# A Predictive and Visualization Tool for Filmmakers

## CSE 6242 A    Team 20: CineSage Insights

Aditya Vikram
avikram33@gatech.edu

Kshitij Pathania
kpathania3@gatech.edu

Shrestha Mishra
shrestha.mishra@gatech.edu

Sultan Syed
ssyed71@gatech.edu

## 1 Abstract

This project aims to uncover a novel pathway for filmmakers to make data-driven decisions using a predictive and visualization tool. This tool takes into account variables such as plot premise, budget, star cast, release time-frame, and targeted geographies to predict and visualize a movie's success and revenue. This has the potential to redefine movie production planning by combining data analysis and visualizations.

## 2 Introduction

Production houses allocate significant time, energy, and financial resources to the pre-production phase of a movie. This phase involves conducting detailed market analysis to gauge the movie's potential reception at the box office, estimate its revenue, and identify the ideal cast members. However, these analyses are typically done manually and demand considerable effort. Estimating a film's potential is critical for planning aspects of the production phase such as script refinement, budget allocation, and casting decisions.

Our project is motivated by a desire to revolutionize the conceptualization, planning, and production of films. Traditionally, these decisions have relied on intuition, past experiences, and subjectivity, sometimes overshadowing the importance of data-driven decision-making.

By leveraging machine learning and data analysis, our application can analyze extensive datasets, including plot summaries, cast choices, budgets, and release dates, to provide production houses with an accurate data-driven revenue prediction model.In doing so, we aim to:

- _Minimize Risk_: Our tool estimates revenue potential, enabling production houses to make informed decisions, allocate resources efficiently, and minimize risks.
- _Enhance Creativity_: Our application's abilities liberate creative minds to focus on storytelling and artistic expression, with the assurance of data-driven guidance.
- _Improve Efficiency_: This tool streamlines decision making, saving time and resources by pinpointing optimal casting choices and ideal release windows.

## 3 Problem Definition

The aim of this project is to create a tool that helps filmmakers understand how well their movie might do in terms of audience reception and revenue. More formally we aim to predict the following:

- _Revenue_: We'll predict potential earnings for the movie in each country based on factors like the budget, plot, cast, and release date.
- _Audience Preference_ Our tool categorizes audiences by analyzing user ratings for various films. This data guides filmmakers in strategic marketing and efficient film promotion, with support from interactive t-SNE visualizations of user clusters.

## 4 Literature survey

Our project can be divided into two broad categories. first would be training a simple deep learning model that will help us predict revenue collection of a movie and second is visualizing those data. We have explored papers that are using deep learning models on several movies datasets to solve, clustering and labeling related problems and differnt techniques used to visualize those results. Numerous studies have been done on various movie-related datasets and deep learning models. The most prevalent ones are sentiment analysis on movie reviews. Performance of neural networks like long short-term memory(LSTM) and convolutional neural networks have been compared on the IMDB dataset with 50k reviews[2]. This paper talks about how a combination of CNN and LSTM models outperform LSTM, CNNs, SVMs and Naive bayes by approximately 2.5%. Support vector Machines along with index-based feature selection has been shown to reduce error rate in sentiment. [7] analysis[14].

People have also applied Markov model based approach[18] to predict violence in a movie using Violent Scene Dataset(VSD) [8]. Work has also been done around summarizing the plot of a movie and predicting it's genre by analysing posters, trailers and key scenes from a movie [5].

Several studies have proposed movie recommendation system that leverages deep learning networks and uses the MovieLens, Book-crossing, Each-Moviedataset[20].[13] uses k nearest neighbour to predict user rating which is then used to recommend movies to them. [9]uses memory based collaborative filtering algorithm to capture weighted average of similar user rating that is then used to recommend movies to users.

A lot of techniques have also been incorporated to predict the box-office revenue of a movie. LSTM models have been used to predict daily box office collections based on social

media interaction and reviews[17].[3] and [15] use data mining techniques to predict lifetime box office collection of a movie. predictions are dependent on historical database and sentiment analysis on Twitter and YouTube.

