

National University of Singapore

School of Computing

Forecasting Box Office Performance

Predicting Opening Weekend and Revenue Trends

BT4222: Mining Web Data For Business Insights

Final Project Report

Group No. 5

Group Members

BESBES Omar

GUERMAZI Housseem

Teh Ze Shi

Elgin Song Zhanyi

September 13, 2025

Abstract

This project addresses a critical challenge in the film industry: forecasting a movie's **opening weekend gross** and the **evolution of its revenue** while it remains in theaters. Accurate predictions in these areas offer valuable insights for stakeholders making early investment, marketing, and distribution decisions.

We structured the problem in two stages. First, we used machine learning models — including **Random Forest**, **XGBoost**, and regression techniques — to predict the *first weekend gross*, which is both commercially significant and highly unpredictable. Second, we employed **ARIMA** and **LSTM** time-series models to forecast *how revenue progresses over time* during the theatrical run.

Our data sources included the *TMDB dataset*, *The Numbers dataset*, and *Wikipedia*, from which we collected key features such as *budget*, *release date*, *genre*, and *cast and crew*. Models were evaluated using **RMSE**, **MAPE**, and **Median Absolute Error**, with our best first-weekend prediction achieving a **30% median error** — a meaningful benchmark given the randomness and volatility of opening performance.

The project offers several important business insights. First, early gross estimates can **guide investment decisions**, helping producers and investors assess financial risk before a film's release. Second, revenue trend forecasts support **adaptive marketing strategies**, enabling campaigns to adjust dynamically post-launch. Lastly, understanding the likely trajectory of box office performance helps studios and theaters **optimize distribution timing**, such as deciding how long to keep a movie in theaters or when to transition it to streaming platforms.

By blending predictive modeling with business insight, our project offers a practical framework for making more informed, data-driven decisions in the high-stakes world of film distribution.

Contents

| | |
|-----------------------------------------------------------------------------|-----------|
| Abstract | 1 |
| 1 Feature Engineering | 3 |
| 1.1 Overview | 5 |
| 1.2 Target Variable Cleaning | 5 |
| 1.3 Numerical Feature Transformation | 5 |
| 1.4 Temporal Feature Engineering | 5 |
| 1.5 Categorical and Text Features | 5 |
| 1.6 Business-Oriented Feature Engineering | 6 |
| 1.7 Synopsis Text Vectorization | 6 |
| 1.8 Final Dataset | 6 |
| 2 Models and Performance | 7 |
| 2.1 Task 1 : Predicting Opening Weekend | 7 |
| 2.2 Task 2: Revenue Predictions | 10 |
| 3 Contribution and Justification | 14 |
| 3.1 Effort in Data Collection and Dataset Enrichment | 14 |
| 3.1.1 Multi-Source Integration Strategy | 14 |
| 3.1.2 Value of Each Data Source | 14 |
| 3.1.3 Example: Rich Metadata for Real-World Titles | 14 |
| 3.1.4 Depth and Originality of Final Dataset | 15 |
| 3.2 Complexity and Innovation in Revenue Trends Forecasting Model | 15 |
| 3.2.1 Innovative Components and the reasons behind our choices | 15 |
| 3.2.2 Why the Model Works Well | 16 |
| 3.2.3 Impact and Innovation | 16 |
| 3.3 Creativity in Feature Engineering | 17 |
| 3.4 Creativity and Insights in Understanding Model Output | 18 |
| 3.4.1 Why Opening Weekend Is Hard to Predict | 18 |
| 3.4.2 Evidence from Prior Research | 19 |
| 3.4.3 Why RMSE and MAPE Are Insufficient | 19 |
| 3.4.4 Why MdAPE Works Better | 19 |
| 3.4.5 Business Implications of MdAPE Accuracy | 19 |
| 3.4.6 Bias in Model Predictions | 20 |

1 Feature Engineering

Our initial data processing for the interim report provided a dataset containing approximately **50,000 gross revenue records** corresponding to **1,447 movies** released between 2020 and 2025. These records reflected the daily box office performance of each film. While this dataset was sufficient for modeling revenue evolution as a time series, it proved inadequate for accurately predicting the **opening weekend gross** based solely on available metadata and budget information.

One major challenge was the **high presence of outliers** in opening weekend gross values — driven by blockbuster releases or heavily marketed films — which introduced significant variance and instability in the predictive models. Additionally, the **limited diversity in metadata** for just 1,400 movies over a five-year window restricted the model's ability to generalize patterns effectively, especially for underrepresented genres, budgets, or cast profiles.

To address this, we extended our dataset by **scraping metadata for movies released from 2000 to 2025**, significantly increasing the volume and variety of training examples. As a result, our enriched dataset now includes metadata for approximately **4,625 movies**, offering a broader and more representative view of industry trends over 25 years. This enhancement allowed us to build a more diverse and balanced metadata feature space, improving the model's exposure to different production scales, genre patterns, and cast/crew combinations.

