

Analyzing the Penetration of Southern Heroes in Northern India Using Real - Time Data Pipelines

Kishor Yadav Kommanaboina

Independent researcher, The Ohio State University Alumni

Email: [kkishore.iith\[at\]gmail.com](mailto:kkishore.iith[at]gmail.com)

Abstract: *The increasing popularity of Pan - India films have facilitated the cross - regional appeal of leading actors from South India into Northern parts of the country. This study aims to analyze the penetration and level of acceptance shown towards these actors in real - time by leveraging a comprehensive data integration pipeline. Understanding the dynamics between different regions is crucial for production studios and celebrity management teams to tailor their strategies accordingly. A scalable data integration system was architected to assimilate and refine inputs from diverse sources on a real - time basis including social networks, box office proceeds, streaming platform viewership metrics, search trends, traditional media coverage, and audience surveys. Key performance indicators such as the Regional Popularity Index (RPI), Cultural Acceptance Indicator (CAI), and Engagement Index (EI) were computed by applying sentiment analysis, feature extraction, and temporal modeling techniques on the consolidated data. The findings generate valuable insights regarding varying regional preferences and cultural acceptance levels, underscoring the importance of fusing heterogeneous data for a holistic perspective on market dynamics. The real - time data integration pipeline enhances decision - making capability. Future work will focus on refining the system architecture for improved accuracy and expanding the ambit of analysis to include additional cultural factors. This research contributes to a deeper understanding of cross - regional relationships within the Indian film fraternity and presents a scalable approach for real - time exploration of audience behavior and evolving market trends.*

Keywords: Pan - India Movies, Southern Heroes, Northern India, Real - Time Data Pipeline, Audience Sentiment Analysis, Box Office Performance, Streaming Viewership, Social Media Metrics, Cultural Acceptance, Regional Popularity Index (RPI), Cultural Acceptance Indicator (CAI), Engagement Index (EI), Cross - Regional Dynamics, Movie Marketing Strategies, Film Industry Analysis, Audience Behavior, Sentiment Analysis, Feature Engineering, Temporal Analysis, Data Integration

1. Introduction

The rise of Pan - India movies has significantly transformed the Indian film industry, promoting the cross - regional appeal of actors and films across the country. Southern heroes, traditionally popular in their native regions, are now gaining traction in Northern markets, driven by the success of these Pan - India productions. Understanding the penetration and acceptance of Southern heroes in Northern India is crucial for movie production houses and the backend teams of heroes to tailor their marketing and strategic efforts effectively.

The existing literature on audience engagement and market dynamics in the film industry has primarily focused on individual aspects such as social media sentiment, box office performance, and streaming viewership. For instance, M. et al. provided insights into cinema audience segmentation based on empirical findings from Indonesia, highlighting the importance of understanding diverse audience preferences [1]. Similarly, Chuan et al. emphasized the potential of social media analysis for forecasting movie box offices, proposing an improved Bass model for better accuracy [2].

Moreover, Sinha et al. explored the role of culture in international auxiliary channels, illustrating how cultural signals influence movie success across different regions [3]. In the Indian context, Agarwal et al. developed a dataset specifically for regional movies, aimed at enhancing recommender systems for better audience targeting [4]. Additionally, S. et al. examined the transition from single - screen theaters to online platforms like YouTube, demonstrating how regional blockbusters adapt and thrive in the digital age [5].

Despite these valuable contributions, there remains a notable gap in the literature concerning the integration of real - time data from diverse sources to analyze cross - regional dynamics, particularly for Southern heroes in Northern India. Existing models often lack the capability to incorporate real - time processing and comprehensive data integration, limiting their effectiveness in capturing the nuanced cultural and regional preferences of audiences. By addressing these gaps, our research aims to develop a robust framework that leverages a real - time data pipeline to compute key metrics such as the Regional Popularity Index (RPI), Cultural Acceptance Indicator (CAI), and Engagement Index (EI).

This study provides a comprehensive analysis of the penetration of Southern heroes in Northern India, offering valuable insights into regional preferences and cultural acceptance. The proposed framework not only enhances predictive accuracy but also supports timely and informed decision - making for strategic marketing and content creation in the film industry.

Problem Statement

The increasing success of Southern Indian films in Northern markets, propelled by the rise of pan - Indian blockbusters, demands intricate scrutiny of viewership tendencies and social acknowledgment across regions. Conventional analyses of industry movements commonly fall deficient, as they tend to focus on individual facets like social networking emotions or box office earnings. Such approaches lack the ability to combine real - time information from diverse resources, limiting their effectiveness in capturing the

nanced and evolving preferences of Northern audiences toward Southern stars.