As an improvement over using LSTMs [4] for encoding the movie plot, we plan to explore BERT [10] which uses transformer architecture with bidirectional attention. BERT is shown to perform better than LSTMs for NLP tasks like sentiment analysis, and we plan to fine-tune pre-trained BERT models on our dataset.

Pertaining to data visualization in the film industry, [12] introduced a static 2D multi-layered visualization method to offer a global perspective on movies and the actor space. In their approach, the base layer organizes movies in columns, color-coding them based on their genre, and superimposes a force-directed graph representing the co-actor network. While their work provides valuable insights for movie enthusiasts, our project focuses on creating interactive visualizations—such as Choropleth maps and node graphs—to empower filmmakers to fine-tune their movie plots and budgets by leveraging our predictive models, ultimately facilitating data-driven decision-making for movie production studios. [11] offers a comprehensive sentiment network visualization of movie reviews, connecting a movie node to popular sentiment nodes extracted from semantic data within its reviews. Our project differs in its focus on movie ratings, which provide a substantial advantage in terms of quantity (a significantly larger volume of ratings are available compared to reviews). We believe that the combination of rating numbers, user-movie interactions, and the network of users and their rated movies will suffice for accurately predicting the user categories likely to appreciate an upcoming movie.

[16] employed a range of visualizations, including heatmaps, choropleth maps, and scatter plots, to visualize events like kills and conquests in the James Bond series. They also utilized statistical models to predict the revenues of future James Bond movies. In contrast, we aim to leverage a significantly larger movies dataset for training our predictive models and predict revenue based on an upcoming movie's plot. Furthermore, our key differentiation lies in creating interactive visualizations tailored to the needs of filmmakers.

## 5    Proposed Method

We present an advanced data-driven visualization suite designed to enhance decision-making for filmmakers, leveraging a comprehensive dataset that includes movie plot summaries, budget information, genres, cast details, and historical revenue data from multiple countries. Initially sourced from 'The Movies Dataset' [6], it was augmented with international revenue data via web scraping techniques (detailed in Subsection 5.3.1).

### 5.1    Data Preprocessing and Analysis

To convert textual data into a usable format for our models, we utilize a fine-tuned BERT model from HuggingFace, known for its superior performance in text-related tasks. Numerical representations of the text are obtained by embedding movie plot summaries through BERT, and these embeddings serve as input features for our predictive models.

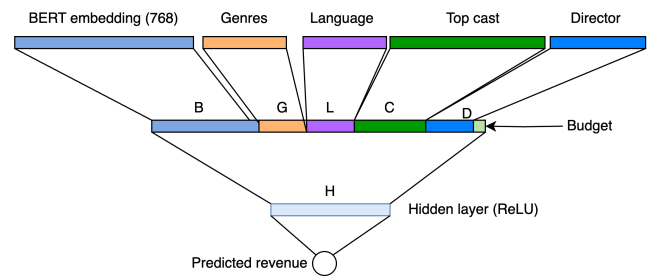### 5.2    Predictive Revenue Visualization

The predictive revenue visualization aims to project revenue distribution across different countries. A neural network architecture, as shown in Figure 1, takes as input:

- the BERT-generated plot embeddings
- multi-hot vector for movie's genres
- one-hot vector for the movie's original language
- movie's budget
- multi-hot vector denoting the top cast members
- multi-hot vector denoting the top crew members

These inputs are first individually embedded using linear layers to reduce dimension, then stacked and embedded again into a hidden layer. Another linear layer maps the hidden layer to the predicted revenue. Each hidden layer uses ReLU activation, and we also test whether dropout improves performance.

Top cast and crew members are identified from the pre-processed dataset using their popularity, and we pick the top 1000 cast members and top 500 crew members for one-hot encoding. This was done so that the training data fits in GPU memory.