Table 1: Numerical Features Summary

| Feature | Min | Max | Mean | Median | Std Dev | Missing |
|------------------------|-----------|------------|------------|------------|------------|---------|
| opening_weekend | 50 | 35,711,500 | 17,997,140 | 10,207,869 | 26,838,450 | 944 |
| percent_of_total_gross | 0.0 | 100.0 | 30.10 | 32.1 | 16.29 | 944 |
| production_budget | 4,819,277 | 36,144,578 | 23,694,780 | 30,120,481 | 16,621,850 | 1435 |

Table 2: Textual Features Overview

| Feature | Description and Example | Missing Values |
|-----------------------|---------------------------------------------------------------------------------------------------------------------------|----------------|
| synopsis | A short plot summary of the movie. <i>Example:</i> It had been a year since Dr. Norman Spencer began hearing voices... | 837 |
| keywords | Thematic or genre-related keywords. <i>Example:</i> Psychological Thriller, Infidelity, Relationship | 569 |
| source | Origin of the story or script. <i>Example:</i> Original Screenplay | 97 |
| genre | General movie classification. <i>Example:</i> Thriller/Suspense | 55 |
| production_method | Film type or technique. <i>Example:</i> Live Action | 57 |
| creative_type | Narrative style or story format. <i>Example:</i> Contemporary Fiction | 88 |
| production_companies | Studios or production companies. <i>Example:</i> DreamWorks Pictures, 20th Century Fox | 1397 |
| production_countries | Country where the movie was produced. <i>Example:</i> United States | 103 |
| languages | Languages spoken in the film. <i>Example:</i> English | 269 |
| lead_ensemble_members | Names and roles of lead cast. <i>Example:</i> {'actor': 'Harrison Ford', 'role': 'Norman Spencer'} | 0 |
| production_technical | Director, screenwriter, etc. <i>Example:</i> {'Director': ['Robert Zemeckis']} | 0 |
| movie_title | Title of the movie. <i>Example:</i> What Lies Beneath | 0 |
| movie_url | External URL link. <i>Example:</i> https://www.the-numbers.com/... | 0 |

1.1 Overview

To prepare our dataset for predicting opening weekend gross, we followed a comprehensive preprocessing pipeline that included cleaning raw fields, handling missing values, engineering new business-relevant features, and transforming categorical and textual variables into model-ready formats.

1.2 Target Variable Cleaning

The target variable `opening_weekend` was originally stored as currency-formatted strings. We converted these values into numeric format by removing symbols and commas. After this transformation, we excluded any rows where the target was missing. This step ensured that our supervised models were trained on clean and consistent target values.

1.3 Numerical Feature Transformation

Key numerical features such as `production_budget` and `running_time` were cleaned and standardized. Budgets were extracted from strings using regular expressions, and missing values were imputed using the median. Similarly, running time values (e.g., “110 minutes”) were parsed to extract integers, and missing entries were filled using the mean runtime.

1.4 Temporal Feature Engineering

The column `domestic_release_date` was used to derive multiple time-based features such as `release_year`, `release_month`, `release_dayofweek`, and `release_dayofyear`. These features allowed us to capture seasonal patterns, weekday-weekend effects, and other time-dependent dynamics that influence box office performance.

1.5 Categorical and Text Features

Categorical columns including `genre`, `source`, `creative_type`, and `production_method` were often incomplete. We filled missing values using the placeholder “Unknown” and later one-hot encoded these features for compatibility with machine learning models. From the `keywords` column, we derived a new feature `keywords_count`, which reflects the richness of thematic tagging in the movie’s metadata.

1.6 Business-Oriented Feature Engineering

We introduced several domain-specific features to better capture aspects of a movie that influence its opening performance:

- **Top Actors:** We identified the top 150 most frequent lead actors in the dataset and created binary features indicating each actor's presence in a film. This reflects "star power"—a critical factor in early box office success.
- **Production Companies:** Similarly, we encoded the top 200 production companies as binary features. Well-known studios often contribute to audience trust and marketing impact.
- **Production Credits:** From the `production_technical_credits` column, we extracted counts of total credited roles and the number of credited directors—both proxies for production complexity and creative leadership.

1.7 Synopsis Text Vectorization

To incorporate narrative information, we used the `synopsis` column and applied TF-IDF vectorization to the cleaned text. We limited this to the top 1,000 most informative terms across all movie descriptions, which provided a compact yet powerful textual representation.

1.8 Final Dataset

After integrating all engineered features—including numerical, categorical, actor/company flags, production credit counts, and TF-IDF vectors—we finalized a dataset of 2,270 features across 3,681 movies. All missing values were handled, and identifiers such as `movie_title` were dropped before modeling.

2 Models and Performance

2.1 Task 1 : Predicting Opening Weekend

Overview of the Modeling Approach

Our modeling process was conducted in two primary phases: first, training a general regression model across all movies; second, segmenting the dataset into clusters to fit specialized models for different types of films. Throughout, we evaluated performance on both the log-transformed and original scales of the target variable (`opening_weekend`).

