

AMERICAN INTERNATIONAL UNIVERSITY- BANGLADESH



Project: IMDB 5000 Movie Dataset

Course Title: INTRODUCTION TO DATA SCIENCE.

Section: C

Date of Submission: 16-08-2025

Semester: Summer, 2024-2025

Course Teacher: DR. ASHRAF UDDIN

Group No: 9

No	Name	ID	Program
01	MD TANJIL TASHRIK ZIM	22-48021-2	BSc in CSE
02	MD. AL-IMRN SAYEM	22-48023-2	BSc in CSE
03	MD. ABRAR RAFID SHARIAR	22-48055-2	BSc in CSE
04	S.M. RASEL	22-48039-2	BSc in CSE

Faculty use only

FACULTYCOMMENTS

	Marks Obtained	
	Total Marks	

1. Introduction

This project analyzes the IMDB 5000 Movie Dataset, focusing on key attributes such as IMDb scores, budget, gross revenue, duration, and content ratings. The objective is to uncover patterns and trends that influence both a movie's financial success and critical reception. To achieve this, the dataset was cleaned and preprocessed to ensure accuracy, followed by statistical analysis, visualizations, and predictive modeling techniques. Through descriptive statistics, correlation analysis, and regression models, this study provides valuable insights into how different factors interact within the film industry. The findings are intended to support filmmakers, researchers, and industry enthusiasts in understanding the dynamics behind successful movies.

Key column:

- `imdb_score` (numeric): IMDb rating of the movie.
- `budget` (numeric): Movie production budget in USD.
- `gross` (numeric): Gross revenue in USD.
- `duration` (numeric): Movie runtime in minutes.
- `content_rating` (categorical): MPAA or equivalent rating .
- `genres` (categorical): List of genres (e.g., Drama, Comedy).
- `num_user_for_reviews`, `num_critic_for_reviews` (numeric): Engagement metrics.
- `num_voted_users` (numeric): Total number of user votes.
- `movie_facebook_likes` (numeric): Facebook popularity.
- `title_year`, `language`, `country` (categorical): Production metadata.

2. Data Preprocessing Steps

Data Inspection

1. Dataset size: **5,043 rows × 28 columns**.
2. Data types checked: Numeric (`imdb_score`, `budget`, `gross`, `duration`) and Categorical (`content_rating`, `genres`, `language`).
3. Initial review showed missing values in several important columns.

Handling Missing Values

- Checked using `colSums(is.na(movie_data))`.
- Missing values found:
 - `gross` = 884
 - `budget` = 492
 - `aspect_ratio` = 329
 - `content_rating` = 303
- Action: Removed rows with NA values using `na.omit()`.

- Justification: Ensures clean, complete dataset for reliable analysis.

Removing Duplicates

- Code removed duplicate rows using `!duplicated()`.
- Justification: Prevents repeated records from biasing results.

Exploratory Data Analysis (EDA)

Univariate Analysis (single variable):

- **IMDb Score:** Mean = 6.44, Median = 6.60, Mode = 6.7, SD = 1.13, IQR = 1.40.
- **Budget:** Strongly skewed; most budgets under \$50M.
- **Gross Revenue:** Right-skewed, dominated by few blockbusters.
- **Profit:** Wide variance; some low-budget movies achieved massive ROI.
- **Duration:** Mean \approx 107 minutes; most clustered between 90–120 min.
- Histograms plotted for IMDb score, budget, and other numeric variables.

Multivariate Analysis (multiple variables):

- Scatterplot matrix: Relationships among IMDb score, budget, gross, and duration.
- Correlation heatmap:
 - `num_voted_users` \leftrightarrow `imdb_score`
 - `num_critic_for_reviews` \leftrightarrow `imdb_score`
 - `budget` \leftrightarrow `gross` weak

Data Wrangling

- Filtered only movies with `gross > $100M` (~12%).
- Selected key attributes (`title`, `imdb_score`, `budget`, `gross`).
- Created **profit** column (`gross - budget`).