For each country with sufficient data, we develop specific models to account for the unique patterns in movie revenue distribution. Training is carried out using Mean Squared Error (MSE) loss, and we plan to optimize the intermediate layer sizes ($B$, $G$, $C$, $L$, and $H$) to further enhance predictive performance.



**Figure 1.** Revenue prediction model architecture (each trapezoid represents a Linear embedding layer). Larger image available here.

The interactive visualization component allows users to adjust parameters like budget and genres, and observe the

changes in revenue distribution in real-time on a choropleth map.

## 5.3 Real-Time Genre Prediction

For an improved user experience, we incorporated a real-time genre prediction model that can predict genre(s) based on the plot overview provided by the user. This feature allows users to work in a dynamic environment, receiving real-time updates with minimal input requirements.

To predict genres, we employed a pre-trained DistilBERT [19] model, fine-tuned for the multilabel classification task of predicting genres from a set of 20 available genres in the dataset. DistilBERT is a transformer that was pre-trained using knowledge-distillation on BERT to reduce the model size by around 40% while still retaining 97% of BERT's language understanding capabilities. Being around 60% faster than BERT, it is a better suited model for our use case, because we need real-time genre predictions based on user's inputs. The implementation utilized Hugging Face transformers [22] to build, train, and evaluate the model.

### 5.3.1 Web Scraping. Extensive scraping was conducted to gather country-specific revenue data, enriching our dataset with more granular insights. IMDb IDs were used to extract data from BoxOfficeMojo [1], yielding country-wise revenue figures for a significant number of movies. Models are specifically trained on countries with ample data to ensure robust predictive accuracy.

## 5.4 User Clustering

To facilitate effective movie promotion strategies, we utilize user ratings to generate genre-centric user preference vectors. Applying K-means and Gaussian Mixture Models, we cluster these preferences, using silhouette scores to ascertain the optimal number of clusters.

t-SNE embedding visualizes the clusters in a two-dimensional space, with node size indicating cluster size and distance representing similarity in preferences. Detailed statistics about the clusters' preferred genres are accessible upon interaction with the nodes. This process significantly aids filmmakers in identifying and targeting the audience segments most likely to engage with their content.

## 5.5 Innovations List

Incorporating cutting-edge technologies and data-driven strategies, our approach introduces the following innovative features:

- *Real-Time Predictive Choropleth Map*: Our interactive choropleth map updates revenue predictions in real-time based on user inputs, providing filmmakers with a powerful tool to visualize potential earnings and adapt their strategies.
- *Audience Preference Clustering with t-SNE Visualization*: We utilize t-SNE to transform complex user rating data into an

intuitive visual format, enabling efficient identification of target audience segments for tailored promotional efforts.

# 6 Evaluation and Experiments

In this section, we evaluate CineSage Insights, focusing on its data-driven approach to movie revenue prediction and audience targeting. We detail the experiments conducted and discuss the visualizations and results to be included in the report.

## 6.1 Experiment Setup

To prepare the dataset for modeling, we created multi-hot encoding vectors for movie genres and one-hot encoding for languages. Notably, a significant number of movies lacked budget data. To address this, we introduced a new features, allowing the model to make revenue predictions even when budget information is missing. Additionally, we scaled the budget and revenue columns by dividing them by 100 million to expedite model training with the AdamW optimizer.

