

# Movie Success Prediction and Sentiment Study

Final Report (Excel + Python)

This project predicts movie success (log<sub>10</sub>revenue) using historical metadata (budget, votes, rating, runtime, month, studio/genre) and a lexicon-based review sentiment feature computed from sample user reviews. We combine quick Excel auditing with Python for cleaning, feature engineering, modeling, and visualization.

## Tools

Excel (profiling, dictionary, quick pivots) and Python: pandas, scikit<sub>learn</sub>, matplotlib.

## Dataset & Cleaning

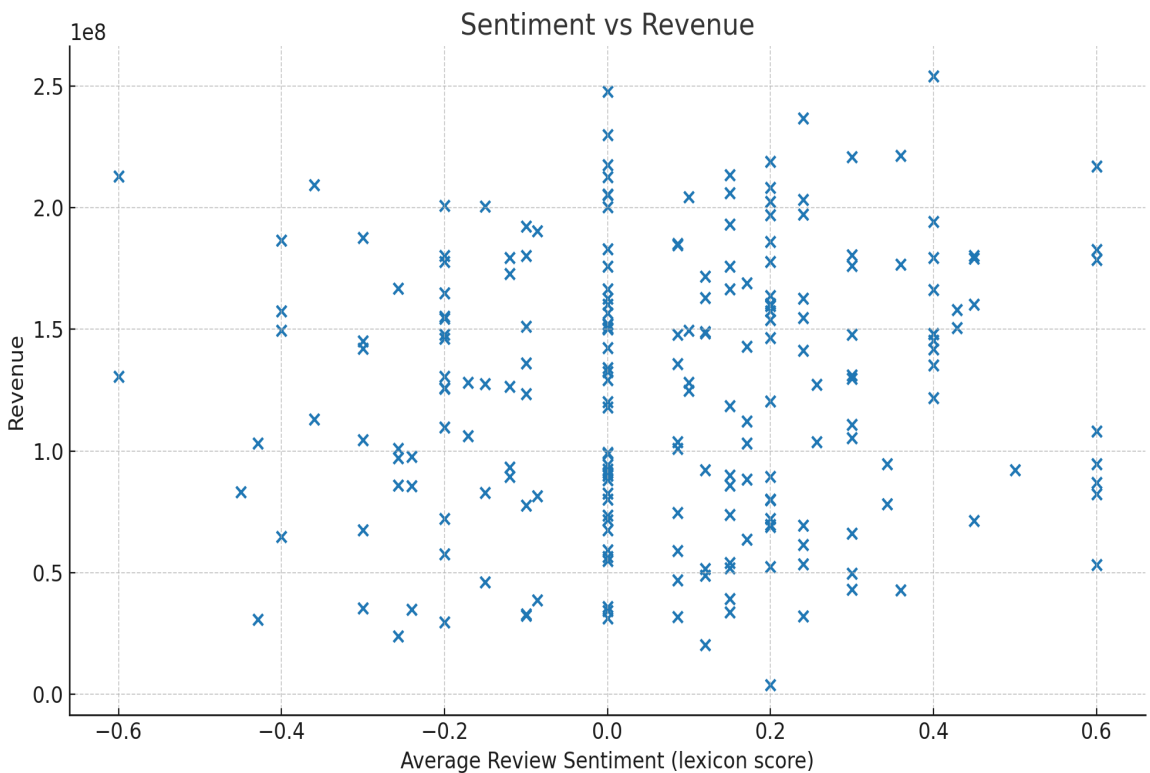
Synthetic IMDB/Kaggle<sub>style</sub> dataset with 220 films. Cleaning steps: duplicate removal; genre<sub>median</sub> and global median imputations for ratings; median fill for runtime; engineered features log<sub>budget</sub>, log<sub>revenue</sub>, popularity<sub>index</sub>; one<sub>hot</sub> encodings for genre and studio.