Normalization & Scaling

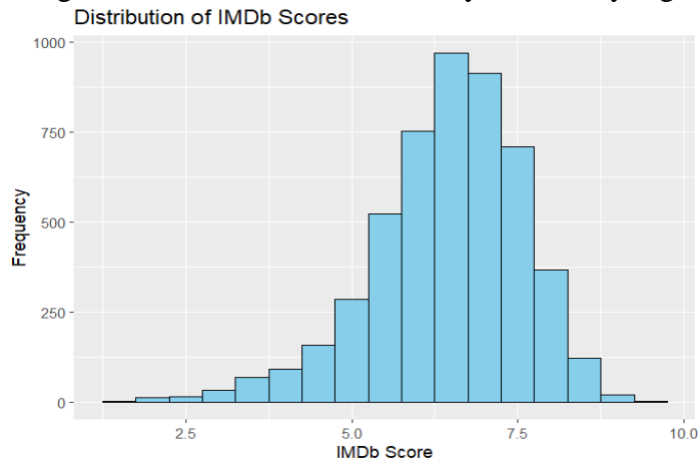
- Scaled numeric features (`imdb_score`, `budget`, `gross`) for consistency.
- Justification: Ensures fair comparison and prepares dataset for machine learning.

3. Key Findings & Visualizations

Finding 1: IMDb Score Distribution

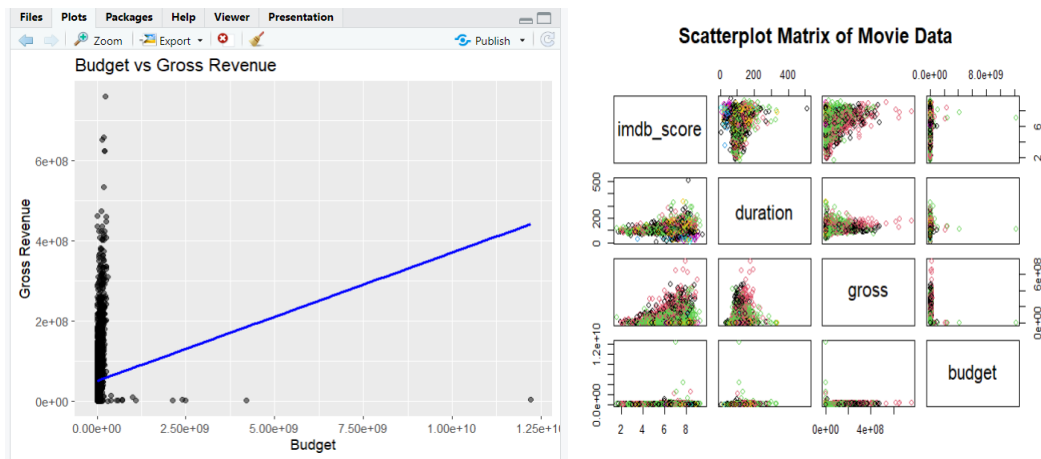
- Histogram shows majority clustered between **6–7**.

- Insight: Most movies are moderately rated; very high or very low scores are rare.



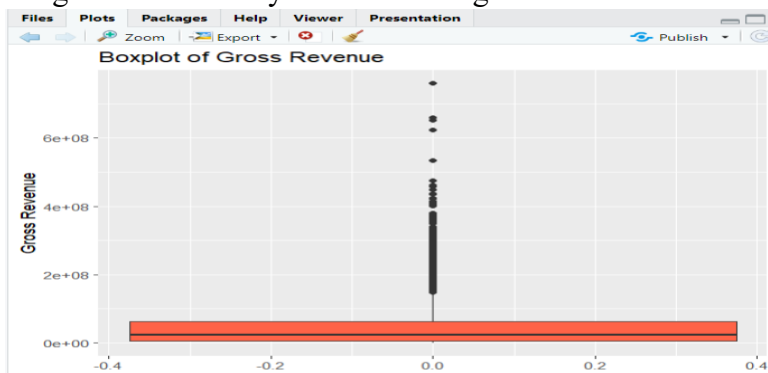
Finding 2: Budget vs Gross Revenue

- Scatterplot with regression line shows weak relationship ($R^2 \approx 0.01$).
- Insight: Large budgets do **not** guarantee commercial success.



Finding 3: Boxplot of Gross Revenue

- Shows strong right-skew; a few blockbusters dominate revenues.
- Insight: Film industry follows a “long-tail” distribution.