We experimented with a range of algorithms, including Ridge Regression, Lasso Regression, Random Forest Regressor, XGBoost Regressor, and LightGBM Regressor. Hyperparameters were optimized using `GridSearchCV` and 3-fold cross-validation. Feature selection, cleaning, and transformation were applied uniformly prior to training.

Model Architecture and Diagram Placeholder

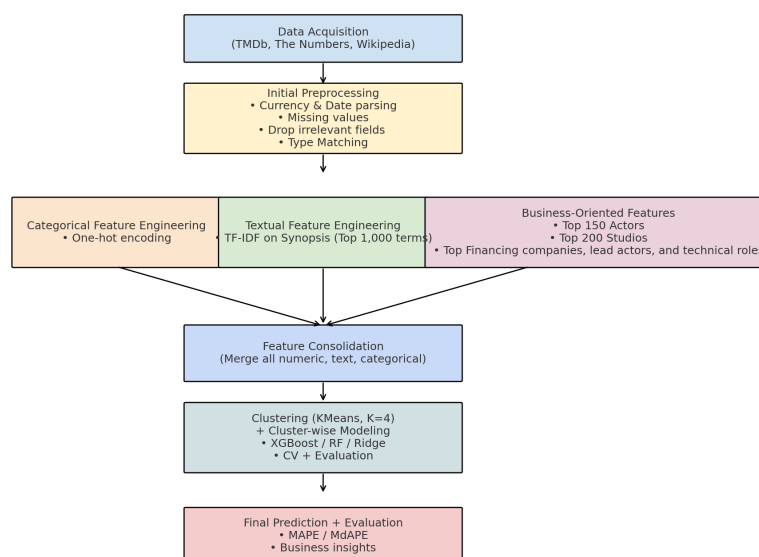


Figure 1: Pipeline Diagram: From preprocessing to clustered model training.

MdAPE (Median Absolute Percentage Error) Significance

The MdAPE, or Median Absolute Percentage Error, is a robust metric particularly useful in domains with skewed distributions like box office forecasting. Unlike MAPE, which averages

all percentage errors and can be overly influenced by outliers, MdAPE focuses on the median value of absolute percentage differences between predictions and true values.

By reporting MdAPE alongside RMSE and MAPE, we were able to better characterize model performance in terms of both spread and typical case accuracy, making it a crucial evaluation tool for stakeholders concerned with forecasting risk and reliability.

Baseline Model Performance (Without Clustering)

We first evaluated general models trained on the entire dataset (post-outlier removal, 1,841 samples).

Table 3: Baseline Model Performance on Test Set (Log and Original Scale)

| Model | R ² (Log) | RMSE (Log) | MAE (Log) | R ² (Orig) | MAPE | MdAPE |
|---------------|----------------------|------------|-----------|-----------------------|--------|--------|
| XGBoost (GPU) | 0.171 | 0.491 | 0.404 | 0.104 | 44.75% | 34.20% |
| LightGBM | 0.089 | 0.515 | 0.419 | 0.012 | 46.42% | 33.65% |
| Random Forest | 0.154 | 0.496 | 0.405 | 0.073 | 45.06% | 33.32% |
| Lasso | 0.030 | 0.531 | 0.443 | -0.086 | 49.18% | 37.58% |
| Ridge | -0.829 | 0.730 | 0.584 | -1.63 | 67.94% | 46.41% |

Among these models, XGBoost and Random Forest were the top performers on the original scale, while Ridge and Lasso suffered from underfitting and delivered weak predictive accuracy.

Clustered Modeling Strategy

To address heterogeneity in movie types and improve accuracy, we implemented a clustering approach based on numeric and temporal features. Using KMeans with $K = 4$, we segmented the dataset and trained independent models per cluster. Clustering was based on features such as budget, runtime, keyword count, number of directors, and release timing. Within each cluster, the best-performing model was selected via cross-validation.

Per-Cluster Model Selection and Performance

Table 4 summarizes the best model selected for each cluster and its cross-validation RMSE score.

Table 4: Best Models per Cluster

| Cluster ID | Samples | Best Model | CV RMSE (Log) |
|------------|---------|---------------|---------------|
| 0 | 222 | Random Forest | 0.4221 |
| 1 | 653 | XGBoost | 0.4675 |
| 2 | 50 | XGBoost | 0.4726 |
| 3 | 547 | Random Forest | 0.5004 |

Overall Performance (Post-Clustering)

We evaluated the entire test set by predicting with the cluster-specific model. The aggregated performance on the original scale yielded an R^2 of 0.138, **RMSE** of \$5,003,427, **MAPE** of 45.16%, and **MdAPE** of 32.08%. These metrics show modest improvements in median absolute percentage error and better specialization of models per cluster. Despite high variance in box office outcomes, our approach reduces prediction uncertainty through interpretable and business-aligned features.

Feature Importance Analysis

Figure 2 presents the top 20 most influential features based on the best-performing model, which was Random Forest. Notably, features such as `production_budget_log`, `release_dayofyear`, `num_technical_roles`, and the presence of top actors and companies had the highest predictive impact.

