# Unveiling Movie Magic:
# Insights into Industry Trends, Audience Preferences, and Box Office Success

-

-

-

## Abstract

ABSTRACT TEST

# Contents

# 1 Introduction

## 1.1 Foreword

The movie industry, a vibrant and ever-evolving domain, offers valuable insights into audience preferences, market dynamics, and creative trends. With an abundance of data available, analyzing movies from multiple dimensions provides opportunities to uncover factors that drive both critical acclaim and commercial success. Our project investigates the intricate relationships between industry trends, movie ratings, and box office performance, providing a comprehensive view of the film landscape.

Our analysis is structured around five core themes: a general overview of movies to establish foundational patterns; industry trends to identify shifts over time; factors influencing box office revenues; determinants of IMDb ratings; and the interplay between high ratings and high box office performance. By leveraging a rich dataset encompassing genres, ratings, revenue, release years, and other critical features, this study employs both descriptive and inferential analyses to generate actionable insights.

The findings aim to benefit stakeholders across the entertainment ecosystem, from producers seeking to align their projects with audience expectations to platforms refining their recommendation algorithms. Through intuitive visualizations and data-driven observations, this report aspires to deepen understanding of what makes movies resonate with audiences while offering a lens into the dynamic forces shaping the industry.

## 1.2 The Data

### 1.2.1 Datasets

The four raw datasets selected for the study, all from Kaggle, were used to analyze movie data and provide personalized recommendations and industry trend predictions from macro movie data, user ratings, and movie-specific information, respectively. After organizing and cleaning the original datasets, they are merged into two .csv datasets for subsequent visual analysis of the movie industry.

- **The Movie Dataset**: The dataset comprises metadata for all 45,000 movies included in the Full MovieLens Dataset, which encompasses movies released on or before July 2017. The dataset includes information on the cast, crew, plot, keywords, budget, revenue, posters, release dates, languages, production companies, countries, TMDB vote counts and vote averages. Additionally, it contains files with 26 million ratings from 270,000 users for all 45,000 movies, with ratings on a scale of 1-5 obtained from the official GroupLens website.
- **Top 1000 Highest Grossing Movies:** The dataset comprises information about the 1,000 highest-grossing films produced by Hollywood studios. It has been updated to reflect the most recent data as of 25 September 2023. The data has been collated from a range of sources, including the Internet Movie Database (IMDb), Rotten Tomatoes and other similar platforms, and has been aggregated for the purpose of performing various data operations.
- **Movie Dataset: Budgets, Genres, Insights:** The movies dataset is a comprehensive collection of information about 4,803 movies. It provides a wide range of details, sourced from github.com/, about each film, including budget, genres, production companies, release date, revenue, runtime, language, popularity, and more.
- **IMDB 5000 Movie Dataset:** The dataset comprises detailed information about over 5,000 films sourced from the Internet Movie Database (IMDb). It encompasses a range of data

points, including the cast, keywords, reviews, budgets, and other pertinent information. Of particular note is the inclusion of data from the cast's Facebook pages and associated data.

**1.2.2 Data Integration Strategy**

Since the above datasets provide rich information in different dimensions respectively, this study needs to organize and merge these datasets to create a comprehensive data framework covering multiple dimensions such as basic information, ratings, box office, genres, production companies, and so on, to support a wide range of analytical needs.

For the specific steps of integration, data cleansing is first required to remove duplicates, fill in missing values, and standardize the content format, followed by merging and de-duplicating the data tables based on the movie title and IMDb ID fields as key fields. Integrating multiple datasets allows for in-depth analysis across multiple dimensions of movie industry-specific information, and is more conducive to exploring the relationship between production budgets and box office revenues.

The processed dataset is divided into the following two, focusing on macro-level information and movie-specific details, respectively. The larger dataset mainly contains basic information such as movie title, release date, duration and other basic information and production information such as movie ratings and box office budget, as well as its social media situation outlining the popularity and specific performance of the movie. This data is mainly used to analyze the overall trend of the movie industry from a macro point of view, to explore the specific performance of movies in different periods and genres, and to better understand the industry development trend and user preferences. The other data is smaller and expands on the previous dataset with information on specific box office situations, distribution companies and main actors, to more deeply analyze the impact of specific characteristics of a movie on its box office and ratings.

