

# Predicting Movie Success through IMDb Ratings and Box Office Earnings

A Data Mining Approach

Nathan Harris

College of Engineering and Applied Science  
University of Colorado Boulder  
Boulder CO, USA  
nharris@colorado.edu

Mayumi Shimobe

College of Engineering and Applied Science  
University of Colorado Boulder  
Boulder CO, USA  
mayumi.shimobe@colorado.edu

## 1. Problem Statement and Motivation

Our project aims to analyze trends and patterns in the film industry by leveraging comprehensive IMDb datasets to predict movie success—focusing on IMDb ratings as a proxy for success and incorporating box office earnings where possible. By applying data mining techniques to publicly available datasets, we hope to uncover hidden relationships among movie attributes such as runtime, release year, genre, and cast, and ultimately develop predictive models that can aid decision-makers in the industry.

## 2. Motivation and Literature Survey

Recommendation engines appear to be the most popular use for Movies and TV data sets. Eighty percent of Netflix streaming time originates from a user clicking on a recommended title [3]. These recommendation engines typically work as content-based, collaborative recommendation, or a hybrid approach [5, 8]. The content-based approach focuses on similarities between the attributes of a movie or television show whereas collaborative approach focuses on the similarities between the behavior of the user and a class of users [7]. IMDb utilizes its own recommendation engine which can be observed when scrolling to the bottom of a given title under "More like this" which leads us to believe they are focused on a content-based approach [4]. Behavior can be analyzed as explicit or implicit, explicit meaning when a user gives a positive rating to a title while implicit is when they binge watch it in one weekend [3, 9]. As far as predicting the success of film at the box office or with average reviews, it seems that there

is no consensus on the indicators for predicting a film's success. Some indicators have been combinations of characters, plot complexity, star power, budget, or "buzz," social media chatter on a particular film [1, 2]. While other indicators are external from a movie's data set, such as search engine results, social media content, and ratings on expert websites [10]. Evaluating previous movie titles is often used to support these more external data indicators.

For our project we will be focused on a content-based approach and utilize ratings as a measurement of success. Utilizing a collaborative approach would involve obtaining demographic user data in addition to implicit/explicit behavior. While we found one data set that includes this, it is far below the threshold for minimum data objects.

## 3. Proposed Work

### 3.1 Data Collection and Integration

We will use datasets from The Internet Movie Database, IMDb, and Box Office Earnings Data. Details are depicted in Section 4.

### 3.2 Data Preprocessing and Cleaning

In this section, we describe the steps we will take to transform the raw IMDb and box office data into a clean, integrated dataset suitable for predictive modeling.

Data loading: We will use MySQL's bulk-loading method (`LOAD DATA LOCAL INFILE`) to import the IMDb datasets (e.g., `title.basics.tsv` and `title.ratings.tsv`) and the box office data from Kaggle.

Any issues related to file access will be resolved by verifying working directory settings and ensuring the ``local_infile`` parameter is enabled on both the client and server sides.

Data cleaning: We will handle missing values by replacing IMDb's special marker (`\N`) with SQL NULLs, removing duplicate records, and standardizing data types, ensuring that numeric fields such as `runtime`, `startYear`, `averageRating`, and `numVotes` are stored as integers or decimals. Outlier detection will be performed using statistical methods introduced in the class to identify extreme values in IMDb ratings and box office earnings. Where necessary, outliers will be either transformed or excluded from the training set: removing outliers where the rating is only provided by a small number of reviews, i.e. two people gave 5 stars versus two thousand people who gave 4.5 stars. Socioeconomic factor like covid affecting movie theater attendance around 2020 will also be considered for outliers.

Data Integration and Transformation: The core IMDb tables—`title_basics` and `title_ratings`, along with any related attributes from chosen datasets—will be joined on the primary key (`tconst`). We will then integrate box office data from Kaggle to enrich our analysis. Categorical attributes such as genres will be preprocessed by extracting the primary genre and converting multi-valued strings into dummy variables, thereby facilitating their use in regression models. This step will also include a preliminary analysis of multicollinearity to guide subsequent feature selection.

