

Comparing Critics to the General Public: Who is a better predictor of a movie's success?

Jacob Inwald and Ollie Jones

1 Overview

This report analyzes the relationship between different aspects of a movie and its gross income. These aspects include director list, actor list, runtime and so on. To investigate these relationships, we looked at a comprehensive dataset compiled by the movie ranking website IMDB.com. This dataset includes information about a movie like the director, the actors, the genre.

We found ..

2 Introduction

Context and motivation The global movie industry is immense, with a total of 629,807 titles listed on IMDb as of December 2022 [1], some dating as far back as 1894 [2]. With such a vast selection, it is impossible for all movies to become successful blockbusters. This raises the question: what factors contribute to a movie's success? There is an abundance of data available regarding movies, their cast, crew, and even their audiences. This data can be used to gain insight into the production, reception, and overall success of films. This study aims to explore how different metrics can predict the success of a movie before it gets released to the general public. Success in this context is defined as the amount of revenue generated by the movie.

When it comes to movie ratings, there are often two different scores available: one from the general public and the other from professional movie critics. At times, scores may be nearly identical. For instance, "The Godfather" [3] received nearly identical scores from both groups. On the other hand, "The Last Jedi" [4] saw a significant difference between the two. The disparity between ratings for certain movies can make it challenging to determine whether they were successful.

Previous work Recent research has explored the various factors that contribute to the success or failure of films. Several studies and papers have been conducted in this field, examining a range of elements. Previous research has investigated the potential of utilizing IMDb data to predict the success of films [5]. Other studies have sought to identify the factors that contribute to the making of a blockbuster movie [6]. Moreover, further studies have attempted to develop their own mathematical models to predict the success of upcoming movies [7].

Objectives This study seeks to identify the most influential factors in determining a movie's success. With the multitude of factors that could potentially affect a movie's performance at the box office, it is not feasible to investigate them all. Therefore, we will focus our research on the director, actors, genre and runtime of a movie, and compare these with its viewer rating, critic score and box office revenue. We will be exploring the following question in depth: who is more accurate in predicting a movie's success, critics or the general public? Additionally, we will investigate the influence of a movie's director and actors on its success, as well as the genres that are most likely to produce successful films.

3 Data

Data provenance We utilized two movie datasets: an IMDb dataset[8] and TMD[9] (The Movie Dataset). The IMDb dataset contains data about movies made from 2006-2016, while the TMD dataset contains data about movies released on or before July 2017. Both datasets were obtained from kaggle.com and are shared under the CC0 1.0 Universal Public Domain Dedication. TMD is merged data from TMDb (The Movie DataBase) and grouplens.org, a movie ranking site. However, we only used the part that was obtained from TMDb. It is worth noting that the provenance of the IMDb dataset has been subject to controversy as it was scraped from IMDb.com. However, for this investigation, we used only a sample of this dataset which is publicly available on kaggle.com.

Data description The IMDb movie dataset contains 12 columns consisting of string or floating point values. One unique column is the Genre, which lists the arbitrary genres assigned to each movie in the dataset such as Action, Adventure, Sci-Fi, Mystery, Horror, Thriller, Animation, Comedy, Family, Fantasy, Drama, Music, Biography, Romance, History, Crime, Western, War, Musical and Sport. The column summary is shown in Table 1.

Column Name	Description	Data Type
Rank	The rank the movie has in the IMDb database	Integer
Title	The name of the movie	String
Genre	The genres that apply to the movie, there can be anywhere from 1-3 genres. A genre can be any from: Action, Adventure, Sci-Fi, Thriller, Animation, Comedy, Family, Fantasy, Drama, Music, Romance, History, Crime, Western, War, Musical, Sport, Horror, Mystery, Biography.	Genre
Description	The description of the movie	String
Director	The person who directed the movie	String
Actors	The lead roles in the movie	String
Year	The year the movie was released	Integer
Runtime (Minutes)	The runtime in minutes of the movie	Integer
Rating	The mean rating of the movie, taken from IMDb.com	Float
Votes	The amount of users that voted on a movie to give it that rating	Integer
Revenue (Millions)	The gross income the movie made at the US box office	Float
Metascore	The rating of movie, determined using aggregated weighted Critics scores	Integer

Table 1: The different columns in the IMDb-Movie-Data.csv file

We only worked with the crew table in TMD, which is a small part of the whole database. This table has three columns - cast, crew, and id. However, the cast and crew columns are not composed of discrete datapoints, instead being json files representing the entire cast or crew list. As such, when working with it, we had to extract the data using string parsing methods. The column summary is shown in Table 2.

Column Name	Description	Data Type
cast	The cast list of the movie, including all actors who appeared in it.	.json file
crew	The entire crew list of the movie, including all the people who made it.	.json file
id	The movie id - used in the larger dataset to connect tables together	Integer

Table 2: The different columns in the credits.csv file

Data processing To assist in parsing the TMD dataset, we used a Python script to count the amount of movies each individual director or actor has helped make. The resulting datasets were saved as CSV files named actor_counts.csv and director_counts.csv respectively. This newly acquired data was then merged with the original IMDb dataset, which involved dropping the Description column while replacing the Director and Actor columns with Director Exp. and Mean Lead Roles Exp. This resulted in a merged data set with structure shown in Table 3.