All revenue prediction models were trained using PyTorch and Hugging Face transformers, using the Mean Squared Error (MSE) loss function and AdamW optimizer. All revenue prediction models were trained using PyTorch and Hugging Face transformers, using the Mean Squared Error (MSE) loss function and AdamW optimizer. Model training was done on free T4 GPUs provided by Colab.
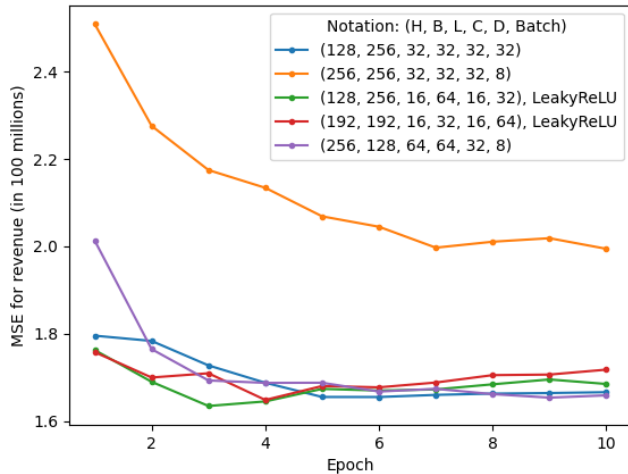
The real-time genre prediction model proposed in Section 5.3 was trained and evaluated on the pre-processed dataset using Hugging Face libraries. Binary cross-entropy loss per genre was used to fine-tune the transformer, and Micro-F1 score was used to pick the best fine tuned model. We used silhouette score to evaluate and find the optimal number of clusters for the second task. We now discuss each task in more detail.

- **Data Characteristics**: The primary dataset used is 'The Movies Dataset', augmented with extensive web-scraped revenue data. This dataset encompasses plot summaries, budget, genres, cast details, and revenue data across multiple countries, comprising approximately 950MB, with data on around 45,000 movies and over 20 million movie ratings.
- **Data Augmentation**: We conducted comprehensive scraping to gather country-specific revenue data, using IMDb IDs to extract information from BoxOfficeMojo. This provided granular insights into country-wise revenue figures for a significant number of movies.
- **Data Preprocessing**: BERT models from HuggingFace were utilized to convert textual data (like plot summaries) into numerical representations. These embeddings served as input features for predictive models. For cluster analysis we processed data by merging user ratings and movie metadata. The emphasis was on genre-based preferences, thus genres were extracted and converted into numerical

representations. Users' genre preferences were then aggregated and normalized to form a comprehensive user-genre dataset.

## 6.2 Predictive Revenue Visualization

- **Model Architecture**: The neural network architecture for revenue prediction integrates BERT-generated plot embeddings, multi-hot genre vectors, and one-hot language vectors, alongside budget data.
- **Country-Specific Models**: Separate models were trained for each country with sufficient data to capture unique revenue distribution patterns. We chose a threshold of having at least 500 movies for a country to be eligible for training. This gave us 17 countries.
- **Worldwide revenue prediction**: In addition, we also utilize the power of big data, and train one additional revenue prediction model on worldwide revenue. This prediction provides a holistic view for the user, since we can currently predict country-wise revenue for countries with sufficient data. The original dataset had this attribute for around 20% of the movies, and we performed additional scraping from BoxMojo to get the worldwide revenue collection for 50% of the movies. This provides ample data to train a large neural network that can accurately predict worldwide revenues.
- **Training and Optimization**: Figure 2 below displays the MSE loss curves for the model trained on the worldwide revenue column for 10 epochs (around 1.5 hrs), with varying batch sizes and other hyperparameters.[1]
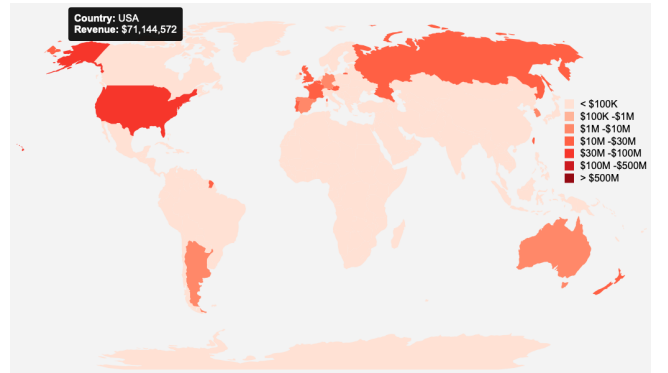


**Figure 2.** Loss curves for different hyperparameter combinations (*H*, *B*, *L*, *C* and *D* are detailed in Section 5.2)

- **Observations**: As expected, larger models tend to perform better in terms of converged MSE values because they are capable of expressing more complex functions.