### 3.3 Predictive Modeling

With the preprocessed dataset, we will pursue a dual modeling strategy to both predict continuous outcomes (IMDb ratings and box office earnings) and classify movies as successful or not. We will consider two possible modeling approaches.

Logistic Regression: To classify movies as "successful" (e.g., an IMDb rating above a threshold, such as 7.0 or top 20% of the total vote), we will build a logistic regression model. Evaluation metrics may include accuracy and confusion matrices. This classification approach will complement the

continuous prediction by providing actionable insights into categorical success.

Additional Machine Learning Algorithms: Depending on preliminary results, we may also experiment with non-linear models such as Random Forests or Decision Tree [11] to capture complex interactions among variables. These methods will help us validate whether advanced algorithms offer improvements over traditional regression techniques.

Our evaluation will employ k-fold cross-validation to ensure robust performance estimates [11]. We will generate scatter plots (for regression) and residual and diagnostic plots to visualize and interpret model performance. Finally, model outcomes will be compared against external benchmarks (industry reports and expert analyses) to validate our findings described in section 5 Evaluation Methods.

### 3.4 Distinction from Previous Work

Various research has been done with the content-based approach in the past to predict success in the movie industry. Nithin V.R. et al [12], for instance, explores logistic regression, SVM regression, and linear regression on IMDb and Rotten Tomatoes to predict both gross revenue and IMDb ratings. They discussed handling duplicates and missing data, yet outliers were not handling. In this study we intend to work on outliers to see the prediction is improved. We chose IMDb and Box Office Earnings to have the data availability, completeness, and consistency among data. In addition, we would make comparison with results from Kristianto et al of Random Forest and XGBoost [13] to see further investigation of machine learning is worthwhile.

## 4. Data Set Details

Data Source and URL:

IMDb Datasets: Publicly available at <https://datasets.imdbws.com/>

Box Office Data: Obtained from <https://www.kaggle.com/datasets/harios/box-office-data-1984-to-2024-from-boxofficemojo>.

Description and Dimensions:

From IMDb Datasets: two data files, `title.basics` and `title.ratings`, are considered. Preliminary cleaning process of other files listed on Part 1, such as

title.principles, implies too many outliers and empty datasets to have fair number of datasets for predictive analysis.

*title.basics.tsv.gz*: Approximately 25.3 million records with 9 attributes:

- tconst: Nominal (string) – Unique identifier.
- titleType: Nominal (string).
- primaryTitle / originalTitle: Nominal (string).
- isAdult: Nominal (binary integer).
- startYear / endYear: Interval (integer).
- runtimeMinutes: Ratio (integer).
- genres: Nominal (string, multi-valued).

*title.ratings.tsv.gz*: Over 15 millions records with 3 attributes:

- tconst: Nominal (string). Ordinal (integer)
- averageRating: Ratio (decimal).
- numVotes: Ratio (integer).

From Box Office Earnings Dataset:

Year: Date (Date)

Title: Nominal (string)

Gross: Ratio (Integer)

## 5. Evaluation Methods

Once study of the predictive models is conducted, the result would be evaluated by its accuracy and confusion matrices. Any visual plots (predicted vs. actual values) would be provided for better analysis. In addition, the results will be compared with literature. Nithin VR et al found accuracy of roughly 51% for linear regression, 42% for logistic regression, and 39% for SVM [12]. We will investigate the effect of outlier handling as well as necessity for the machine learning algorithms by comparing with accuracy of ~84% reported by Kristianto et al using Random Forest [13].

## 6. Tools

Our project will employ the following tools:

Programming Languages: Python (with Pandas, NumPy, scikit-learn, and StatsModels for analysis and modeling).

Database: MySQL on AWS RDS for data storage, integration, and querying.

Visualization: Matplotlib and Tableau for generating both static and interactive visualizations.