Solution

To resolve this gap, we put forth an encompassing real - time data integration pipeline that assimilates inputs from online media, ticket sales reports, streaming viewership records, search patterns, traditional press coverage, and audience surveys. By leveraging state - of - the - art sentiment examination, characteristic extraction, and temporal modeling, our framework computes key performance metrics including the Regional Popularity Index (RPI), Cultural Acceptance Indicator (CAI), and Engagement Index (EI). This integrated approach not only improves predictive precision but also offers applicable insights, allowing film production houses and celebrity administration to tailor their strategies appropriately. The scalable design of the data pipeline ensures ongoing observation and customization, thereby supporting informed decision - making and fostering deeper comprehension of cross - regional viewership dynamics.

2. Methodology

Data Collection and Integration

To analyze the penetration of Southern heroes in Northern India, we designed a comprehensive data pipeline that integrates multiple data sources into a unified storage system. This pipeline ensures continuous collection, transformation, and analysis of data from various platforms, providing real - time insights into audience behavior and market dynamics.

Data Sources

1) Social Media Metrics:

- a) **Platforms:** Twitter (via Twitter4j), Facebook (via Pages API), Instagram (via Instagram Graph API), Reddit (via PRAW), YouTube (via YouTube Data API), TikTok (via TikTok API for Business)
- b) **Data Collected:**
 - **Twitter:** Tweets, retweets, likes, mentions, hashtags related to Southern heroes.
 - **Facebook:** Posts, comments, likes on official pages.
 - **Instagram:** Posts, comments, likes, follower counts.
 - **Reddit:** Posts, comments, upvotes, downvotes from relevant subreddits.
 - **YouTube:** Video views, likes, comments, and channel statistics.
 - **TikTok:** Video views, likes, comments, and follower counts.
- c) **Frequency:** Real - time ingestion for Twitter, Facebook, Instagram, YouTube, and TikTok; periodic batch ingestion for Reddit due to rate limits.
- d) **Processing:** Data is ingested using appropriate APIs and tools, and streamed into the data lake (S3/HDFS/GCS) using Apache Kafka and batch processing jobs.

2) Box Office Data:

- a) **Sources:** Regional movie theaters, box office reporting services such as Box Office India, Bollywood Hungama, and Koimoi.
- b) **Data Collected:** Regional box office collections, number of screenings, ticket sales, and revenue.

- c) **Frequency:** Collected daily but processed weekly for trend analysis.
- d) **Processing:** Data is ingested into the data lake (S3/HDFS/GCS) using Apache Sqoop, ensuring secure and reliable transfer from source databases.

3) Streaming Service Data:

- a) **Sources:** FlixPatrol, Parrot Analytics, SimilarWeb.
- b) **Data Collected:** Viewership statistics, regional watch patterns, user ratings.
- c) **Frequency:** Updated daily.
- d) **Processing:** Data is scraped from third - party websites and stored in the data lake (S3/HDFS/GCS). Apache Spark processes the data to derive viewership trends and user engagement metrics.

4) Search Trends:

- a) **Sources:** Google Trends, Bing Trends.
- b) **Data Collected:** Search volumes for Southern heroes and related keywords in different regions.
- c) **Frequency:** Updated daily.
- d) **Processing:** Data is collected via APIs and stored in the data lake (S3/HDFS/GCS). Apache Airflow orchestrates the workflow for regular data extraction and loading into the analysis environment.

5) Traditional Media Coverage:

- a) **Sources:** Newspapers, magazines, TV channels.
- b) **Data Collected:** Articles, reviews, interviews, mentions of Southern heroes.
- c) **Frequency:** Collected daily.
- d) **Processing:** Text data is extracted using web scraping tools like Scrapy and stored in the data lake (S3/HDFS/GCS).
- e) **Audience Surveys:**
- f) **Methods:** Online surveys, face - to - face interviews, and ads on YouTube, Facebook, Instagram.
- g) **Data Collected:** Audience preferences, perceptions, cultural attitudes, regional sentiment towards Southern heroes.
- h) **Frequency:** Collected quarterly.
- i) **Processing:** Survey data is stored in the data lake (S3/HDFS/GCS). Apache Spark processes this data for statistical analysis and integration with other data sources.