[1]All trained models can be viewed and downloaded from OneDrive folder.

Replacing ReLU with LeakyReLU activation boosts the performance slightly. This could be because LeakyReLU prevents dead neurons in the network by providing a small negative output even if the preactivation is negative [21]. A larger batch size works better as it provides more stable gradients for the optimizer, leading to faster convergence. We similarly tune hyperparameters for country-wise prediction models too. Their analysis follows a similar trend and the plots are skipped in the interest of space.

- **Interactive Component**: An interactive visualization of the revenue prediction allows users to adjust parameters like budget and genres, observing real-time changes in revenue distribution on a choropleth map. We suggest genres based on the genre prediction model, allowing the user to approve suggestions. An example of the visualization is shown in Figure 3.
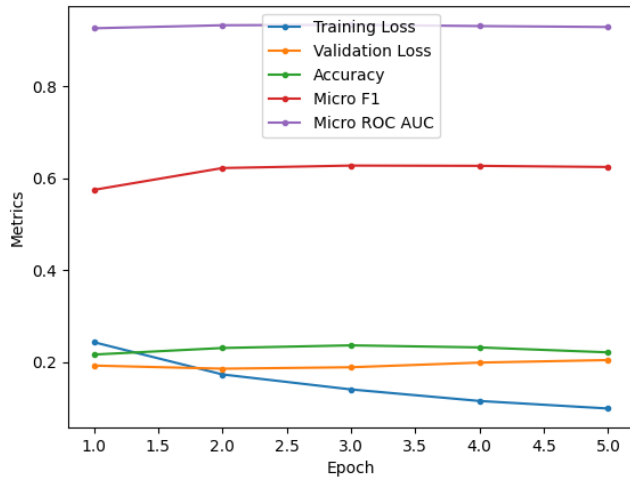


**Figure 3.** Choropleth Map showing revenue predicted by our model for various countries

**6.2.1 Genre prediction model results.** Figure 4 shows the variation of the F1 score, training and validation loss and ROC AUC score as the model is trained. We note that the training and validation loss seem to have converged and are close enough, implying that the model is not overfitting and training has converged. We achieve a high ROC AUC score (0.93) means that the model can effectively distinguish between positive and negative examples of each genre, based on the plot overview. This is supplemented by the accuracy of predicted genres for plot overviews from newer movies.

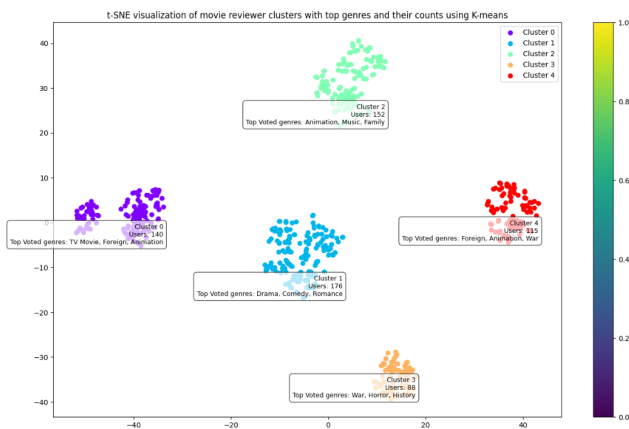## 6.3 Cluster Analysis and Visualization

Two clustering techniques were employed: K-Means and Gaussian Mixture Model (GMM). Each method aimed to cluster users based on their genre preferences. For K-Means, we experimented with different cluster numbers to find the optimal configuration using silhouette scores. The optimal silhouette score was achieved at 5 clusters. Figure 5 below visualizes those clusters using t-Distributed Stochastic Neighbor Embedding (t-SNE) to embed the data points in two dimensions. From the figure, we note that the clusters are