Column Name	Description	Data Type
Rank	See table 1	Integer
Title	See table 1	String
Genre	See table 1	Genre
Director Exp.	The number of movies that the director of the movie has made	Float
Mean Lead Roles Exp.	The mean number of movies that the lead actors have been in	Float
Year	See table 1	Integer
Runtime (Minutes)	See table 1	Integer
Rating	See table 1	Float
Votes	See table 1	Integer
Revenue (Millions)	See table 1	Float
Metascore	See table 1	Integer

Table 3: The different columns in the merged data set

To check this data was properly normalised we made a histogram plot of all the numeric variables, shown in Figure 1. As expected, a few variables did not appear to be normally distributed, namely: Revenue (Millions), Votes, Runtime (Minutes), Director Exp., and Rating.

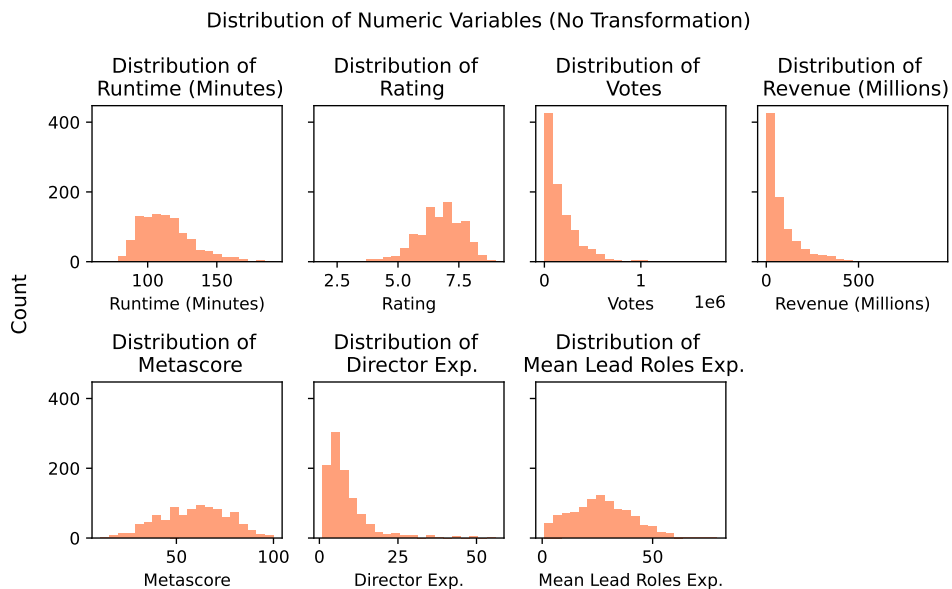


Figure 1: The distributions of the numeric variables in the merged dataset

As shown in the plot, Revenue (Millions) and Votes are severely right skewed, implying an exponential distribution; Runtime (Minutes) and Director Exp. appear to be less severely right-skewed, implying a lognormal distribution; Rating seems to be left-skewed. The transformations for these variables that gave the best approximations to a normal distribution were:

- Revenue (Millions) : Cube Root transform
- Director Exp. : Log transform
- Votes : Cube Root transform
- Runtime (Minutes) : Log transform
- Rating : Square transform

Figure 2 shows the transformed and normalised numeric variables. The label has the p-values from testing whether the transformed distribution is normal, using the Kolmogorov-Smirnov test[10] for goodness of fit. One interesting note is that although the Director Exp. column fails the Kolmogorov-Smirnov test, there is clearly missing data; around 3 columns are missing from the histogram shown. With this in mind, and as the histogram does follow the normal distribution curve, there is sufficient evidence to assume $\ln(\text{Director Exp})$ has a normal distribution.

