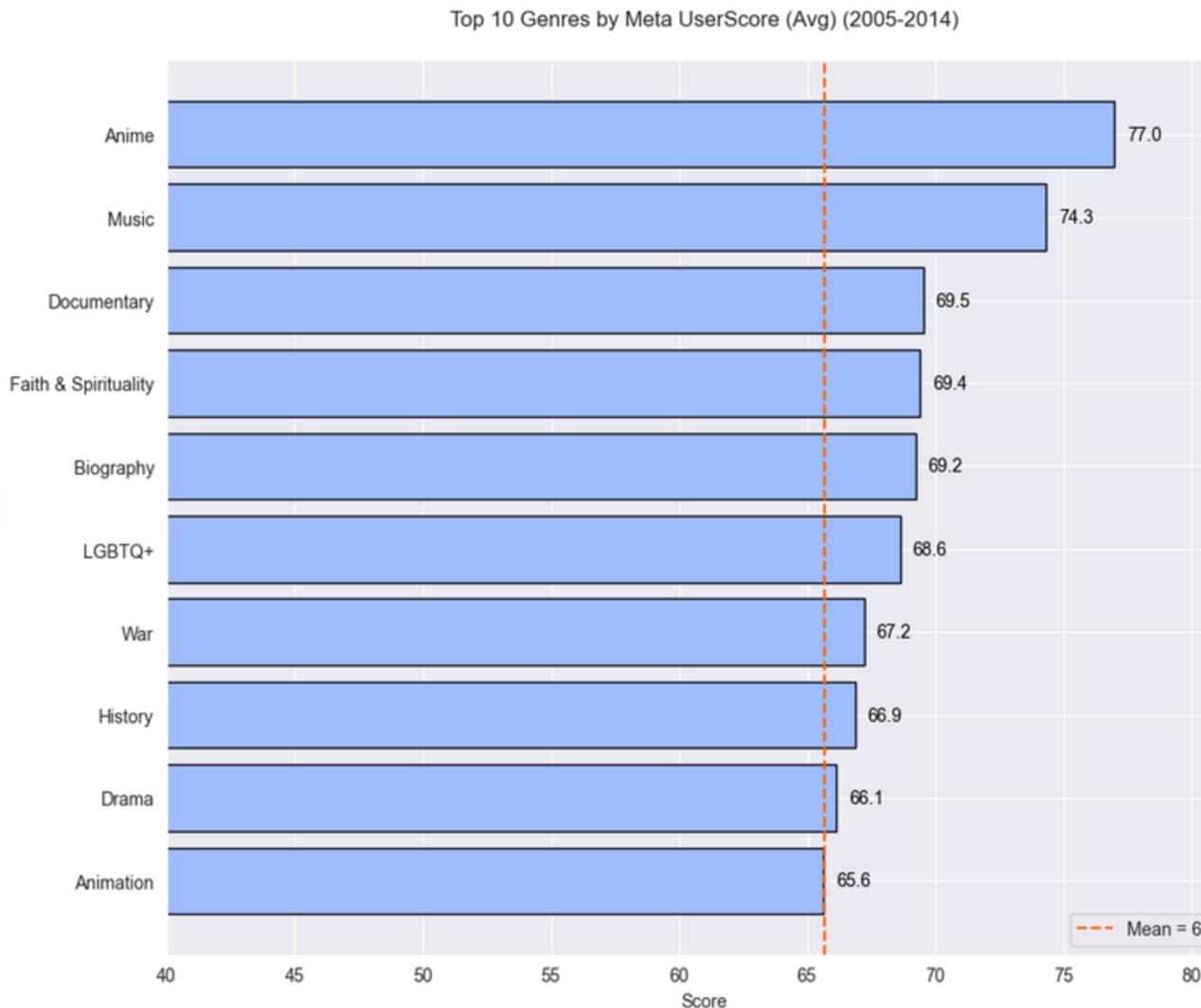