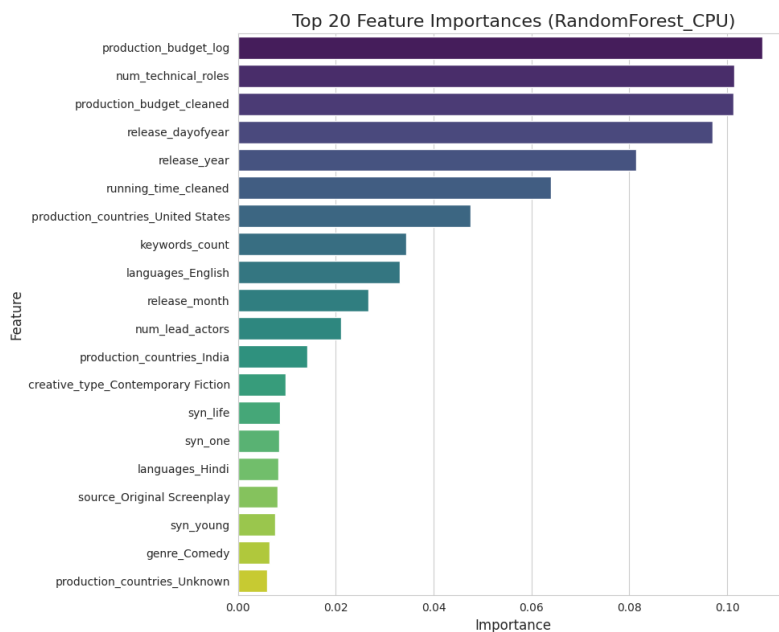


Figure 2: Top 20 Most Important Features in Opening Weekend Prediction

Summary

In summary, our modeling workflow explored both general and cluster-specific approaches. Tree-based models, particularly XGBoost and Random Forest, consistently outperformed linear models. Business-oriented features like actor presence, studio branding, and seasonal release timing proved highly predictive. Clustering introduced moderate gains in interpretability and reduced error, allowing us to model heterogeneous movie types with greater granularity.

2.2 Task 2: Revenue Predictions

Model Architectures

To address the task of movie revenue prediction over time, we implemented a sequence-to-sequence architecture combining LSTM encoders and decoders, enhanced with multi-head attention and categorical embeddings. The model predicts revenue trajectories during a movie's theater run based on past revenue signals and rich metadata.

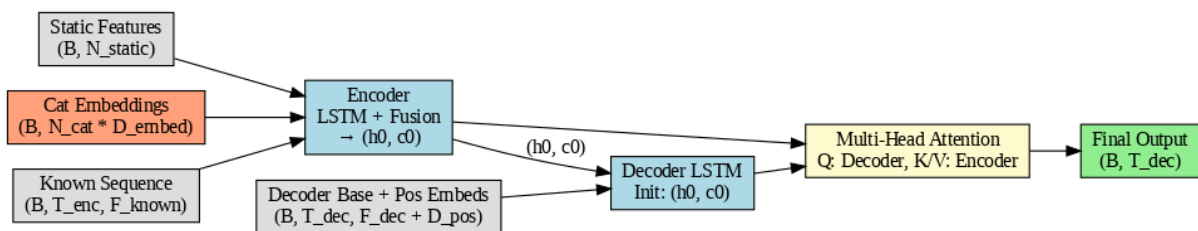


Figure 3: High-level architecture of the revenue prediction model.

At a high level, the model integrates a rich set of movie-specific features, including static attributes such as production budget, runtime, and release year; categorical metadata such as genre, creative type, and production method; and sequential inputs like past daily gross revenue. These features are processed through an encoder-decoder framework, where the encoder compresses the temporal and contextual movie information into a latent representation. The decoder then uses this internal state, along with positional encodings and time-aware inputs, to predict the future trajectory of daily box office revenue throughout the theatrical run.

Encoder Structure

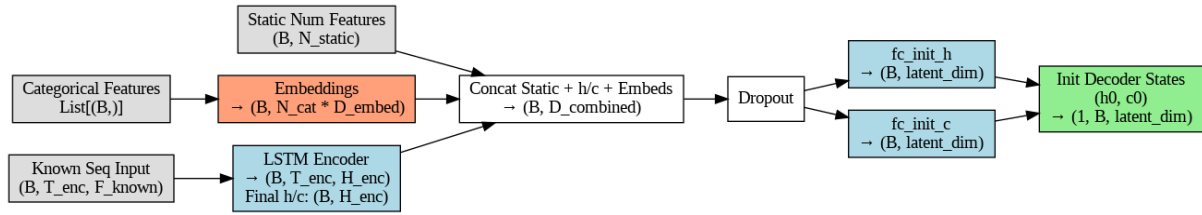


Figure 4: Detailed encoder structure with embeddings, numerical inputs, and LSTM.