# 2 Tasks

This study divides the visual analysis about the movie industry into the following five sections, which analyze the movie industry in terms of information such as the overall trend of the movie industry, changes in capital investment, and key factors affecting the box office and ratings of movies.

## 2.1 General Overview of Movies
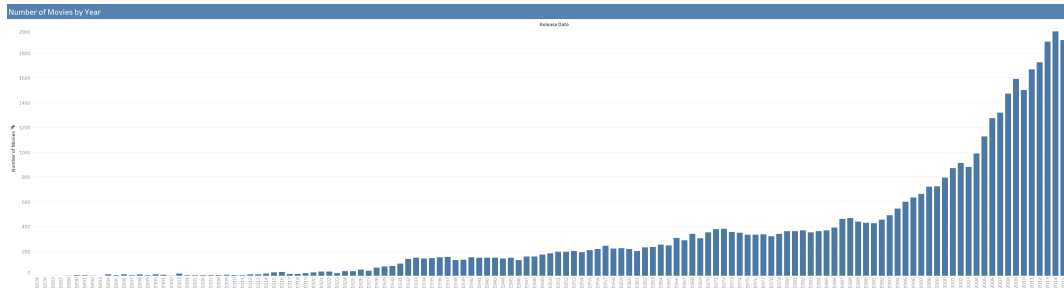
### 2.1.1 Introduction

The global film industry is vast and diverse, encompassing movies from different genres, production scales, and cultural backgrounds. To better understand the overall structure and patterns of this industry, it is crucial to first gain a comprehensive overview of the dataset. This section aims to provide insights into the foundational aspects of the movie dataset, such as the total number of films, their distribution over time, and the diversity in genres and production origins.

By exploring this general overview, we can identify key trends, such as the growth of the industry, the trend of each year. These insights lay the groundwork for the more detailed analyses that follow, offering context to the trends and patterns uncovered in subsequent sections.

Through visualizations, this section sheds light on the fundamental characteristics of the dataset, serving as the first step in unraveling the complexities of the global movie industry.

**2.1.2 Number of Movies by Year**

We use a bar chart to illustrate the number of movies each year and its trend. The x-axis represents the year and the y-axis represents the number of movies. It can be seen that the number of movies increased steadily before the 21st century, and that the increase became significantly larger after the 21st century.
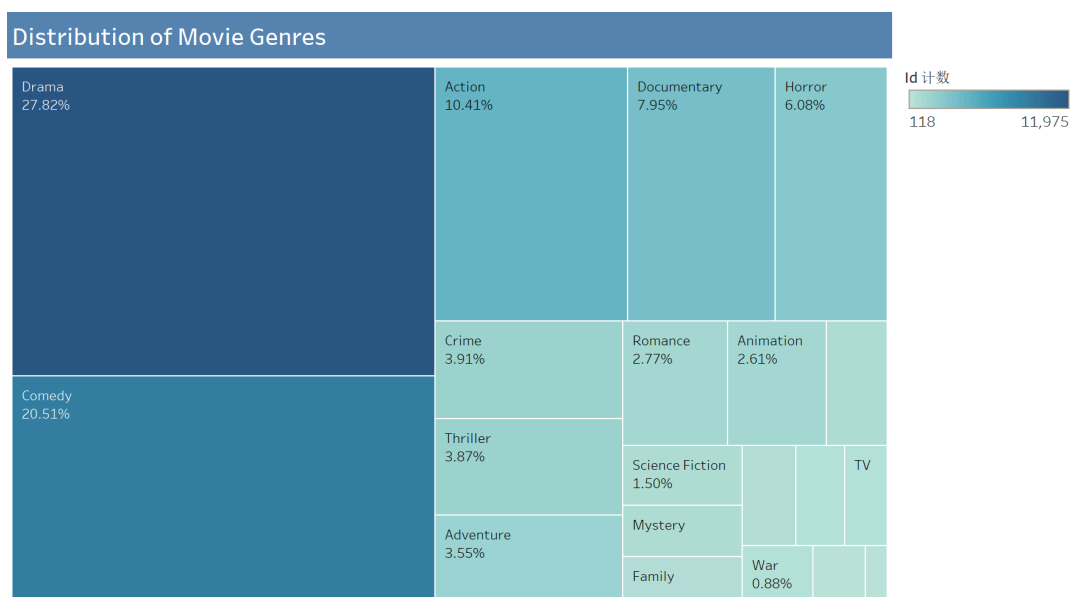