**Figure 4.** Evaluation metrics for the Genre Prediction model

well-separated and tightly packed, showing the efficacy of K-Means clustering on our dataset. Our client can interact with the visualization by hovering over the clusters to know the number of users in that cluster and the top-voted genres.
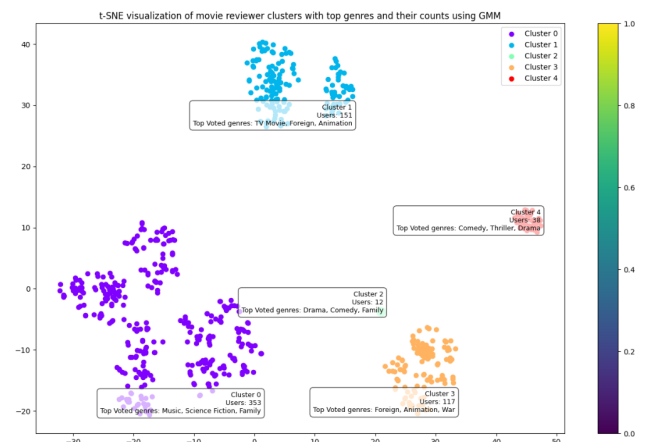


**Figure 5.** K-means clusters of movie critics based on their rating trends in different genre

GMM was similarly applied, with the intent of capturing more nuanced user groups and the results are visualized using t-SNE in Figure 6. This visualization provided a clear representation of how users were grouped according to their genre preferences. It also illustrated the density and separation of clusters, giving insights into the distinctiveness of each cluster. Our experiments revealed distinct clusters with unique genre preferences. The t-SNE visualizations effectively illustrated the separation and characteristics of these clusters. K-Means and GMM yielded different cluster formations, offering varied perspectives on user preferences. The analysis of genre-specific appeal in clusters provided

actionable insights into audience segmentation and content strategy.

**Visualization**: The clusters were visualized using t-SNE embedding, with no of nodes indicating the number of users in each cluster and distances representing the similarity in genre preferences. Users were able to interact with the visualizations to access detailed statistics about the preferred genres within each cluster. Each cluster contains the information of the reviewers population and their three top rated genres.



**Figure 6.** GMM clusters of movie critics based on their rating trends in different genre

To further illustrate the practical application of CineSage Insights, a live demonstration of the visualization suite is available at CineSage Insights Live Demo. This demo showcases the interactive features and real-time capabilities of our tools in action. The CineSage Insights platform is hosted on Azure's free tier (F1) under the App Service plan (ASP-DefaultResourceGroupeastus2-ab7b) with a Linux operating system. This setup, while cost-effective, comes with certain limitations, particularly in deploying heavy machine learning models. As a result, we have utilized dummy data for the live demonstration to showcase the functionality of the suite. However, the complete machine learning pipeline is fully integrated in our local setup. For a more comprehensive experience, the suite can be deployed locally by following steps in our Github repository's README. This enables users to interact with the actual machine learning models and obtain real-time predictions and insights.

The insights garnered from this analysis are invaluable for filmmakers and marketers, enabling them to tailor their content and marketing strategies effectively to different audience segments, thereby maximizing engagement and revenue potential.

# 7 Conclusion and Discussion

## 7.1 Conclusion

This study introduced CineSage Insights, an advanced data-driven visualization suite, designed to revolutionize decision-making in the film-making industry. By leveraging a comprehensive dataset including movie plot summaries, budget information, genres, cast details, and historical revenue data, we developed predictive models and interactive visual tools. The suite's core components — the predictive revenue visualization and user clustering analysis — provide filmmakers with unprecedented insights into potential earnings and audience preferences.

Our predictive revenue model, integrating BERT embeddings and neural network architectures, demonstrated robust performance in forecasting movie revenues across different countries. The addition of an interactive choropleth map further empowered users to visualize and adjust revenue predictions based on varying inputs like budget and genres. This real-time adaptability is crucial for strategic planning in film production and distribution.

In audience analysis, our application of K-Means and Gaussian Mixture Model clustering to user preference data revealed distinct audience segments. The t-SNE visualizations provided a clear, intuitive representation of these clusters, highlighting their unique genre preferences. This clustering not only uncovers the prevailing tastes in movie genres but also aids in crafting targeted marketing strategies.