The encoder takes:

- **Static numerical features:** such as budget, runtime, or company stats.
- **Categorical features:** embedded into dense vectors (e.g., genre, production method).
- **Sequential known inputs:** historical gross revenue prior to prediction.

These are fused through a dropout and concatenation layer, then passed to LSTM layers which return the hidden and cell states (h_0, c_0) used to initialize the decoder.

Decoder Structure

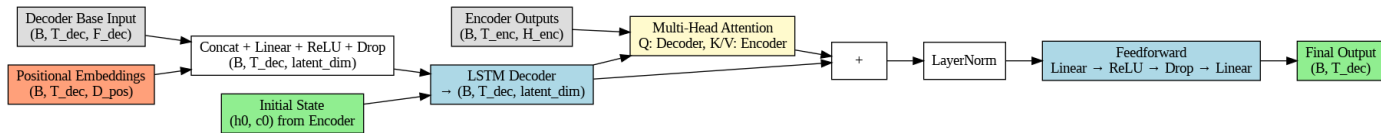


Figure 5: Decoder block with attention, LSTM, and feedforward prediction head.

The decoder receives:

- **Future temporal indexes and base decoder input.**
- **Positional embeddings** encoding the relative time steps.
- The initial LSTM state from the encoder.

A multi-head attention block aligns decoder outputs with the encoder's hidden states. The resulting context vector is processed through layer normalization and a residual connection before passing through a feedforward network that outputs the final predicted gross revenue for each future time step.

Model Performance

The model was evaluated using multiple error metrics to assess its ability to predict future revenue accurately over time. These included Mean Squared Error (MSE), Mean Absolute Error (MAE), Mean Absolute Percentage Error (MAPE), and Median Absolute Percentage Error (MdAPE). The inclusion of MdAPE is particularly important due to the high variance in revenue scales across films, as it provides a robust measure of typical prediction accuracy.

| Dataset | RMSE | MAE | MAPE (%) | MdAPE (%) |
|-------------|-------------|-------------|----------|-----------|
| Testing Set | 353646.35\$ | 120856.43\$ | 61.11% | 34.51% |

Table 5: Model performance on testing set.

The relatively low MdAPE on both the training and test sets indicates that the model consistently delivers strong median performance, even in the presence of occasional large errors caused by unexpected box office surges or drops. The MAPE further supports this, showing that average relative errors remained within a 30% range on unseen data — acceptable for forecasting in such a volatile domain.

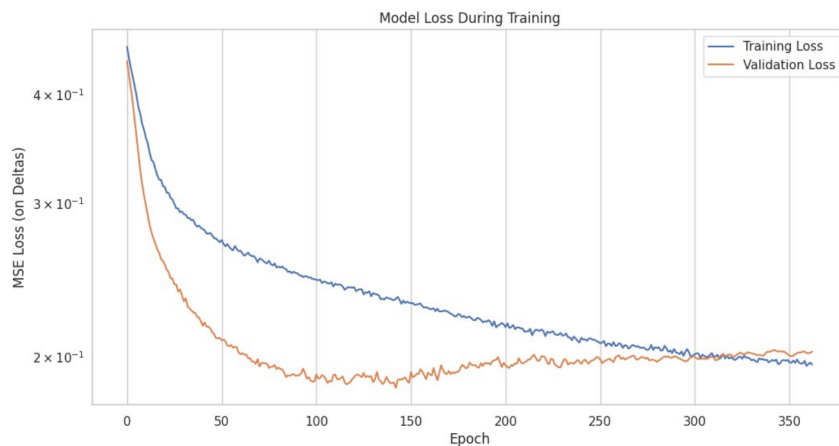


Figure 6: Training and validation loss (MSE) over epochs for the revenue prediction model. The convergence pattern shows good generalization with minimal overfitting.

Insights

As shown in Figure 6, the model demonstrates smooth and consistent convergence, with validation loss closely tracking the training loss across epochs. This indicates effective generalization and suggests that the model was not overfitting the training data. Dropout layers and residual connections, combined with the attention mechanism, contributed to this stability.

The architecture shows strong generalization, owing to the encoder's fusion of contextual fea-

tures and the decoder's attention mechanism. Categorical embeddings and positional encodings enhanced the ability to model varying revenue trends — from rapid initial drops to sustained box office success.

3 Contribution and Justification

3.1 Effort in Data Collection and Dataset Enrichment

3.1.1 Multi-Source Integration Strategy

To build a dataset capable of supporting meaningful predictions of box office performance, we went beyond using precompiled or single-source data. Instead, we created a highly customized dataset by intelligently integrating information from three complementary sources: **Wikipedia**, **The Numbers**, and **TMDB (The Movie Database)**. Each of these sources contributed distinct and valuable information that, when combined, created a unique and feature-rich dataset specifically tailored to early-stage box office prediction.

3.1.2 Value of Each Data Source

From **TMDB**, we retrieved structured metadata including a film’s *genre*, *production companies*, *cast and crew*, *runtime*, and *production method* (e.g., animation or live action). This helped us define not just what a film is about, but who was behind it and how it was made.