There are three main reasons causing this trend:

1. **Flat Growth Before the 21st Century:** Limited technological advancements, high production costs, and localized markets contributed to a steady but slow increase in movie production.
2. **Rapid Growth After the 21st Century:** Advances in digital filmmaking, globalization of cinema, the rise of streaming platforms, and increased accessibility significantly boosted production. Emerging markets and government incentives also played a role.
3. **Implications of Growth:** While this surge democratized filmmaking and diversified content, it raises concerns about market saturation and the balance between quantity and quality.

Additionally, we have designed this chart to be interactive, meaning that users can select a specific year to view the details of the overall movie statistics for that year in the following charts.

**2.1.3 Distribution of Movie Genres**

A tree map is employed to demonstrate the distribution of movie genres. There are possibly multiple genres for a movie, and we choose the first genre tag as the main genre of the movies. The size and the color depth represent the proportion of the genre.

We can see that Drama accounts for 27.82%, reflecting its versatility and broad storytelling range. Comedy, at 20.51%, showcases its universal appeal and ability to entertain diverse audiences. There are also smaller genres, such as Action, thriller, romance, and niche genres like horror and sci-fi, which occupy smaller proportions, appealing to specific audience groups.

It is possible that drama's flexibility to merge with other genres and comedy's accessibility are driving its popularity. However, focusing on the first genre tag may overlook the impact of hybrid genres while simplifying the analysis. Therefore, it is important to conduct further data mining and in-depth visualization.
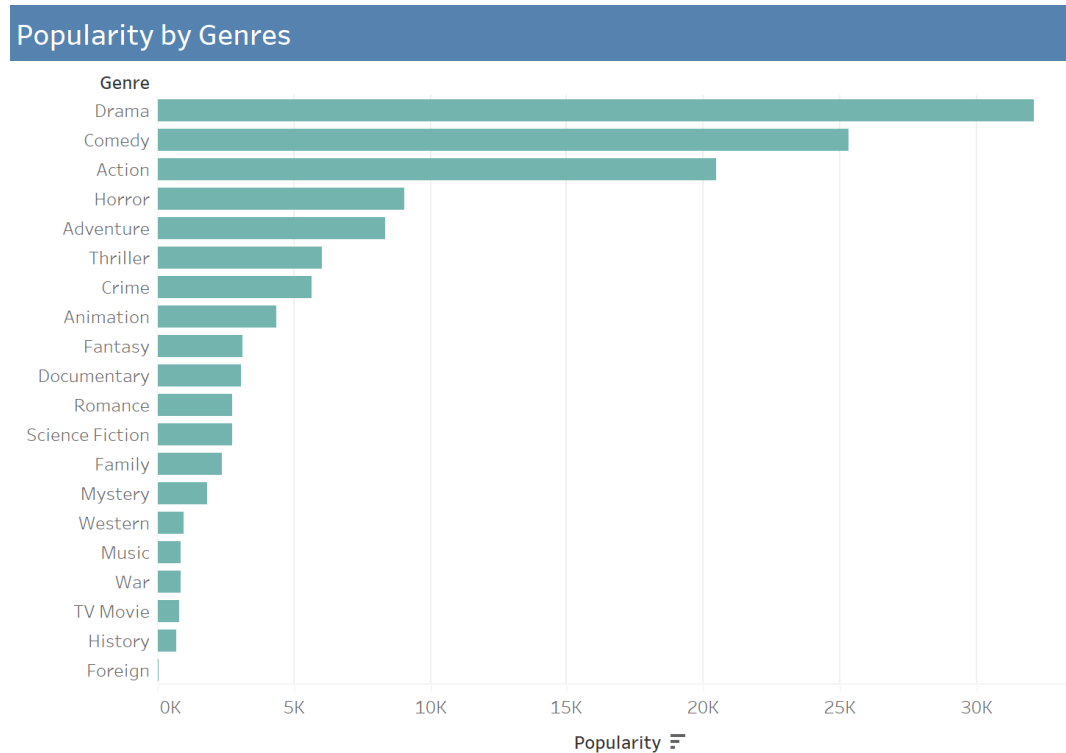
### 2.1.4 Number of Movies by Region