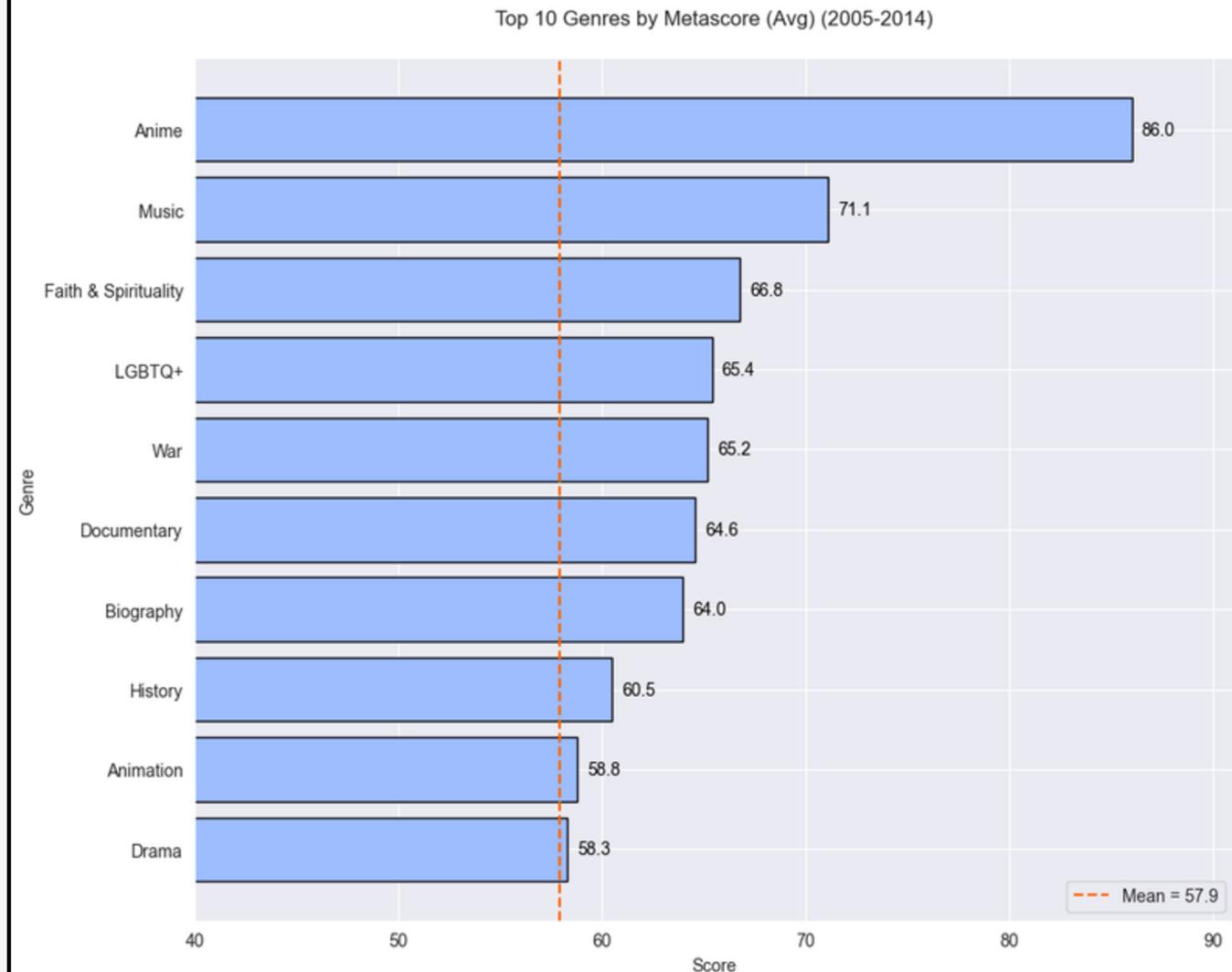
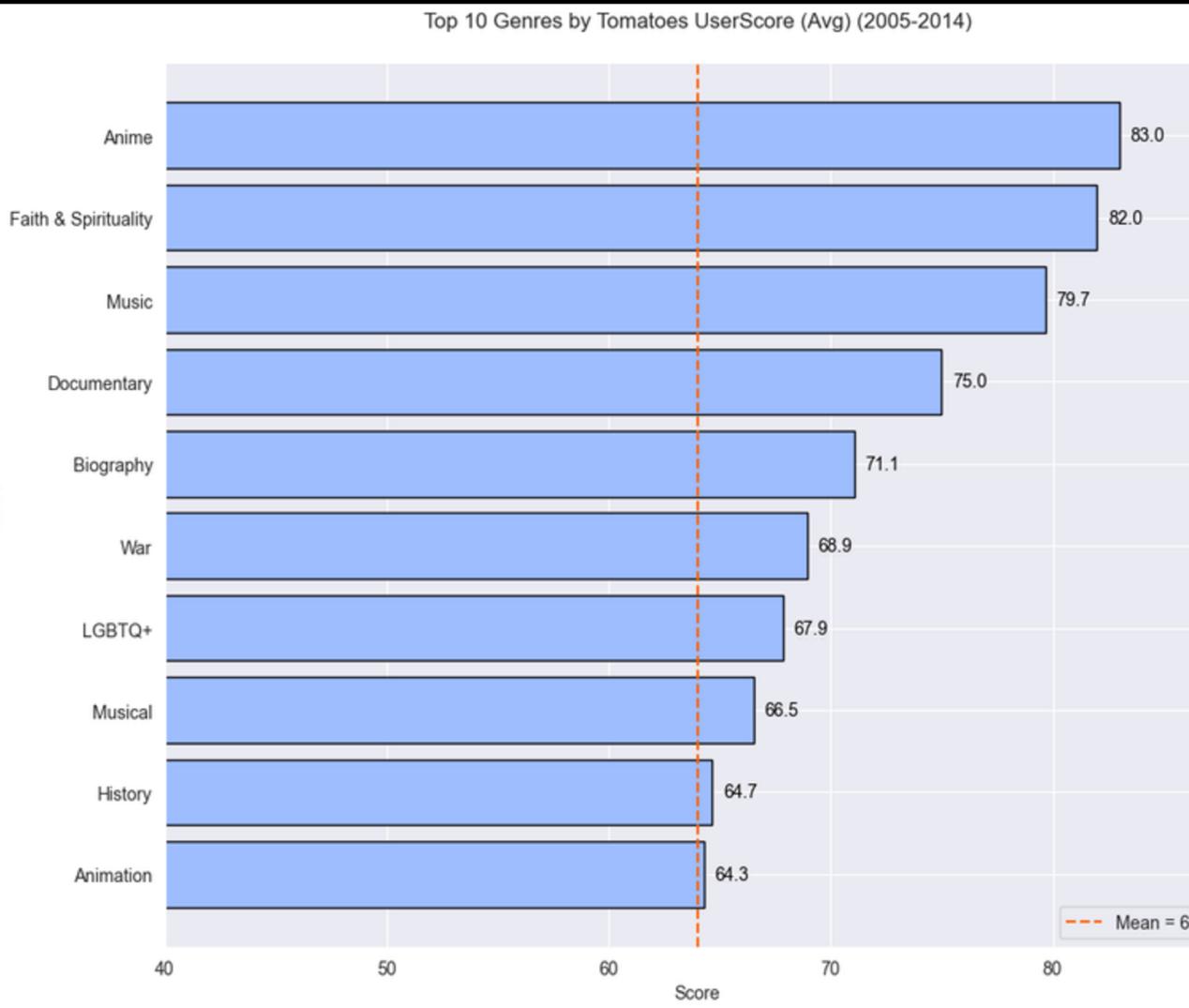