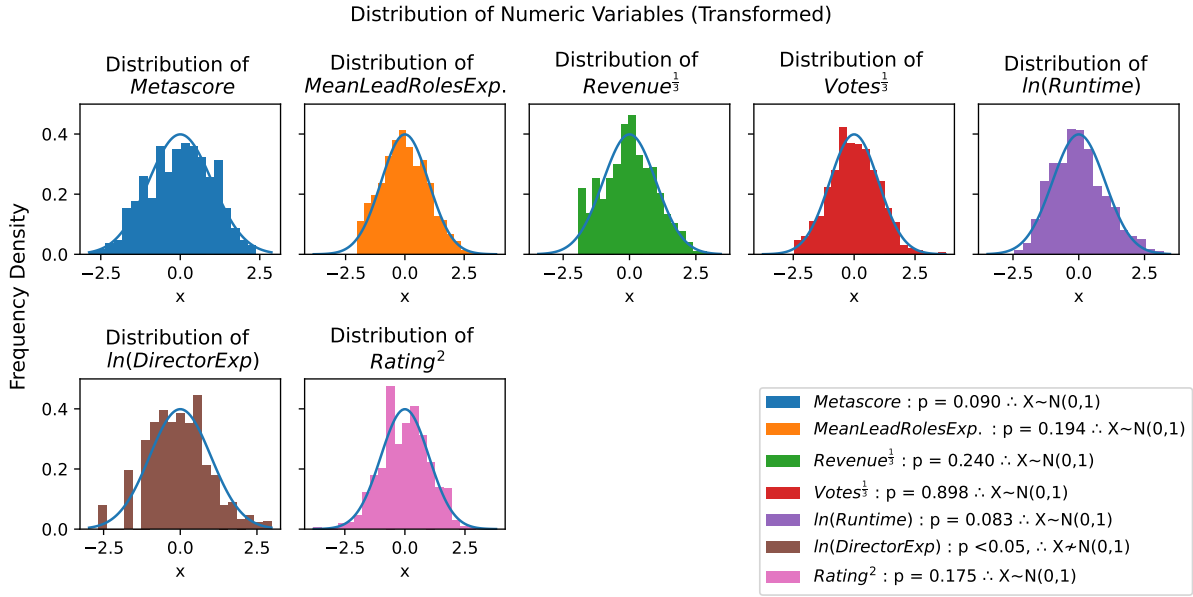


Figure 2: The distributions of the standardised and normalised numeric variables in the merged dataset, the legend has the p-values from testing with $H_0 : X \not\sim N(0,1)$. Also shown on the plots is the Gaussian Distribution with the columns mean and standard deviation

Figure 2 shows the transformed and normalised numeric variables. The legend contains the p-values from testing the data with $H_0 : X \not\sim N(0,1)$, using the Kolmogorov-Smirnov normality test[10] for goodness of fit. It's worth noting that even though the Director Exp. column failed the Kolmogorov-Smirnov test, there is missing data present in the histogram shown. Given that the histogram follows a normal distribution curve, there is convincing evidence to assume $\ln(\text{Director Exp.})$ has a normal distribution. Overall, these transformations helped to normalize the dataset and enhance its suitability for further analysis.

4 Exploration and Analysis

A data science analysis of the paper, including:

- Visualisations (for example Figure ??) and tables (for example Table ??). Please make sure that all figures and tables are referred to in the text, as demonstrated in this bullet point.
- Interpretation of the results
- Description of how you have applied one or more of the statistical and ML methods learned in the FDS to the data

- Interpretation of the findings

You can use equations like this:

$$\bar{x} = \sum_{i=1}^n x_i \quad (1)$$

or maths inline: $E = mc^2$. However, you do not need to reexplain techniques that you have learned in the course – assume the reader understands linear regression, logistic regression K-nearest neighbours etc. Remember to explain any symbols use, e.g. “ n is the number of data points and x_i is the value of the i th data point.”.

5 Discussion and Conclusions

Summary of findings

Evaluation of own work: strengths and limitations

Comparison with any other related work E.g. “Anscombe has also demonstrated that many patterns of data can have the same correlation coefficient”.

Wikipedia can also be cited but it is better if you find the original reference it for a particular claim in the list of references on the Wikipedia page, read it, and cite it.

The golden rule is always to cite information that has come from other sources, to avoid plagiarism.

Improvements and extensions

References

- [1] IMDb. *IMDb Statistics*. 2022. URL: <https://www.imdb.com/pressroom/stats/>.
- [2] Alexander Black. *Miss Jerry*. 1894. URL: https://www.imdb.com/title/tt0000009/?ref_=adv_li_tt.
- [3] Francis Ford Coppola. *The Godfather*. 1972. URL: https://www.rottentomatoes.com/m/the_godfather.
- [4] Rian Johnson. *Star Wars: The Last Jedi*. 2017. URL: https://www.rottentomatoes.com/m/star_wars_the_last_jedi.
- [5] Rijul Dhir and Anand Raj. “Movie Success Prediction using Machine Learning Algorithms and their Comparison”. In: *International Conference on Secure Cyber Computing and Communication (ICSCCC)* 1 (2018). URL: <https://ieeexplore.ieee.org/abstract/document/8703320>.
- [6] Martin C. Snell Alan Collins Chris Hand. “What makes a blockbuster? Economic analysis of film success in the United Kingdom”. In: *Managerial and Decision Economics* 23 (2002). URL: <https://doi.org/10.1002/mde.1069>.
- [7] Amr Yousef Javaria Ahmad Prakash Duraisamy and Bill Buckles. “Movie success prediction using data mining”. In: *International Conference on Computing, Communication and Networking Technologies (ICCCNT)* 8 (2017). URL: <https://ieeexplore.ieee.org/abstract/document/8204173>.
- [8] Ivan Gonzalez. *1000 IMDB movies (2006-2016)*. Scraped from <https://IMDB.com>. 2023. URL: <https://www.kaggle.com/datasets/gan2gan/1000-imdb-movies-20062016>.
- [9] Rounak Banik. *The Movies Dataset*. 2018. URL: <https://www.kaggle.com/datasets/rounakbanik/the-movies-dataset?resource=download>.

- [10] Laha Chakravarti and Roy. *Handbook of Methods of Applied Statistics*. Vol. 1. John Wiley and Sons, 1967, pp. 392–394.