# Movies Analytics Report

## 1. Executive Summary

This report summarizes exploratory and descriptive analytics performed on a movies dataset. It highlights data description, cleaning steps, genre and rating distributions, director and year trends, and provides actionable insights based on the PySpark analysis.

## 2. Dataset Description

Source: movies_dataset.csv
Sample size: ~1,000,000 records
Typical columns include:
• MovieID — Unique movie identifier
• Title — Movie title
• Genre — Primary or multiple genres
• ReleaseYear, ReleaseDate — Year and date of release
• Country — Country of origin
• BudgetUSD, US_BoxOfficeUSD, Global_BoxOfficeUSD — Financial performance
• IMDbRating, RottenTomatoesScore — Rating metrics
• NumVotesIMDb, NumVotesRT — Popularity indicators
• Director, LeadActor — Creative leads

### 2.1 Data Quality Summary

• Missing values identified in numeric fields were treated as null.
• Empty strings replaced with nulls for categorical columns.
• Duplicates removed based on MovieID.
• Column names normalized to lowercase with underscores.
• Non-parsable release years were corrected using date inference.

## 3. Operations Performed

### 3.1 Data Cleaning & Preprocessing

Data cleaning involved handling nulls, type casting, and splitting multi-valued genres. Date fields were standardized to year format, and key numeric columns were verified for consistency.

### 3.2 Descriptive Analytics & Visualizations

Key analyses performed:
• Genre popularity (movie count by genre)
• IMDb rating distribution
• Top directors by number of films
• Yearly release trends
• Budget vs. global box office correlation

## 4. Key Insights

### 4.1 Genre & Ratings

• Drama, Action, and Comedy emerged as dominant genres by count.
• IMDb ratings exhibit a right-skewed distribution, with most movies clustering between 6–8.
• Rating bucket analysis shows a concentration in the 6–8 range.

### 4.2 Directors & Production Trends

• A small number of directors contribute disproportionately to the dataset's total films.
• The 1990s and 2010s mark significant peaks in global movie releases.
• Top-grossing directors show consistent high IMDb ratings.

### 4.3 Financial Metrics

• Movies with higher budgets tend to exhibit a positive, though non-linear, correlation with box office returns.
• The variance suggests the influence of non-financial factors such as genre, cast, and marketing.

### 4.4 Yearly & Geographic Trends

• Movie releases have steadily increased over the decades, reflecting industry growth.
• The USA, UK, and India lead in production volume.
• Streaming platforms are increasingly represented in post-2015 data.

# 5. Recommendations

### 5.1 Production Strategy

Focus investments on genres with strong audience appeal (Action, Drama) and proven ROI. Encourage diversity in lower-performing genres to explore niche markets.

### 5.2 Marketing & Distribution

• Schedule releases around seasonal peaks to maximize audience engagement.
• Use rating-based segmentation to optimize promotional budgets.
• Leverage data from streaming trends for direct-to-digital strategies.

### 5.3 Analytical Roadmap

• Implement predictive modeling for box office success using regression on budget, genre, and director variables.
• Develop dashboards (Power BI / Tableau) to track trends and KPIs.
• Automate monthly movie analytics pipelines for future datasets.

# 6. Appendix

This report was generated using PySpark-based analysis on the provided movies dataset. Visualizations were created with Seaborn and Matplotlib, showing trends in genres, ratings, directors, and financial performance.