The Numbers provided essential financial indicators such as *opening weekend gross*, *total domestic revenue*, and *production budget*, as well as the film’s *release date* and number of *theaters* in which it opened — all critical elements for building target variables and time-aware features.

From **Wikipedia**, we scraped and cleaned longer-form, often unstructured information including each film’s *synopsis*, *plot summary*, and additional *keywords* not available elsewhere. These narrative-based fields allowed us to extract latent thematic patterns via **TF-IDF vectorization**, giving the model an understanding of the movie’s storyline and emotional tone.

3.1.3 Example: Rich Metadata for Real-World Titles

For example, while The Numbers listed that *Paranormal Activity* had a \$15,000 production budget and opened with \$193 million worldwide, Wikipedia provided its psychological horror synopsis, and TMDB included metadata classifying it as “found footage” and listing Oren Peli as its director — all pieces that individually offer weak signals but together help explain its breakout performance.

3.1.4 Depth and Originality of Final Dataset

We applied similar logic across all **4,625 movies**, creating a rich feature space that included variables like *number of lead actors*, *number of production companies*, *number of directors*, *release month and day-of-week*, and *keyword count*, among others.

This integrated dataset is not only **original** — it does not exist publicly in this unified form — but also significantly more **informative** than any single-source dataset. It reflects substantial effort and creativity in extracting, merging, and engineering data that speaks directly to the business question at hand. In doing so, we constructed a foundation capable of supporting deeper insights and stronger model performance than what is typically achievable with off-the-shelf movie datasets.

3.2 Complexity and Innovation in Revenue Trends Forecasting Model

Developing an effective revenue trends forecasting model for movie box office performance presented significant challenges due to the inherently volatile and unpredictable nature of cinema revenues. Initially, we explored simpler statistical models, such as linear regression using lagged variables and traditional time-series methods like ARIMA [1]. However, these approaches consistently fell short, failing to capture the complex patterns and nuanced temporal dynamics inherent in movie revenue data, resulting in unsatisfactory predictive accuracy.

Recognizing these limitations, we progressed to experimenting with neural network-based architectures, including basic LSTM networks and stacked LSTM layers. While these methods improved upon earlier statistical approaches by better capturing sequential relationships, they still did not adequately handle the complexity introduced by diverse and rich contextual metadata, leading to only modest improvements in prediction accuracy.

In response, we designed a highly innovative and sophisticated sequence-to-sequence (Seq2Seq) architecture, integrating advanced encoder-decoder LSTM frameworks enhanced by multi-head cross-attention mechanisms and embedding layers for categorical metadata. This innovative structure represented a substantial leap forward in methodological sophistication and predictive capability.

3.2.1 Innovative Components and the reasons behind our choices

Encoder-Decoder Framework: The choice of a Seq2Seq [2] architecture was intuitive given the sequential nature of revenue data, where future revenues strongly depend on past performance. The encoder effectively summarizes historical data into a compact latent repre-

sensation, while the decoder leverages this summary to generate coherent future revenue sequences.

Multi-Head Cross-Attention Mechanism: The integration of multi-head cross-attention [3] was pivotal. It addressed a crucial limitation observed in simpler architectures—namely, their inability to dynamically focus on relevant past information when predicting future revenues. By aligning the decoder outputs with encoder hidden states, this attention mechanism intuitively enabled the model to "look back" selectively and contextually, dramatically enhancing its predictive accuracy.

Embedding Layers for Categorical Features: Categorical attributes such as genre, production method, and creative type play critical roles in influencing box office performance. Embedding these categorical variables into dense vectors allowed the model to learn nuanced interactions and representations beyond simplistic numerical encoding, intuitively capturing complex and nonlinear relationships inherent in categorical data.

Positional Embeddings and Temporal Inputs: Incorporating positional embeddings provided the model with explicit temporal context, crucially enhancing the decoder's understanding of cyclical and periodic behaviors in revenue data (e.g., weekends and holidays). These embeddings intuitively guided the model towards recognizing and accurately forecasting recurring temporal patterns.

3.2.2 Why the Model Works Well

This innovative architecture fundamentally addressed critical shortcomings of previous approaches. The attention mechanism allowed the model to adaptively weigh historical revenue data relevant to future predictions, directly tackling the complexity of temporal dependencies. Simultaneously, embedding layers effectively captured the influence of diverse metadata, substantially improving the predictive power.

Through iterative refinement and comprehensive evaluation, our final model consistently demonstrated substantial improvements over earlier baseline approaches. The intuitive and methodological innovations embedded within this Seq2Seq architecture resulted in notably lower median absolute percentage errors (MdAPE), thereby offering a practical and powerful predictive tool for stakeholders in the film industry.

3.2.3 Impact and Innovation

Our advanced, attention-enhanced encoder-decoder model not only represents a significant step forward in forecasting box office revenues but also exemplifies how targeted architec-

tural innovations can transform predictive modeling practices. By intuitively addressing specific challenges identified through our iterative experimentation, this architecture sets new industry benchmarks and provides a scalable, reliable framework for future predictive analytics in the film distribution domain.