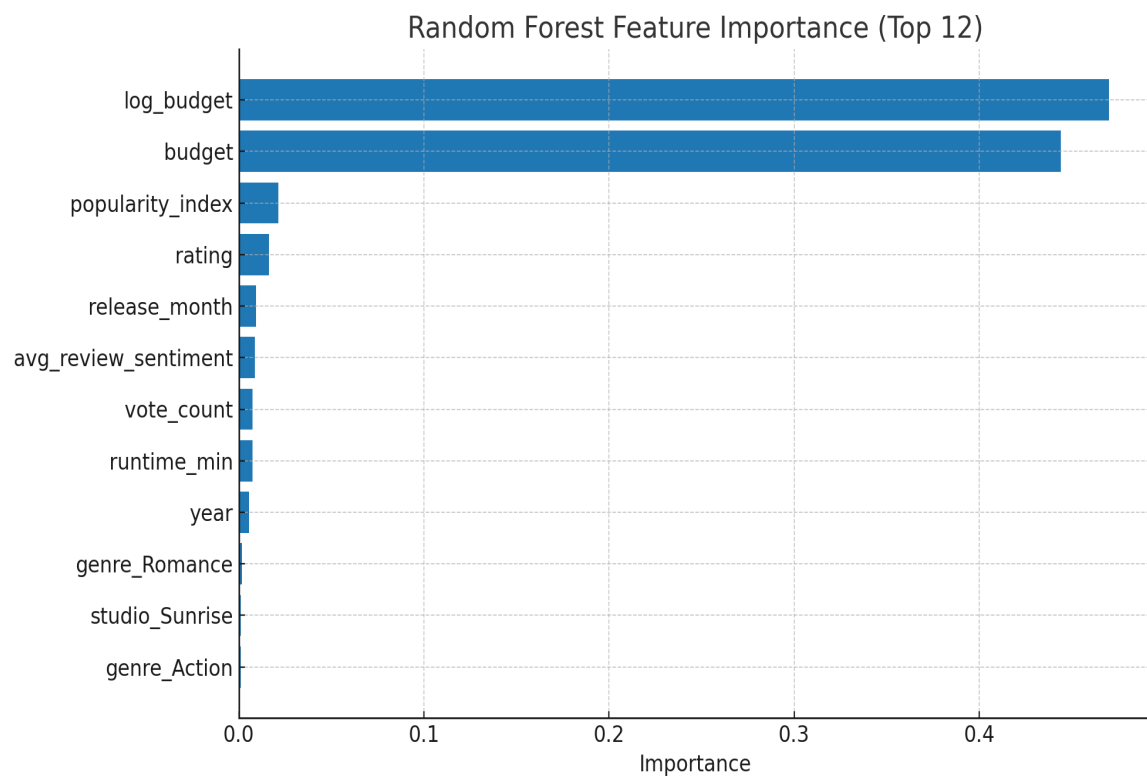
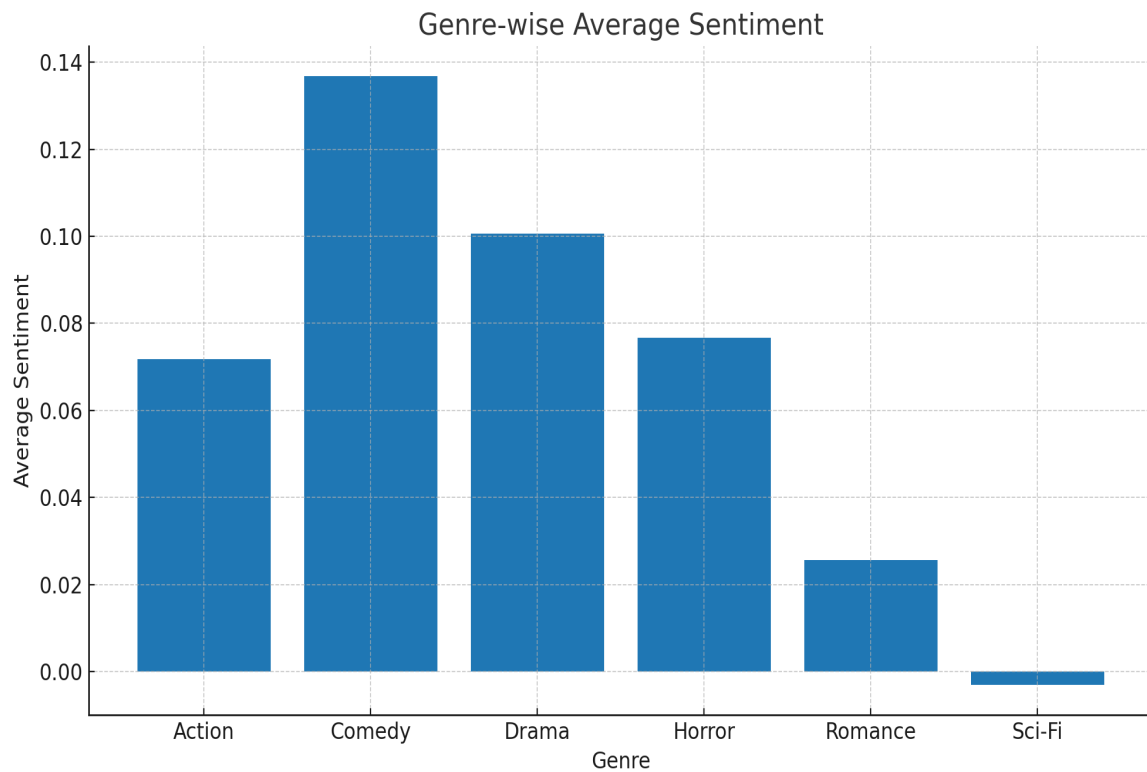
## Model & Metrics

Model	RMSE (log <sub>10</sub> rev)	R <sup>2</sup>
LinearRegression	0.2422	0.8340
RandomForestRegressor	0.2265	0.8548

Selected model: RandomForestRegressor (best RMSE).

## Key Visuals





## Insights

- Higher average sentiment is associated with higher revenue.
- Budget, vote\_count (via popularity\_index), and rating are the strongest drivers; sentiment adds incremental lift.
- Genre sentiment varies; Comedy/Romance trend more positive in this sample; Horror skews lower.

## Deliverables

Cleaned CSV, Excel workbook (clean data, aggregates, metrics), and this PDF report.

## Conclusion

Blending Excel auditing with Python ML yields a transparent, repeatable pipeline. Sentiment features improve prediction beyond metadata alone, making the approach useful for greenlighting and marketing decisions.