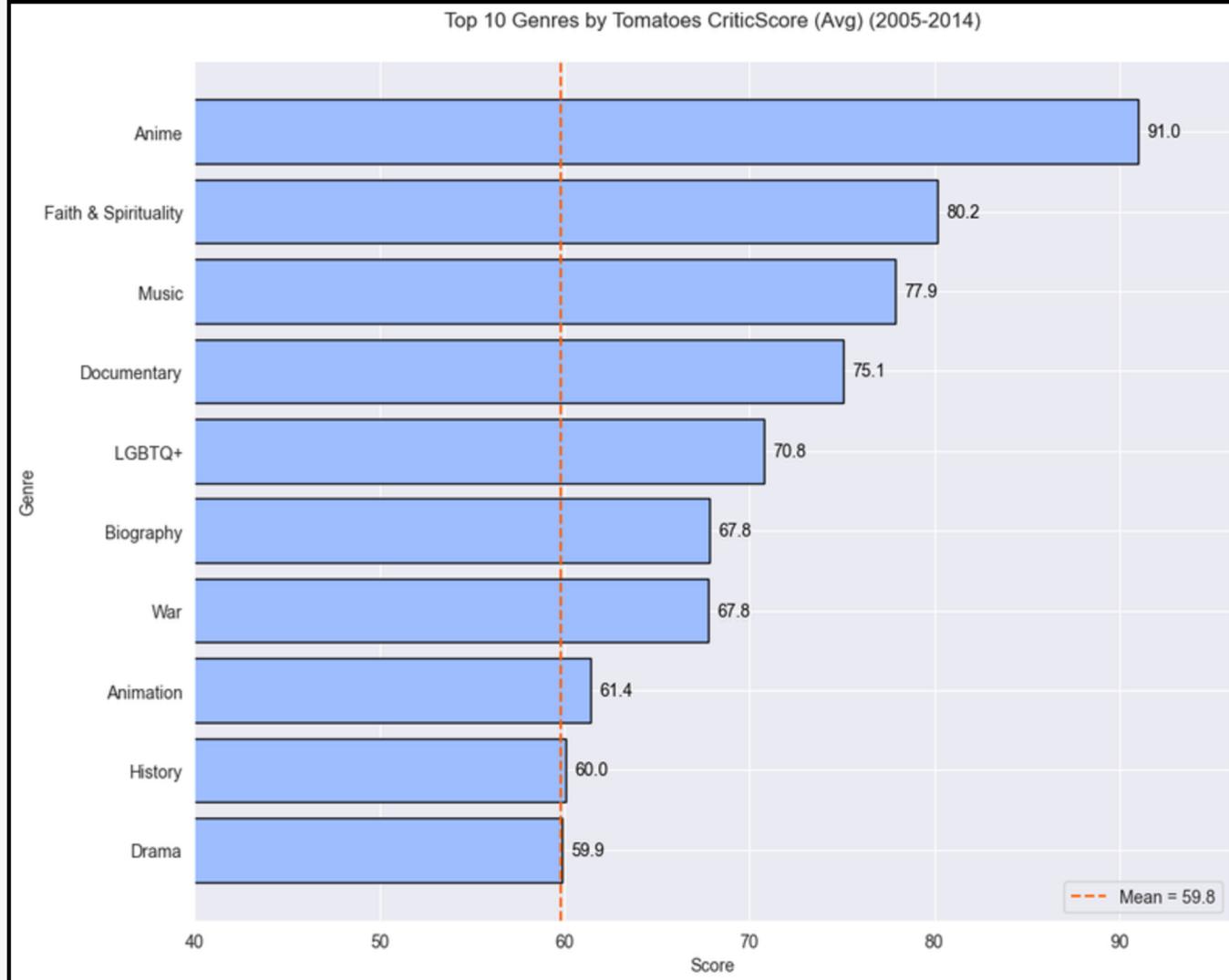