Version Control: GitHub for collaborative coding and version management.

## 7. Milestones

We have set up AWS RDS instance, created database/tables, loaded IMDb and Box Office datasets and conducted initial data cleaning and preprocessing for some of datasets.

By Mar 7, 2025: Perform exploratory data analysis using SQL and Python and begin feature engineering, including encoding categorical variables and assessing multicollinearity.

By March 14, 2025: Develop and test predictive models (linear regression, logistic regression, and possibly additional machine learning algorithms such as Random Forests). Evaluate model performance and refine feature selection based on analysis and correlation matrices.

By April 11, 2025: Finalize visualizations and interpretability analyses. Compare model results with external benchmarks and expert analyses.

By April 18, 2025: Prepare the final report and project presentation, incorporating feedback and ensuring that all aspects of the evaluation framework are clearly documented. We have set our weekly Agile meeting every Thursday to meet deadlines for each project parts and final submission.

## ACKNOWLEDGMENTS

We greatly appreciate feedback and comments from professor and fellow students in Spring 2025 semester of CSPB4502 Data Mining class in University of Colorado Boulder.

## REFERENCES

- [1] Coming to a Screen Near You: How Data Science Is Revolutionizing the Film Industry. Retrieved February 1<sup>st</sup> 2025 from: <https://datacolumn.iaa.ncsu.edu/blog/2023/01/30/coming-to-a-screen-near-you-how-data-science-is-revolutionizing-the-film-industry/>
- [2] Warner Bros Will Now Use AI to Help Decide Which Movies to Make. Retrieved February 1<sup>st</sup> 2025 from: <https://www.ign.com/articles/2020/01/08/warner-bros-will-now-use-ai-to-help-decide-which-movies-to-make>
- [3] Netflix Recommendations and Page Rank. Retrieved February 1<sup>st</sup> 2025 from: <https://blogs.cornell.edu/info2040/2020/11/08/netflix-recommendations-and-page-rank/>
- [4] Deep Dive into Netflix's Recommender System. Retrieved February 1<sup>st</sup> 2025 from: <https://towardsdatascience.com/deep-dive-into-netflixs-recommender-system-341806ac3b48>
- [5] Building a Collaborative Filtering Recommendation Engine. Retrieved February 1<sup>st</sup> 2025 from: <https://datacolumn.iaa.ncsu.edu/blog/2020/02/10/building-a-collaborative-filtering-recommendation-engine/>
- [6] MovieLens 100k Dataset. Retrieved February 1<sup>st</sup> 2025 from: <https://grouplens.org/datasets/movielens/100k/>
- [7] Netflix Research Archive. Available at: <https://research.netflix.com/archive>

- [8] Data Science in the Film Industry. Retrieved February 1<sup>st</sup> 2025 from: <https://www.kdnuggets.com/2019/07/data-science-film-industry.html>
- [9] How Data Analytics Help Movie Success Prediction. Retrieved February 1<sup>st</sup> from: <https://datavisitor.com/how-data-analytics-help-movie-success-prediction/>
- [10] Building Recommender Systems. Retrieved February 1<sup>st</sup> 2025 from: <https://www.kdnuggets.com/2019/04/building-recommender-system.html>
- [11] Jiawei Han, Micheline Kamber, Jian Pei. Data Mining: Concepts and Techniques, 3rd Edition. Morgan Kaufmann, 2011.
- [12] Nithin V R, Sarath Babu P, Pranav M, Lijjiya A, "Predicting Movie Success Based on IMDb Data," International Journal for Research in Applied Science & Engineering Technology (IJRASET), vol. 5, issue X, pp. 503–507, Oct. 2017.
- [13] Michael Kristianto, Deastri Anggie Shanovera, Janette Mirna Putri Trijono, Ican Sebastian Edbert, Derwin Suhartono, "Movie Success Prediction Using Machine Learning Models." Intemational Conference on Technology Innovation and Its Applications, 2024