The table below shows the number of movies produced by different regions. The table shows the top 20 countries in terms of the number of movies produced. The United States produced the most movies (18,429), followed by the United Kingdom (3,072) and France (2,711). It is worth noting that many movies are produced by multinational joint ventures, and here we take the first field. While this allows for a simplification of data organization, it also creates some bias in our analysis.

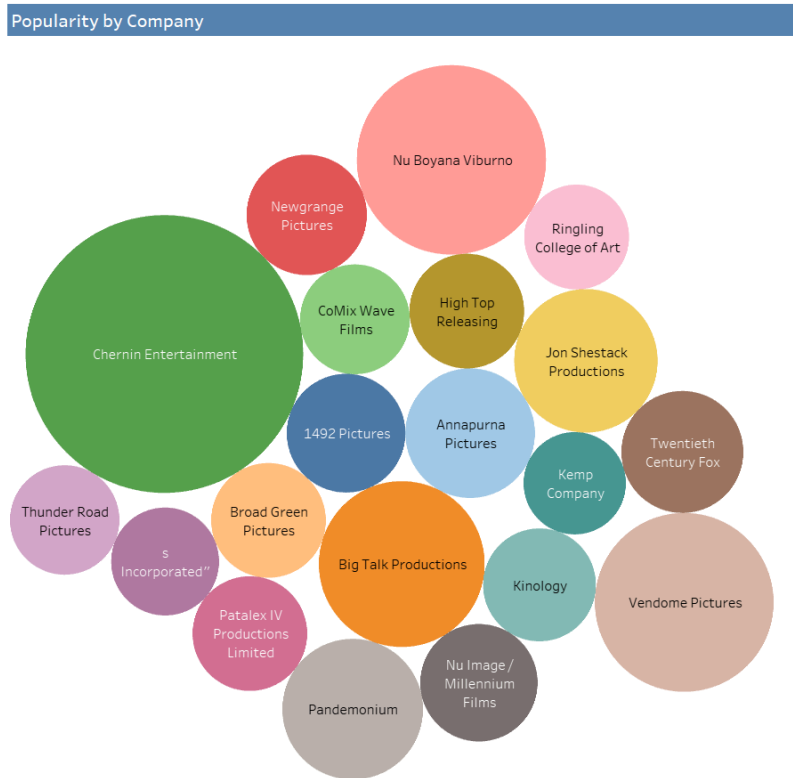| Number of Movies by Region | |
| --- | --- |
| Region | |
| United States of America | 18,429 |
| United Kingdom | 3,072 |
| France | 2,711 |
| Japan | 1,507 |
| Canada | 1,497 |
| Italy | 1,472 |
| Germany | 1,418 |
| Russia | 800 |
| India | 780 |
| Spain | 601 |
| Australia | 506 |
| Hong Kong | 468 |
| South Korea | 457 |
| Sweden | 398 |
| Finland | 324 |
| China | 300 |
| Belgium | 299 |
| Denmark | 297 |
| Brazil | 262 |
| Poland | 245 |

### 2.1.5 Popularity by Genre

The bar chart illustrates the popularity of different movie genres. Drama and comedy are the most popular genres, reflecting their broad appeal and frequent production in the film industry. Action follows closely, driven by its high entertainment value and ability to attract a diverse

audience. Genres such as horror, adventure, and thriller show moderate popularity, appealing to more niche but dedicated audience groups. On the other hand, genres like music, western, and mystery are the least popular, likely due to their narrower audience reach or lower production volume. Overall, the chart highlights the dominance of versatile genres while showcasing the varied preferences of movie audiences.