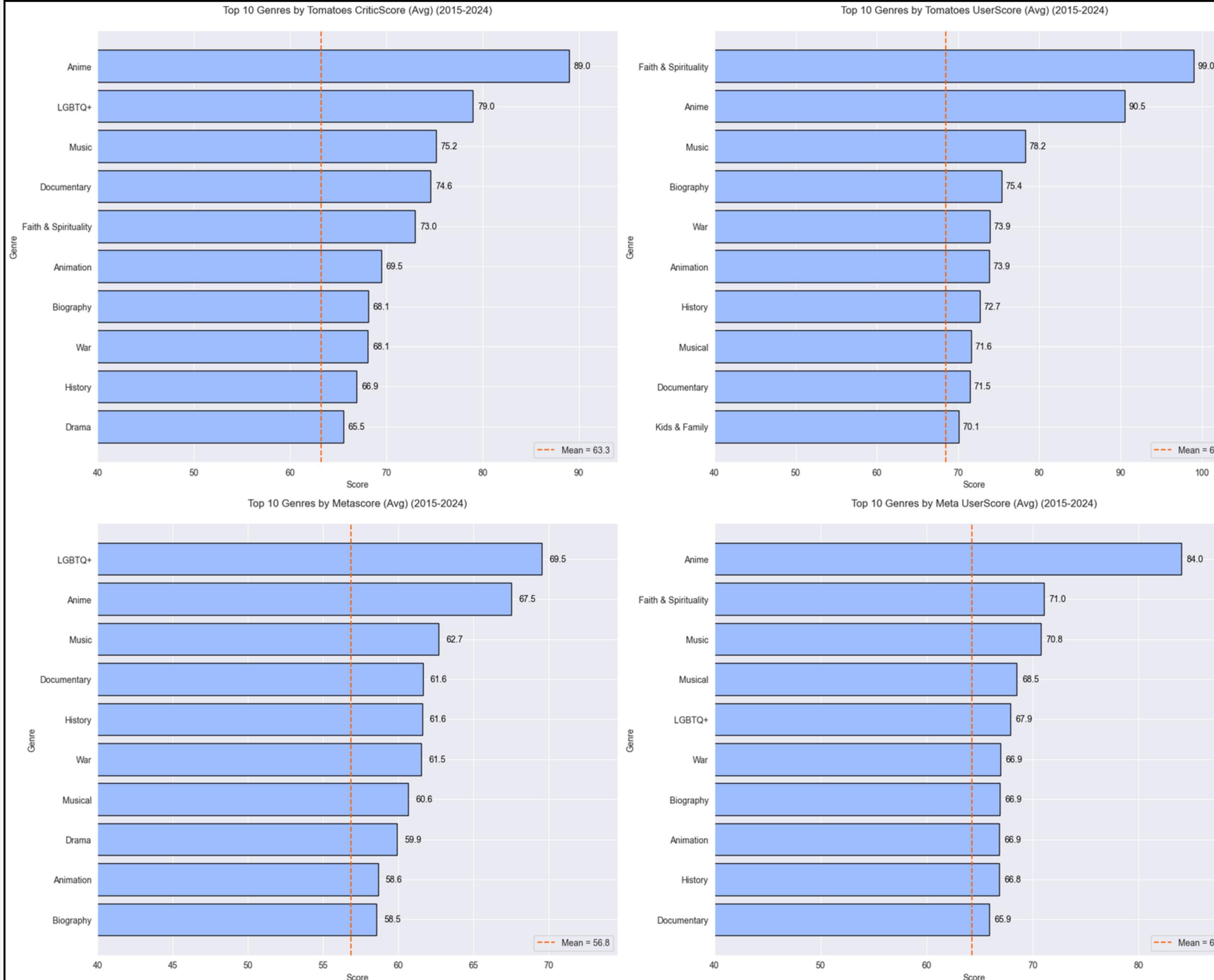
Calculate mean  
score by genres



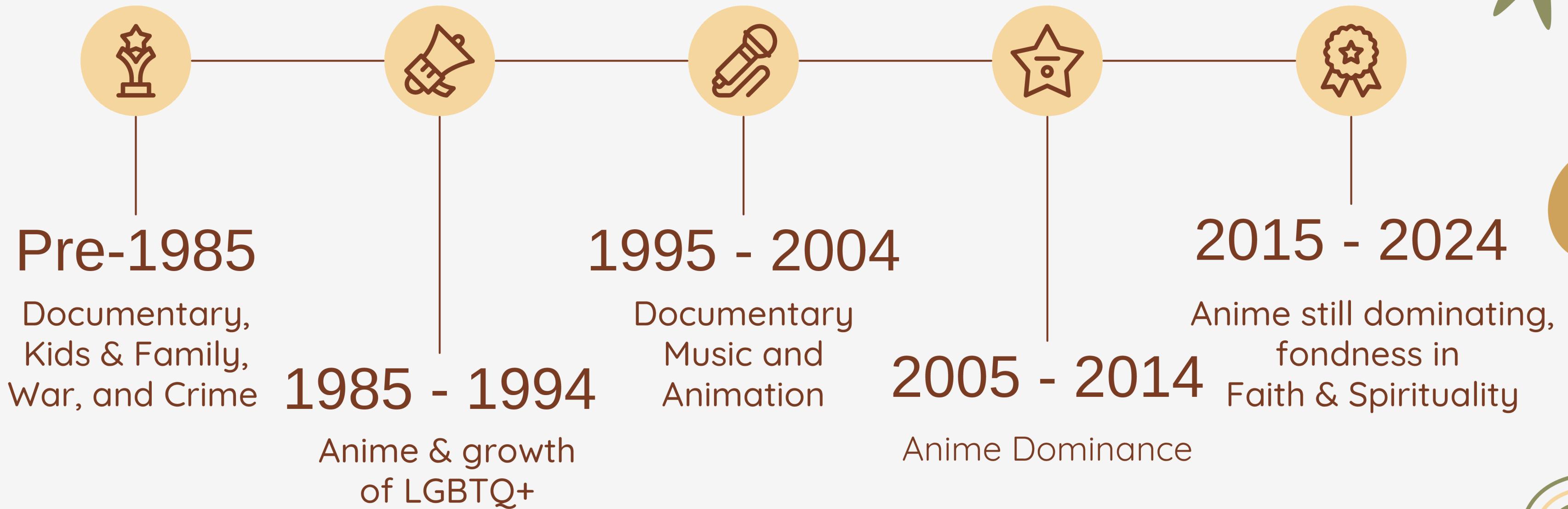






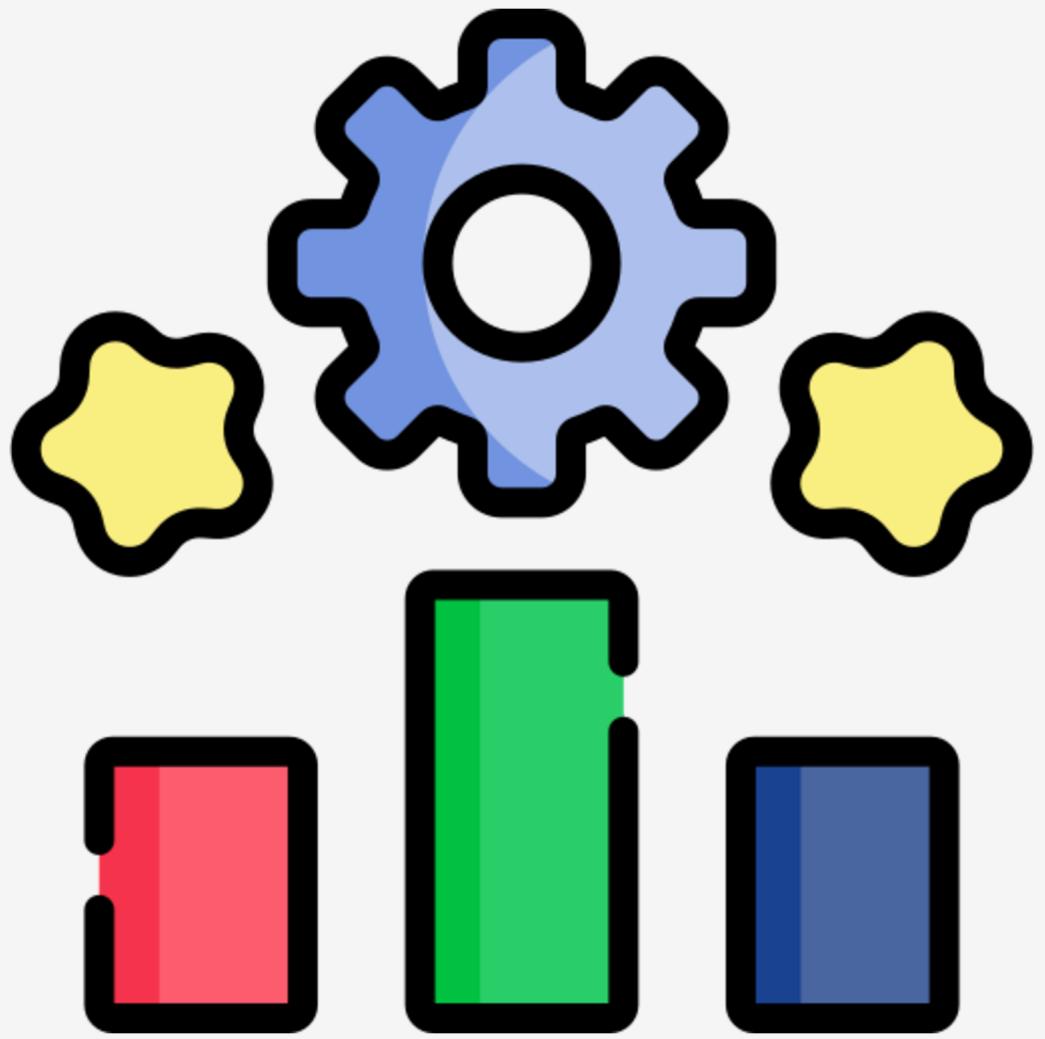


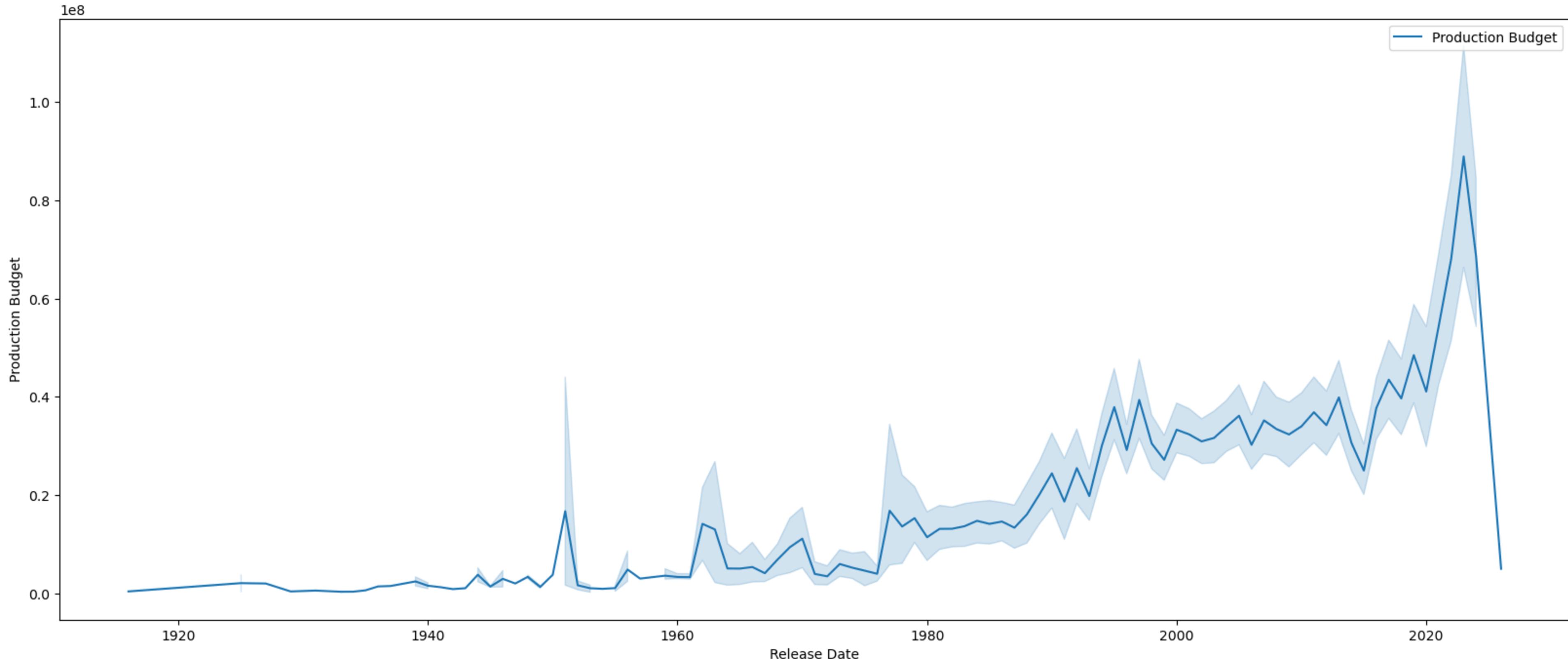
# Conclusions

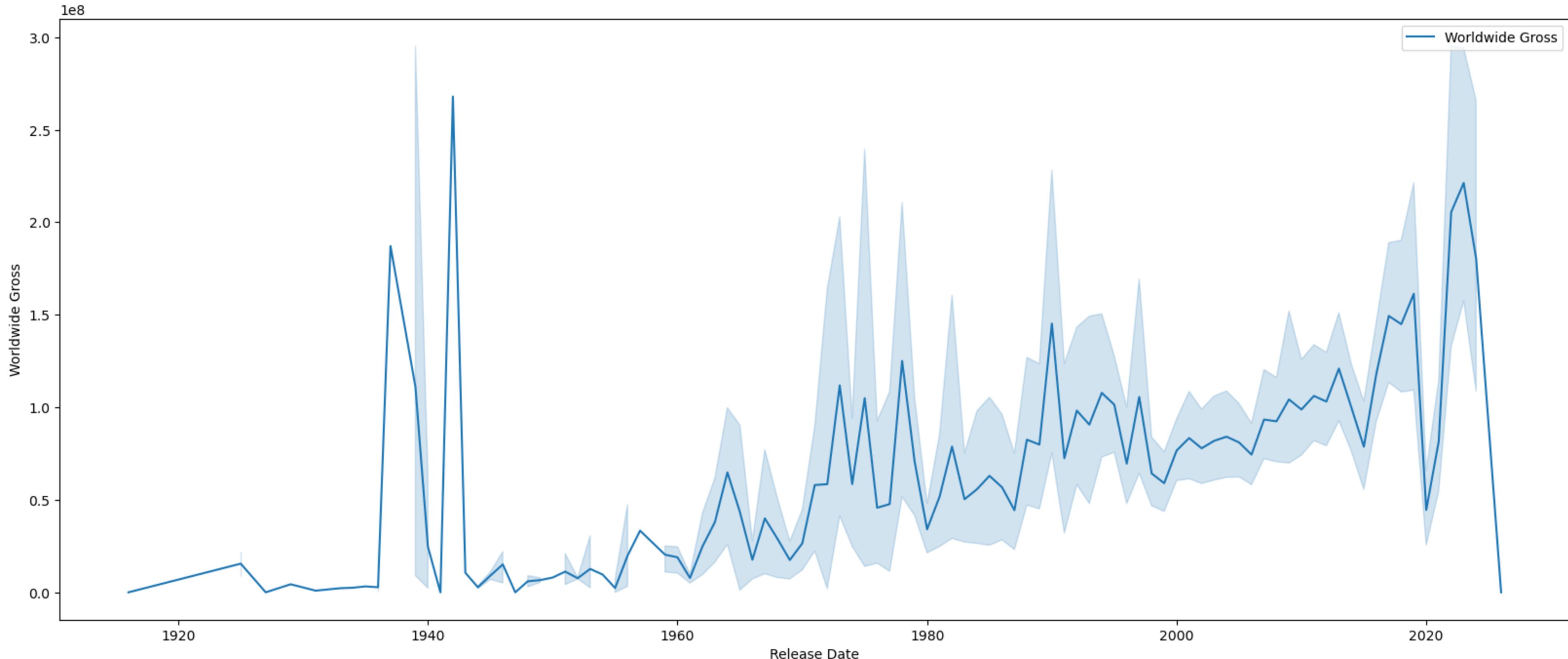


