

Movie Success Prediction and Sentiment Study

Introduction

Movies play a significant role in entertainment and culture, influencing audiences worldwide. Predicting their success before release is a major challenge for producers and analysts. This project aims to develop a **data-driven approach** to predict the success of movies using IMDB/Kaggle datasets. Additionally, it performs a **sentiment analysis of viewer reviews** to understand audience opinions and their correlation with movie ratings and box office outcomes.

Abstract

The Movie Success Prediction and Sentiment Study project combines machine learning and natural language processing (NLP) techniques to forecast a movie’s performance based on features such as genre, rating, and user sentiments. Using VADER Sentiment Analyzer from the NLTK library, the system evaluates the tone of user reviews and classifies them as positive, neutral, or negative. A regression model (using Scikit-learn) is then trained to predict box office success or rating scores. The analysis also includes genre-wise sentiment distribution to understand which types of movies resonate most positively with viewers. This project demonstrates how computational techniques can enhance entertainment analytics and assist stakeholders in making data-informed decisions.

Tools Used:

Tool / Library	Purpose
Python	Core programming language
NLTK (VADER)	Sentiment analysis of movie reviews
Scikit-learn (Sklearn)	Building regression models
Pandas & NumPy	Data preprocessing and analysis
Matplotlib / Seaborn	Visualization of sentiment and prediction trends
Excel	Tabular data review and result compilation

Steps Involved in Building the Project

1. Data Collection:

Imported datasets from **IMDB/Kaggle**, including movie title, genre, ratings, and box office collections. Optionally, user reviews were scraped for sentiment analysis.

2. Data Preprocessing:

Cleaned datasets by handling missing values, removing duplicates, and normalizing numeric fields like revenue and rating.

3. Sentiment Analysis using VADER:

Applied VADER SentimentIntensityAnalyzer from NLTK to extract polarity scores (positive, negative, neutral, compound) from each review.

4. Feature Engineering:

Combined sentiment scores with numerical and categorical movie attributes (genre, duration, rating, etc.) to build input features for prediction.

5. Model Development:

Implemented Linear Regression and Random Forest Regressor using Scikit-learn to predict box office success or rating score based on extracted features.

6. Visualization and Trend Analysis:

- a. Plotted genre-wise sentiment averages using Seaborn bar charts.
- b. Compared predicted vs. actual revenues.
- c. Displayed sentiment distributions using pie and scatter plots.

7. Evaluation:

Used metrics like R^2 Score, Mean Squared Error (MSE), and Sentiment Accuracy to evaluate the model's performance.

Conclusion:

The project successfully integrated sentiment analysis with predictive modeling to estimate movie success. Results showed that viewer sentiment positively correlates with movie ratings and box office revenue. Genres like *Drama* and *Adventure* exhibited the highest positive sentiment scores, indicating stronger emotional engagement from audiences. This study highlights the potential of combining NLP and machine learning for entertainment analytics, providing valuable insights for filmmakers, critics, and audiences alike.