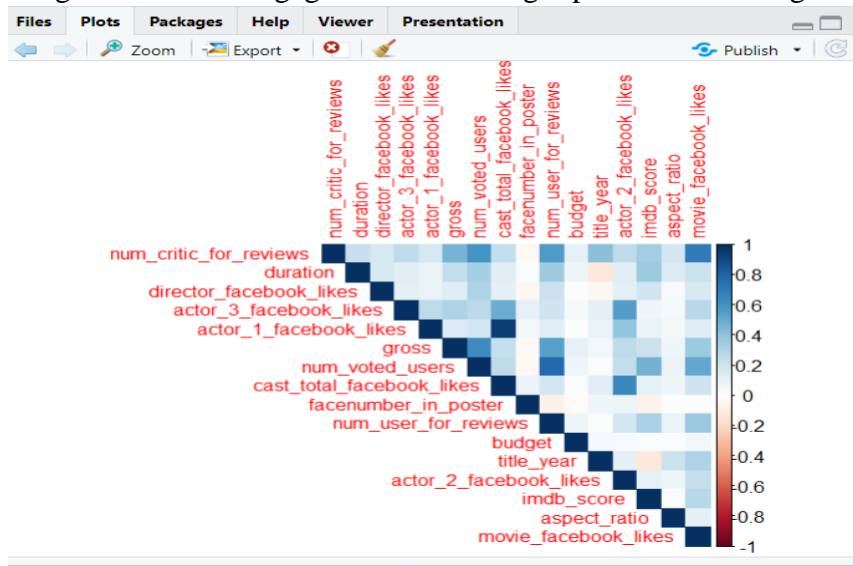
Finding 4: IMDb Score vs Profit

- Scatterplot shows a weak but noticeable trend that higher IMDb scores link to better profitability.
- Suggests critical reception may influence financial outcomes indirectly.



Finding 5: Correlation Heatmap

- Positive correlations between engagement metrics (votes, critic reviews) and IMDb scores.
- Insight: Audience engagement is a stronger predictor of ratings than budget is of revenue.



Finding 6: High-Grossing Movies (> \$100M)

- Only about **12%** of movies cross this mark.
- Derived **profit** metric highlights some low-budget films with very high ROI.

35	color	Marc Webb	493	142
1		actor_2_name	actor_1_facebook_likes	gross
2	Joel David Moore		1000	760505847
3	Orlando Bloom		40000	309404152
4	Rory Kinnear		11000	200074175
5	Christian Bale		27000	448130642
6	James Franco		24000	336530303
7	Donna Murphy		799	200807362
8	Robert Downey Jr.		26000	458991599
9	Daniel Radcliffe		25000	301596980
10	Lauren Cohan		15000	330249062
11	Marlon Brando		18000	200069408
12	Mathieu Amalric		451	168368427
13	Orlando Bloom		40000	423032628
14	Christopher Meloni		15000	291021565
15	PierFrancesco Favino		22000	141614023
16	Robert Downey Jr.		26000	623279547
17	Sam Claflin		40000	241063875
18	Michael Stuhlbarg		10000	19020854
19	Adam Brown		5000	255108370
20	Andrew Garfield		15000	262030663
21	William Hurt		891	105219735
22	Adam Brown		5000	258355354
23	Thomas Kretschmann		6000	218051260
24	Kate Winslet		29000	658672302
25	Scarlett Johansson		21000	407197282
26	Judy Greer		3000	652177721
27	Helen McCrory		883	304360277
28	James Franco		24000	373377893
29	Jon Favreau		21000	408992722
30	Alan Rickman		40000	334185206
31	Kelsey Grammer		20000	234360014
32	Tyler Labine		12000	268488329
33	Kevin Dunn		894	402076689
34	Sophia Myles		974	245428137
35	Mila Kunis		44000	234903076
36	Andrew Garfield		15000	202853933
1				genre
1		Action Adventure Fantasy Sci-Fi		



4. Conclusion

This project explored the Movie Metadata dataset through data cleaning, statistical analysis, and visualization. At the beginning, missing values and duplicate records were removed to improve data quality. New variables such as profit and return on investment were created to better understand the financial performance of movies.

Descriptive statistics (mean, median, mode, standard deviation, and interquartile range) were calculated for key variables including budget, gross revenue, profit, IMDb score, and duration. These measures provided insights into the overall trends and variability in the dataset.

Univariate and bivariate visualizations revealed important patterns. Histograms showed the distribution of IMDb scores, budgets while scatterplots highlighted the relationships between budget and gross revenue, IMDb score and profit, as well as ROI and IMDb score. Correlation analysis further demonstrated strong links between budget, revenue, and profit.

The genre analysis identified the top 15 most common genres, showing which types of movies are most frequently produced. This information, along with IMDb scores and profitability measures, helps to understand both audience preferences and market performance across different genres.

In conclusion, this project demonstrates how data science techniques can be effectively applied to real-world datasets. It highlights the importance of thorough data cleaning, descriptive analysis, and visualization in generating meaningful insights that support informed decision-making.