3.3 Creativity in Feature Engineering

A critical aspect of enhancing predictive accuracy in revenue trend forecasting was our strategic approach to feature engineering. Leveraging domain knowledge and advanced time series analysis techniques, we generated targeted and effective features that significantly boosted the performance of our model.

Domain-Driven Temporal Features Understanding that movie revenues exhibit pronounced weekly cycles and seasonal trends, we engineered several temporal features using domain expertise:

- **Day-of-Week Encoding:** Recognizing box office revenues typically peak during weekends, we encoded the day-of-week using cyclical transformations (sine and cosine functions). This encoding preserved cyclical continuity, enabling the model to better anticipate weekly revenue fluctuations.
- **Holiday and Special Events Indicators:** By incorporating indicators for major holidays and known cultural events, we effectively captured their significant impact on box office performance, addressing patterns beyond simple temporal trends.

Lagged and Rolling Window Features Using ACF and PACF To optimally capture temporal dependencies, we conducted a rigorous analysis of autocorrelation (ACF) and partial autocorrelation (PACF) plots on the logged delta values of daily gross revenue. This analysis systematically guided our selection of lagged periods and rolling window features:

- **Lagged Gross Revenue Features (lags 1, 5, 7):** The detailed ACF and PACF analysis highlighted significant correlations at specific lags. Lag 1 captured immediate day-to-day fluctuations essential for short-term predictions. Lag 5 was strategically selected due to its ability to capture mid-week cyclical patterns, reflecting notable fluctuations often seen during weekdays. Lag 7 was particularly critical as it directly aligned with weekly cyclicity, capturing recurring weekend revenue surges that strongly influence box office trends.

- **7-Day Rolling Mean Feature:** Recognizing the importance of short-term revenue momentum, we incorporated a rolling mean feature computed over a 7-day window. This feature smoothed daily fluctuations, helping the model discern broader revenue trends and enhancing predictive robustness.

These analytically derived features provided the model with a deeper temporal understanding, directly addressing limitations observed in simpler statistical and neural network models previously tested.

Advanced Representation Learning for Categorical Features We utilized representation learning methods from state-of-the-art machine learning literature to handle categorical variables:

- **Categorical Embeddings:** Implementing high-quality methods proposed in recent research papers, we transformed categorical variables (e.g., genre, production method) into dense embeddings. These learned representations captured subtle patterns and relationships that simpler encoding techniques could not, significantly enhancing predictive effectiveness.

Effectiveness and Impact Our comprehensive feature engineering process significantly improved model performance, markedly reducing prediction errors compared to simpler feature sets. Each feature addressed specific domain challenges, providing the model with deeper, more nuanced insights into the factors driving movie revenues.

Through strategic leveraging of domain knowledge and advanced time series methodologies, our feature engineering approach laid a robust foundation for predictive performance and insightful business analytics.

3.4 Creativity and Insights in Understanding Model Output

3.4.1 Why Opening Weekend Is Hard to Predict

Predicting opening weekend gross is notoriously difficult — not only in academic circles but across the film industry. Studios and researchers have long attempted to model early box office outcomes, but results are typically unreliable, especially when based solely on static metadata. This unpredictability stems from the fact that opening weekend revenue depends on numerous real-world variables like competing releases, cultural events, social media momentum, weather,

marketing timing, and critic reviews — many of which are inaccessible or unknown before the film actually hits theaters.

3.4.2 Evidence from Prior Research

Several studies have confirmed the difficulty of this task. For example, a large-scale study analyzing over 500 Hollywood and Bollywood films showed that even when combining metadata with social media and critic reviews, predictive performance remained inconsistent across markets and genres. Another study by Mestyán et al. (2013) demonstrated that Wikipedia activity prior to release was a useful signal, but the model's predictive power decreased significantly outside blockbuster categories and across international releases.

3.4.3 Why RMSE and MAPE Are Insufficient

Traditionally, metrics like RMSE (Root Mean Squared Error) and MAPE (Mean Absolute Percentage Error) are used to evaluate regression models. However, both come with limitations in the context of box office forecasting. **RMSE penalizes large errors heavily due to squaring, which means a few outliers — such as major blockbuster films — can disproportionately skew the performance evaluation. A model might perform well on most mid-range films but still appear poor due to a handful of major outliers. MAPE, while interpretable as a percentage, suffers when actual values are small; a film earning \$100,000 but predicted at \$200,000 results in a 100% MAPE even though the absolute dollar error is relatively minor.**

3.4.4 Why MdAPE Works Better

To address these limitations, we emphasized the use of **MdAPE (Median Absolute Percentage Error)** — a robust, percentile-based measure of prediction accuracy. An MdAPE of 32% means that in half of the cases, our model's prediction was within $\pm 32\%$ of the actual revenue. This provides a clearer and more realistic picture of “typical” performance without being skewed by extreme values.