# Question 4:

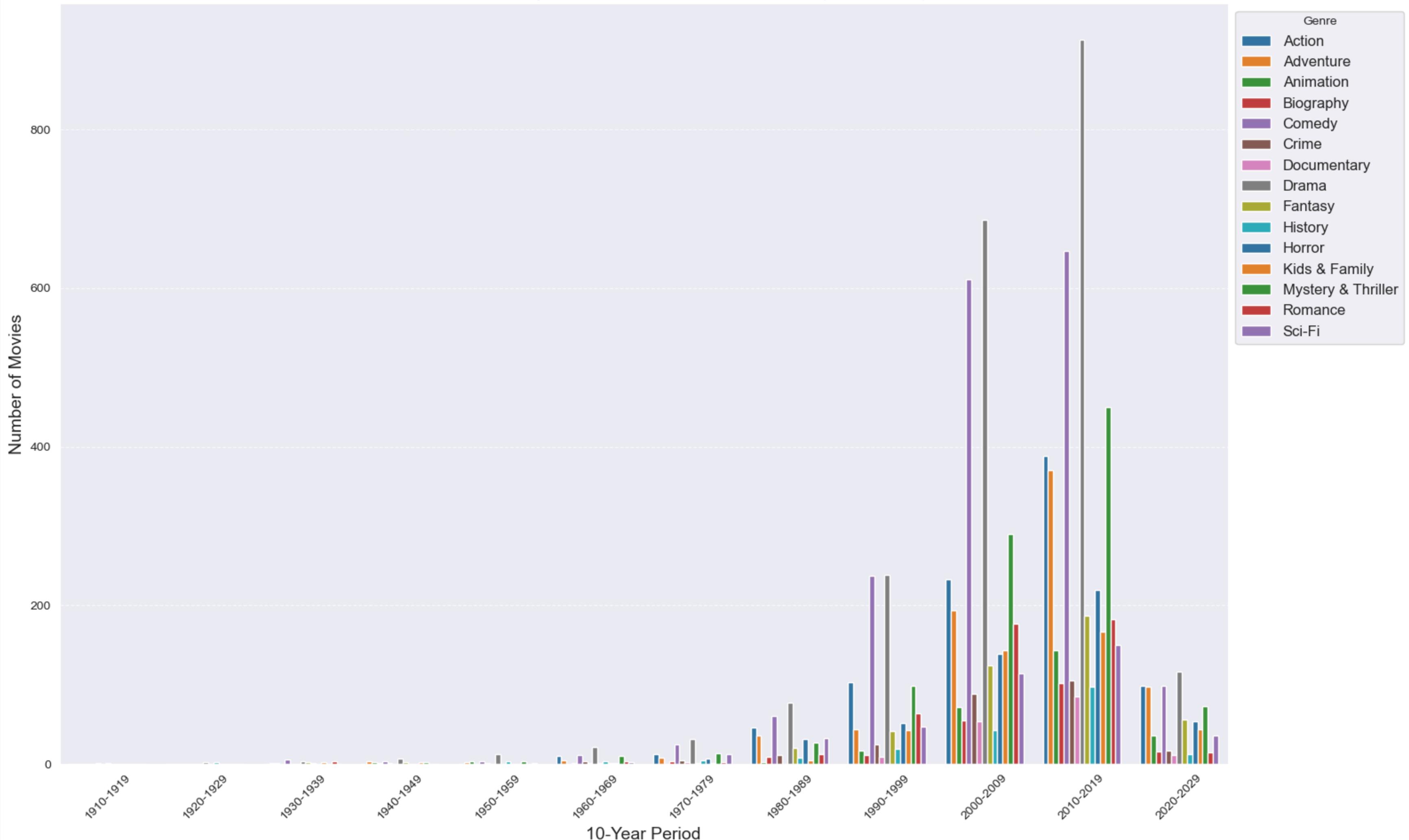
Are there trends in movie production and performance over time?







## Number of Movies by Genre in Each 10-Year Period (1910-2030)



Mean User and Critic Scores Over Time (1910-2030)



# Conclusion

## Financial aspect

- Higher budgets = higher revenue.
- Tech and franchises boosted growth post-1990.
- COVID-19 caused a sharp 2020 decline.

## Genre trends

- Action/Adventure: Always popular.
- Comedy: Steady demand.
- Animation/Family: Growth since 1990s.