**Popularity by Genres**

| Genre | Popularity |
|---|---|
| Drama | (bar to ~32K) |
| Comedy | (bar to ~25K) |
| Action | (bar to ~20K) |
| Horror | (bar to ~9K) |
| Adventure | (bar to ~8K) |
| Thriller | (bar to ~6K) |
| Crime | (bar to ~6K) |
| Animation | (bar to ~4K) |
| Fantasy | (bar to ~3K) |
| Documentary | (bar to ~3K) |
| Romance | (bar to ~2.5K) |
| Science Fiction | (bar to ~2.5K) |
| Family | (bar to ~2K) |
| Mystery | (bar to ~1.5K) |
| Western | (bar to ~0.7K) |
| Music | (bar to ~0.7K) |
| War | (bar to ~0.7K) |
| TV Movie | (bar to ~0.7K) |
| History | (bar to ~0.5K) |
| Foreign | (minimal) |

### 2.1.6 Number and Popularity of Movies Popularity by Company

The visualization illustrates the popularity of various film production companies through circles of differing sizes, where larger circles indicate higher popularity. Cherin Entertainment emerges as the most prominent company, significantly overshadowing its competitors and suggesting a strong influence in the industry. Newgrange Pictures, CoMix Wave Films, and 1492 Pictures are notable mid-tier players, demonstrating solid reputations without reaching Cherin's level of prominence. The diversity of companies, ranging from well-known entities like Twentieth Century Fox to smaller firms such as Pandemonium, highlights the variety within the film sector.

**Popularity by Company**



The visualization depicts the number of movies produced by various film companies, represented through circles of varying sizes. Each circle's size corresponds to the total number of films released by the company, providing a visual representation of their output. Companies such as Warner Bros. and Walt Disney Pictures dominate in terms of film production volume, suggesting a robust capacity for creating content. However, the average popularity of the films produced by these companies may vary significantly. For instance, while Warner Bros. has a large output, the average popularity of its films might be influenced by factors such as genre diversity and marketing strategies. Conversely, companies like Universal Pictures and Twentieth Century Fox, with fewer films, may have a higher average popularity per movie, indicating a focus on quality or blockbuster hits.

This analysis highlights the relationship between production volume and average film popularity, suggesting that both quantity and quality play crucial roles in a company's overall success in the film industry.

## 2.2 Trends of the Movie Industry

### 2.2.1 Introduction

The movie industry has experienced profound shifts in its financial landscape, characterized by changes in funding, production costs, and box office revenues. This section delves into the evolving dynamics of financial investment and returns within the industry. By analyzing the relationship between investment and output, we aim to shed light on the flow of capital and return trends. This exploration will provide a deeper understanding of how financial strategies impact the production and success of films, offering valuable insights into the monetary mechanisms that drive the industry forward.

### 2.2.2 Movie Industry Development

The chart illustrates the trends in budget and revenue in the movie industry from 1925 to 2020. In the early years, both budget and revenue remained low and stable, reflecting modest investment and returns in the mid-20th century. During the 1980s and 1990s, there was a gradual increase, indicating expanded production scales and market growth. The 2000s and 2010s saw significant spikes, highlighting the impact of blockbuster films and technological advancements, leading to substantial financial investments and higher box office returns. However, from the 2010s to 2020, there were fluctuations despite the high levels of budget and revenue, possibly due to changing audience preferences and the rise of digital streaming platforms. Overall, the chart demonstrates a strong correlation between increasing budgets and rising revenues, underscoring the growing financial investments and returns in the movie industry over time.

### 2.2.3 Budget and Revenue

The treemap illustrates the distribution of movie budgets and revenues across different years. Each rectangle represents a year, with the size of the rectangle indicating the level of budget and revenue. Darker shades suggest higher values, indicating years with significant financial investment and revenue generation. The larger and darker blocks, particularly in the 2000s and 2010s, reflect a period of increased spending and higher box office returns, while lighter and smaller blocks in earlier years depict lower financial activity. This visualization effectively highlights the growth in financial scale within the movie industry over time.