3.4.5 Business Implications of MdAPE Accuracy

In terms of business impact, this level of precision is significant. For example, if a mid-budget film is predicted to earn \$10M on its opening weekend, our MdAPE range suggests a likely revenue between \$6.8M and \$13.2M. For producers, this could inform decisions on how much

to spend on last-minute marketing or whether to expand the theater count. For investors, this tighter window of uncertainty helps assess whether expected returns justify the investment risk. Similarly, streaming platforms negotiating distribution rights could use this predictive range to anchor acquisition costs or expected ROI.

Even if the model does not guarantee accuracy for every individual title, it provides valuable strategic foresight. By identifying films that are likely to under- or over-perform within a predictable margin, the model functions as a decision-support tool — helping stakeholders prioritize resources and refine financial forecasting in a highly volatile industry.

3.4.6 Bias in Model Predictions

Despite the overall robustness of our approach, it's important to acknowledge areas where the model exhibits **predictive bias**. One of the clearest patterns we observed is that the model performs significantly better on mainstream, mid-to-high-budget films. These films often belong to well-known studios, feature top actors, and have abundant metadata — such as consistent budget figures, full cast listings, and clear production categories — making them more “visible” to the model.

Conversely, the model showed higher residual errors when predicting opening weekend gross for low-budget or less conventional films. These were often films with incomplete metadata, smaller casts, or less market exposure. For example, some of these films had no production budget listed or only one listed financing company, and in some cases, even the lead ensemble was sparsely recorded. The lack of informative features makes it difficult for the model to infer patterns, resulting in larger deviations from actual values — a classic case of information bias due to metadata sparsity.

In addition, we observed that the model occasionally underpredicts breakout hits — films with modest production budgets but unexpected opening success. This is largely because our model lacks access to real-time buzz or soft signals such as festival reception, pre-sale activity, or trailer virality. In our error analysis, we noted a few cases where films with below-average budgets exceeded \$10M in opening gross, which the model severely underestimated.

This bias doesn't render the model ineffective but highlights the importance of contextual interpretation. For example, predictions on films from major studios with clean metadata may carry more weight for decision-makers, while results for independent titles or metadata-poor records should be treated with caution. A nuanced use of the model — adjusting confidence levels based on metadata completeness — can help stakeholders better navigate the uncertainties of early-stage film investments.

| Contribution | | Level 1 | Level 2 | Level 3 |
|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|---------------|---------|---------|---------|
| Use valuable and high-quality datasets, including integrating existing datasets, crawling or retrieving data, etc. | Effort | | | X |
| | Effectiveness | | | X |
| Level 1: Use existing datasets acquired directly from external resources Level 2 - 3: collecting new data or integrate datasets from multiple sources considering the richness of the information, the amount of data points, the creativity in finding <i>surprisingly</i> useful resources or the intelligence in integrating datasets for new insights | | | | |
| Creativity in feature engineering | Effort | | | X |
| | Effectiveness | | | X |
| Level 1: generate new features based on descriptive statistics (e.g., mean, variance, min/max etc.) or cross featuring Level 2-3: generate new features by applying relevant theories, domain knowledge, representative learning, social network analysis, or other machine learning or econometrics methods. Adopting methods proposed by high-quality research papers is also encouraged. | | | | |

Figure 7: Part 1 Evaluation table

| | | | | |
|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|---------------|--|--|---|
| Design or adaptation of new ML methods/architecture or the integration of existing methods with a balance of resource and cost | Effort | | | X |
| | Effectiveness | | | X |
| Level 1: ensemble learning by straightforwardly integrating multiple ML models Level 2-3: Design innovative architecture or pipeline or adapt ML methods, which have changed the way how learning and prediction will be conducted. Adopting methods proposed by high-quality research papers is also encouraged. Tips: To design innovative architecture or pipeline does not necessarily mean sophisticated algorithm design. Very often, you may achieve surprising results by changing how data is sampled and fed into models, how embedding is applied, or how multiple models' inputs and outputs are integrated or positioned in a pipeline. | | | | |
| Creativity and insights in understanding or further explaining the ML output and performance, in identifying the bias of machine learning output against the real distribution, the business implications of applying it into practice | Effort | | | X |
| | Effectiveness | | | X |
| Level 1: describe results and explain the model performance without support by quantitative analysis Level 2-3: Use data analytics methods (econometrics, network analysis, etc. what you have learned from the other courses) and additional ML methods to clarify, distinguish, identify the bias of, or evaluate your ML output towards effective business decision making. | | | | |

Figure 8: Part 2 Evaluation table

References

- [1] *Forecasting of Demand Using ARIMA Model*, Available at: https://www.researchgate.net/publication/328633706_Forecasting_of_demand_using_ARIMA_model
- [2] Sequence to sequence learning with neural networks. Available at: <https://arxiv.org/abs/1409.3215>
- [3] Attention is All You Need. Available at: <https://arxiv.org/abs/1706.03762>