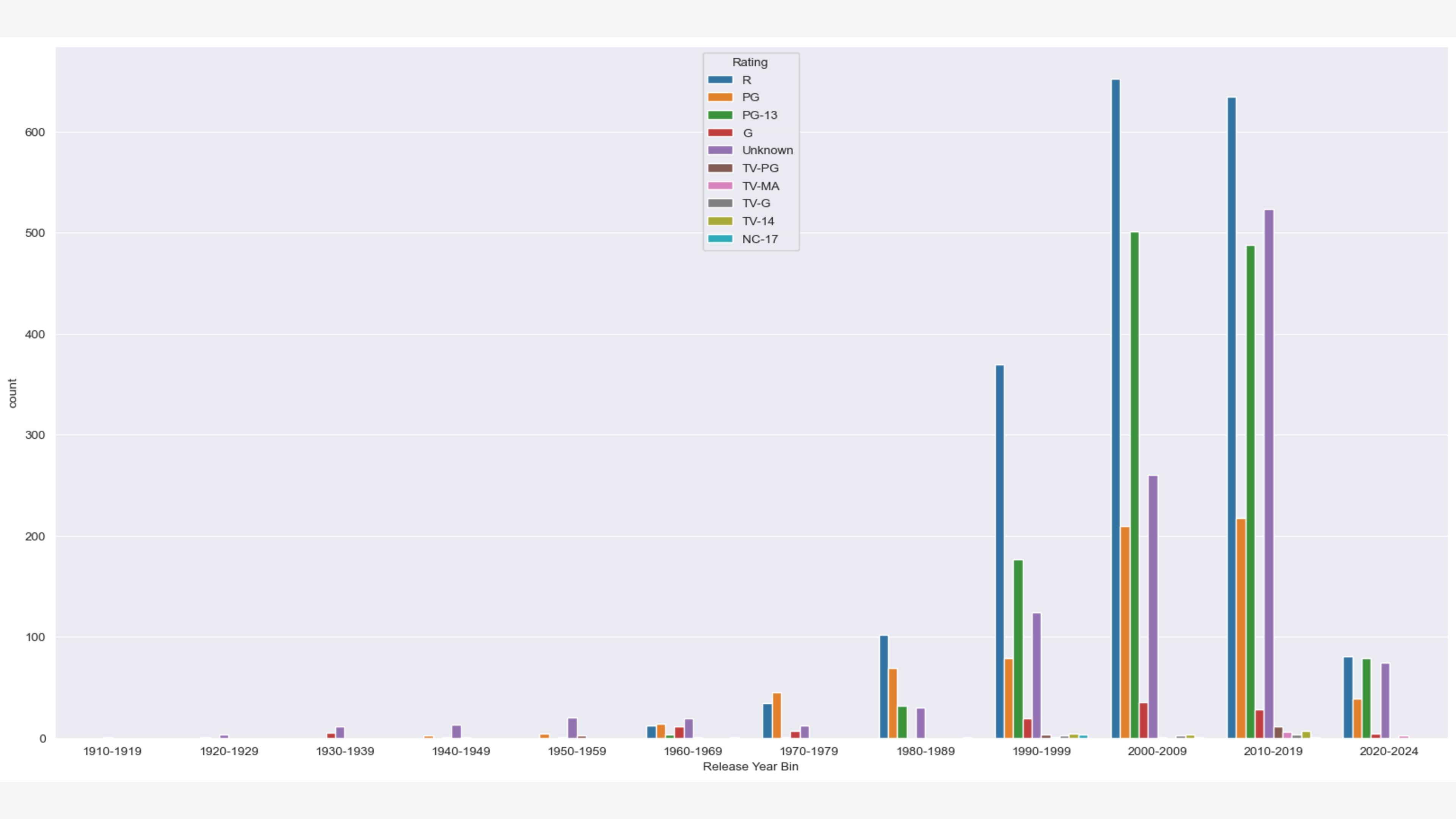
## Critic and audience ratings

- User-Critic gap widens since 1980s.
- User Scores recovering recently.
- Streaming reshapes audience ratings.

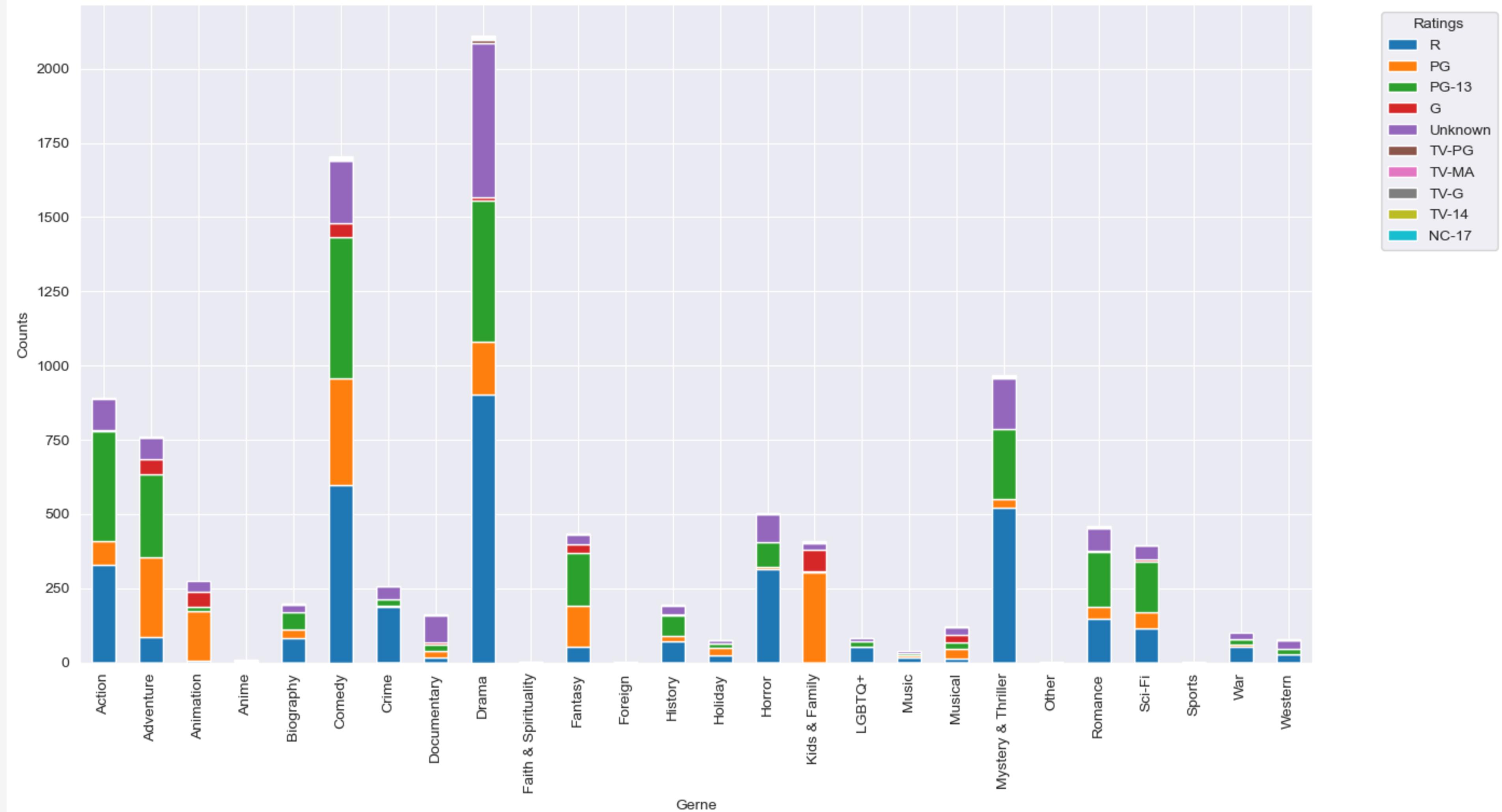
# Question 5:

## What is the trend of rating movies over time?





Distribution of Rating over Gerne



# Conclusion



Reflects changing norms and audience preferences.