The designed interactive feature for this treemap allows users to click on different years, which then links to three additional detailed charts. These charts provide an in-depth look at the geographic distribution of movie budgets and revenues for the selected year. This functionality enables a deeper exploration of how financial resources were allocated and revenue was generated across different regions, offering valuable insights into the global dynamics of the film industry for that specific year.

### 2.2.4 Company and Revenue

The pie chart displays the revenue distribution among different movie production companies, with each segment representing a company's share of total income. Different colors distinguish the contributions of each company, allowing for a clear comparison of their market impact. Major players like Walt Disney Pictures, Universal Pictures, and others are prominent, indicating their significant roles in revenue generation within the industry. This visualization effectively highlights the competitive landscape and dominance of certain studios in the movie market.

### 2.2.5 Budget vs. Revenue Scatter Plot

The scatter plot illustrates the relationship between movie budgets and box office revenues. Each point represents a film, with the x-axis showing the budget and the y-axis indicating the revenue. The plot reveals a general trend where higher budgets can lead to higher revenues, but there is considerable variability. Some films achieve substantial revenue with moderate budgets, while others with large budgets do not perform as well. This indicates that while budget is a factor in box office success, it's not the sole determinant.

### 2.2.6 Geographical Distribution of Budget and Revenue

The two maps depict the distribution of movie budgets and box office revenues across various countries and regions. Darker shades on the maps indicate higher values. The United States stands out with the deepest colors, reflecting its significant contribution to both budgets and revenues in the film industry. Other regions show varying levels of financial activity, highlighting the global nature of movie production and revenue generation. These visualizations effectively demonstrate the geographical disparities in the film industry's financial landscape.

### 2.2.7 Highlights

- **Trends in the movie industry**: The line graph shows that the movie industry has experienced significant growth since the 1970s, especially peaking in the early 21st century.
- **Budget vs. Revenue**: The scatterplot shows the relationship between movie budgets and box office revenues, with most movies concentrating in the lower range of budgets and revenues, but with a few high-budget, high-revenue movies.
- **Company Revenue Distribution**: The pie chart shows the distribution of revenues of major movie production companies, with companies such as Walt Disney Pictures occupying a larger share of the total.

- **Global Market Distribution**: The geographic distribution analysis shows the dominance of the United States, while also identifying the rise of other high-income markets, such as China and other emerging markets.

This part is a continuation of the overview section. A combination of various charts (line graphs, tree charts, pie charts, scatter plots, and world maps) is used to comprehensively explore the changes in the development of the movie industry and to deeply analyze the dynamic relationship between capital investment and box office output. This helps viewers better understand the economic trends of the movie industry. Users can further deepen their understanding of the financial flows of the movie industry by selecting different years through the interactive function to view the data performance of a specific year.

## 2.3 Factors Affecting Movie Box Office

### 2.3.1 Introduction
Understanding the factors that influence movie box office performance is crucial for stakeholders in the film industry. This section explores key elements such as marketing strategies, star power, genre preferences, release timing, and critical reviews. By analyzing these factors, we aim to provide insights into how they contribute to a film's financial success, helping filmmakers and producers make informed decisions to maximize box office returns.

### 2.3.2 Gross of Movies Published Each Year
The bar chart illustrates the total gross of films released each year from 1925 to 2020. It shows a significant upward trend, particularly from the late 1970s onwards, with noticeable peaks in the 2000s and 2010s. This increase reflects the growing scale and financial impact of the film industry over time. The darker shades in recent years indicate higher gross values, highlighting the era of blockbuster films and expanded global distribution. This visualization effectively captures the industry's growth in revenue generation across decades.

### 2.3.3 Gross with Different Genre
The bubble chart displays the gross revenue of films across different genres. Each bubble's size represents the total earnings for that genre. Action films dominate with the largest bubble, followed by Comedy and Adventure, indicating their strong box office performance. Drama, while smaller, also contributes significantly to revenue. Other genres like Crime, Horror, and Animation show varied but notable earnings. This visualization highlights the popularity and financial success of different film genres in the industry.