**ALL TEAM MEMBERS HAVE CONTRIBUTED A SIMILAR AMOUNT OF EFFORT.**

## 7.2 Discussion

The integration of advanced machine learning techniques and data visualization tools in CineSage Insights represents a significant stride in the application of data science to the film industry. However, there are several areas for future exploration and improvement:

- **Data Expansion**: Incorporating additional data sources, such as social media sentiment analysis and advanced demographic information, could further refine the predictive accuracy and audience insights.
- **Model Enhancement**: Continuous refinement of the predictive models, especially in handling sparse or missing data, could improve accuracy. Experimentation with different architectures and feature sets may yield better results.
- **User Experience**: Enhancing the interactivity and user-friendliness of the visualization tools will make the insights more accessible to a broader range of industry stakeholders.
- **Ethical Considerations**: As with any data-driven approach, ethical considerations around privacy and data usage must be rigorously addressed, particularly when expanding data sources.

In conclusion, CineSage Insights stands as a testament to the power of combining data analytics with intuitive visualization in making informed decisions in the film industry. While there are areas for enhancement, the current suite provides valuable tools for filmmakers and marketers, offering a glimpse into the future of data-driven decision-making in cinema.

## 7.3 Future Work

Looking forward, the potential of CineSage Insights can be expanded through:

- **Integration with Real-Time Data Sources**: Linking the suite with real-time box office data feeds could enable dynamic updating of predictions and insights.
- **Customizable User Interfaces**: Developing interfaces tailored to specific roles within the film industry could make the tool more relevant and useful for various stakeholders.
- **Machine Learning Advances**: Leveraging emerging machine learning techniques and algorithms could offer even more refined predictions and insights.
- **Global Market Analysis**: Expanding the model to encompass a wider array of international markets could provide a more global perspective on film performance.

In essence, CineSage Insights has the potential to not only guide current filmmaking and marketing decisions but also to shape the future of how the film industry leverages data for strategic advantage.

# References

[1] Box Office Mojo [n. d.]. *Box Office Mojo*. Box Office Mojo. https://www.boxofficemojo.com/

[2] Nehal Mohamed Ali, Marwa Mostafa Abd El Hamid, and Aliaa Youssif. 2019. Sentiment analysis for movies reviews dataset using deep learning models. *International Journal of Data Mining & Knowledge Management Process (IJDKP) Vol* 9 (2019). https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3403985

[3] Krushikanth R. Apala, Merin Jose, Supreme Motnam, C.-C. Chan, Kathy J. Liszka, and Federico de Gregorio. 2013. Prediction of Movies Box Office Performance Using Social Media. In *Proceedings of the 2013 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining* (Niagara, Ontario, Canada) *(ASONAM '13)*. Association for Computing Machinery, New York, NY, USA, 1209–1214. https://doi.org/10.1145/2492517.2500232

[4] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473* (2014). https://arxiv.org/abs/1409.0473

[5] Max Bain, Arsha Nagrani, Andrew Brown, and Andrew Zisserman. 2020. Condensed movies: Story based retrieval with contextual embeddings. In *Proceedings of the Asian Conference on Computer Vision*. https://openaccess.thecvf.com/content/ACCV2020/html/Bain_Condensed_Movies_Story_Based_Retrieval_with_Contextual_Embeddings_ACCV_2020_paper.html

[6] Rounak Banik. 2017. *The Movies Dataset*. https://www.kaggle.com/datasets/rounakbanik/the-movies-dataset

[7] Shuvamoy Chatterjee, Kushal Chakrabarti, Avishek Garain, Friedhelm Schwenker, and Ram Sarkar. 2021. JUMRv1: a sentiment analysis dataset for movie recommendation. *Applied Sciences* 11, 20 (2021), 9381. https://www.mdpi.com/2076-3417/11/20/9381