MPAA ratings (1968) set age-appropriate guidelines.



Rise of R-rated films for mature themes in action, drama, and horror.



PG and PG-13 grew with family-friendly and teen-focused films.



PG-13 dominates modern blockbusters, balancing mature themes with broad appeal.

# 06

# Modelling

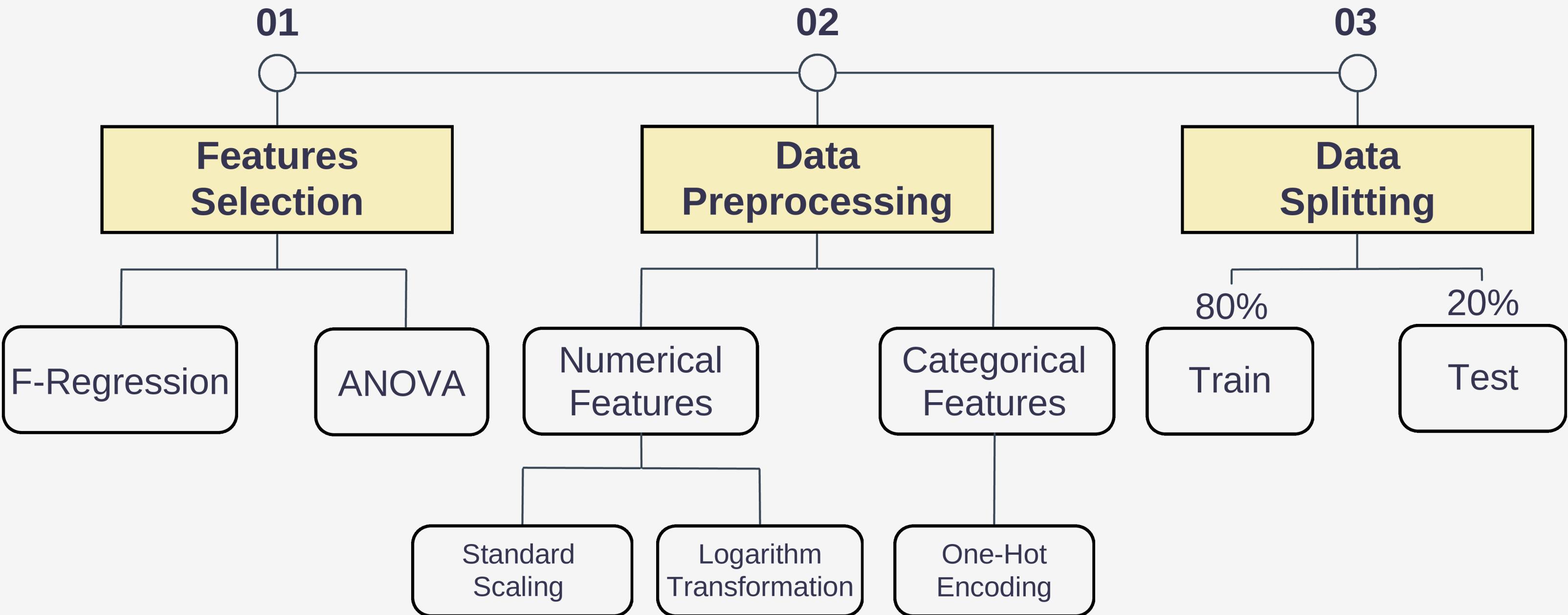


# Main Reason

Predict the gross  
of movies



# Data Preparation



# Models Training

A total of 8 models  
(including and excluding GridSearchCV)

Decision Tree

Random Forest

XGBoost

Gradient Boosting

# Models Testing

Model Name	MAE	MAE (GridSearchCV)	Mean R-Squared	Mean R-Squared (GridSearchCV)
XGBoost	0.425	0.417	0.599	0.607
Gradient Boosting	0.429	0.423	0.596	0.6
Decision Tree	0.575	0.493	0.177	0.474
Random Forest	0.425	0.441	0.567	0.573

# Models Evaluation

## Gradient Boosting vs. XGBoost (with GridSearchCV)

Model Name	MAE	Mean R-Squared	Traning Time
Gradient Boosting	0.423	0.6	35s
XGBoost	0.417	0.607	18s
Difference (%)	+ 1 %	- 1%	+ 95%

# Models Evaluation

XGBoost with  
GridSearchCV is  
the best model



# Prediction

**Movie Gross Prediction Form**

Critic Score:

User Score:

Metascore:

Meta User Score:

Budget (\$):

Genres:

Rating:

Studio:

Year:

Month:

**✓ Predict**

Prediction result will appear here



**Movie Gross Prediction Form**

Critic Score:

User Score:

Metascore:

Meta User Score:

Budget (\$):

Genres:

Rating:

Studio:

Year:

Month:

**✓ Predict**

**Predicted Total Gross: \$77,391,885.07**

# Reflection

## Challenges

- Overcoming anti-scraping measures (CAPTCHA, dynamic content).
- Handling complex data structures with multiple values in columns like Cast, Director, Genre.
- Resolving notebook conflicts on GitHub and distributing team tasks.
- Handling missing values, heavy one-hot encoding, and categorical inconsistencies.

## Lessons Learned

- Proficient in data crawling (Selenium, Scrapy) and imputation techniques (KNN, Decision Tree).
- Gained insights into data science workflow, hyperparameter tuning, and feature impacts (genres, ratings, etc.).
- Learned to create interactive form interfaces in Jupyter Notebook for smoother workflows.
- Improved collaboration, real-world problem-solving, and data cleaning skills.



Thanks  
for  
listening!