[8] Claire-Hélène Demarty, Cédric Penet, Mohammad Soleymani, and Guillaume Gravier. 2015. VSD, a public dataset for the detection of violent scenes in movies: design, annotation, analysis and evaluation. *Multimedia Tools and Applications* 74 (2015), 7379–7404. https://link.springer.com/article/10.1007/s11042-014-1984-4

[9] Maunendra Sankar Desarkar, Sudeshna Sarkar, and Pabitra Mitra. 2010. Aggregating Preference Graphs for Collaborative Rating Prediction. In *Proceedings of the Fourth ACM Conference on Recommender Systems* (Barcelona, Spain) *(RecSys '10)*. Association for Computing Machinery, New York, NY, USA, 21–28. https://doi.org/10.1145/1864708.1864716

[10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. 4171–4186. https://arxiv.org/abs/1810.04805

[11] Hyoji Ha, Hyunwoo Han, Seongmin Mun, Sungyun Bae, Jihye Lee, and Kyungwon Lee. 2019. An Improved Study of Multilevel Semantic Network Visualization for Analyzing Sentiment Word of Movie Review Data. *Applied Sciences* 9, 12 (2019). https://doi.org/10.3390/app9122419

[12] Bruce W. Herr, Weimao Ke, Elisha Hardy, and Katy Borner. 2007. Movies and Actors: Mapping the Internet Movie Database. In *2007 11th International Conference Information Visualization (IV '07)*. 465–469. https://doi.org/10.1109/IV.2007.78

[13] Jeffrey Lund and Yiu-Kai Ng. 2018. Movie Recommendations Using the Deep Learning Approach. In *2018 IEEE International Conference on Information Reuse and Integration (IRI)*. 47–54. https://doi.org/10.1109/IRI.2018.00015

[14] Asha S Manek, P Deepa Shenoy, and M Chandra Mohan. 2017. Aspect term extraction for sentiment analysis in large movie reviews using Gini Index feature selection method and SVM classifier. *World wide web* 20 (2017), 135–154. https://link.springer.com/article/10.1007/S11280-015-0381-X

[15] Márton Mestyán, Taha Yasseri, and János Kertész. 2013. Early prediction of movie box office success based on Wikipedia activity big data. *PloS one* 8, 8 (2013), e71226. https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0071226

[16] Vahan Petrosyan. 2014. Visualizing and Forecasting Box-Office Revenues: A Case Study of the James Bond Movie Series. (2014). https://digitalcommons.usu.edu/cgi/viewcontent.cgi?article=1413&context=gradreports

[17] Yunian Ru, Bo Li, Jianbo Liu, and Jianping Chai. 2018. An effective daily box office prediction model based on deep neural networks. *Cognitive Systems Research* 52 (2018), 182–191.

[18] Khaled Saad, Mai El-Ghandour, Ahmed Raafat, Reem Ahmed, and Eslam Amer. 2022. A Markov Model-Based Approach for Predicting Violence Scenes from Movies. In *2022 2nd International Mobile, Intelligent, and Ubiquitous Computing Conference (MIUCC)*. 21–26. https://doi.org/10.1109/MIUCC55081.2022.9781703

[19] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108* (2019).

[20] Katrien Verbert, Hendrik Drachsler, Nikos Manouselis, Martin Wolpers, Riina Vuorikari, and Erik Duval. 2011. Dataset-Driven Research for Improving Recommender Systems for Learning. In *Proceedings of the 1st International Conference on Learning Analytics and Knowledge* (Banff, Alberta, Canada) *(LAK '11)*. Association for Computing Machinery, New York, NY, USA, 44–53. https://doi.org/10.1145/2090116.2090122

[21] Aditya Vikram. 2023. Navigating Neural Networks: Exploring State-of-the-Art Activation Functions. Medium article.

[22] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Perric Cistac, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. Transformers: State-of-the-Art Natural Language Processing. Association for Computational Linguistics, 38–45. https://www.aclweb.org/anthology/2020.emnlp-demos.6