6) Forum Data:

- a) **Platforms:** Online forums and discussion boards (via web scraping).
- b) **Data Collected:** Posts, comments, mentions of Southern heroes.
- c) **Frequency:** Periodic batch ingestion.
- d) **Processing:** Data is ingested using web scraping tools like BeautifulSoup, Scrapy, and Selenium, and stored in the data lake (S3/HDFS/GCS).

Data Storage and Processing

All data collected from various sources is stored in a unified data lake, which can be Amazon S3, Google Cloud Storage (GCS), or Hadoop Distributed File System (HDFS), depending on the specific needs and compatibility of the tools used for analysis.

a) Data Cleaning:

- **Tasks:** Removing duplicates, handling missing values, standardizing formats.
- **Tools:** Apache Spark, Python (Pandas).
- **Details:** Standardization includes normalizing numerical data, converting text data to lowercase, and ensuring consistent date formats. Apache Spark handles large - scale data cleaning operations to prepare datasets for analysis.

b) Data Transformation:

- **Tasks:** Normalizing data, extracting features.
- **Tools:** Apache Spark, Python (Numpy, Scikit - learn).
- **Details:** Feature extraction involves identifying key metrics such as sentiment scores from text data, viewership patterns from streaming data, and search trends from search engines. Apache Spark's MLlib is used for feature extraction and transformation tasks.

Sentiment Analysis**a) Lexicon - Based Sentiment Analysis:**

- **Tools:** VADER, AFINN.
- **Tasks:** Sentiment scoring of social media posts and comments.
- **Details:** VADER and AFINN are utilized to perform initial sentiment scoring, which is then aggregated to gauge overall sentiment trends.

b) Advanced NLP Analysis (BERT):

- **Tools:** Hugging Face Transformers.
- **Tasks:** Context - aware sentiment analysis.
- **Details:** BERT models are fine - tuned on domain - specific data to capture nuanced sentiments from social media and traditional media text. This analysis provides deeper insights into audience perceptions.

Feature Engineering**a) Combining Features:**

- **Tasks:** Merge data from various sources, create composite indicators.
- **Tools:** Python (Pandas, Numpy).
- **Details:** Features from social media sentiment, box office data, streaming viewership, and surveys are combined to create metrics like the Regional Popularity Index (RPI), Cultural Acceptance Indicator (CAI), and Engagement Index (EI).

b) Temporal Features:

- **Tasks:** Generate lagged variables, rolling averages.
- **Tools:** Apache Spark, Python (Pandas).
- **Details:** Temporal features help capture trends and seasonality in the data, enabling more accurate predictive modeling.

c) Interdependencies with VAR:

- **Tasks:** Analyze relationships between variables.
- **Tools:** Statsmodels in Python.
- **Details:** Vector Autoregression (VAR) models are used to study the interdependencies between various features, such as the impact of social media sentiment on box office performance.

Computing Composite Indices

The following indices are computed by combining various features and using weights calculated through methods like Expert Judgment, Analytic Hierarchy Process (AHP), Principal Component Analysis (PCA), and Regression Analysis to determine the relative importance of each component.

a) Regional Popularity Index (RPI):

- **Components:** Social Media Sentiment Scores, Regional Box Office Performance.
- **Formula:** $RPI = (W1 \times \text{Sentiment Score}) + (W2 \times \text{Box Office Performance})$

b) Cultural Acceptance Indicator (CAI):

- **Components:** Survey Positive Response Rate, Search Volume Index.
- **Formula:** $CAI = (W1 \times \text{Survey Positive Response Rate}) + (W2 \times \text{Search Volume Index})$

c) Engagement Index (EI):

- **Components:** Streaming Views, Media Mentions.
- **Formula:** $EI = (W1 \times \text{Streaming Views}) + (W2 \times \text{Media Mentions})$

Validation and Continuous Iteration**a) Validation:**

- **Tasks:** Apply indices to a test dataset, compare results with known outcomes, adjust weights if necessary.
- **Details:** Ensure the accuracy and reliability of the computed indices by validating against historical data and expert evaluations.

b) Continuous Iteration:

- **Tasks:** Regularly update data, retrain models, monitor performance.
- **Tools:** Apache Kafka, Apache Airflow, TensorFlow, Keras.
- **Details:** Implement a feedback loop to continuously refine and improve the models, ensuring they adapt to new data and changing market dynamics.

Visualization Layer

To enable the hero backend team and production houses to monitor performance and take strategic actions based on the data, a visualization layer is implemented. This layer provides an intuitive interface for real - time data analysis and decision - making.

Visualization Layer

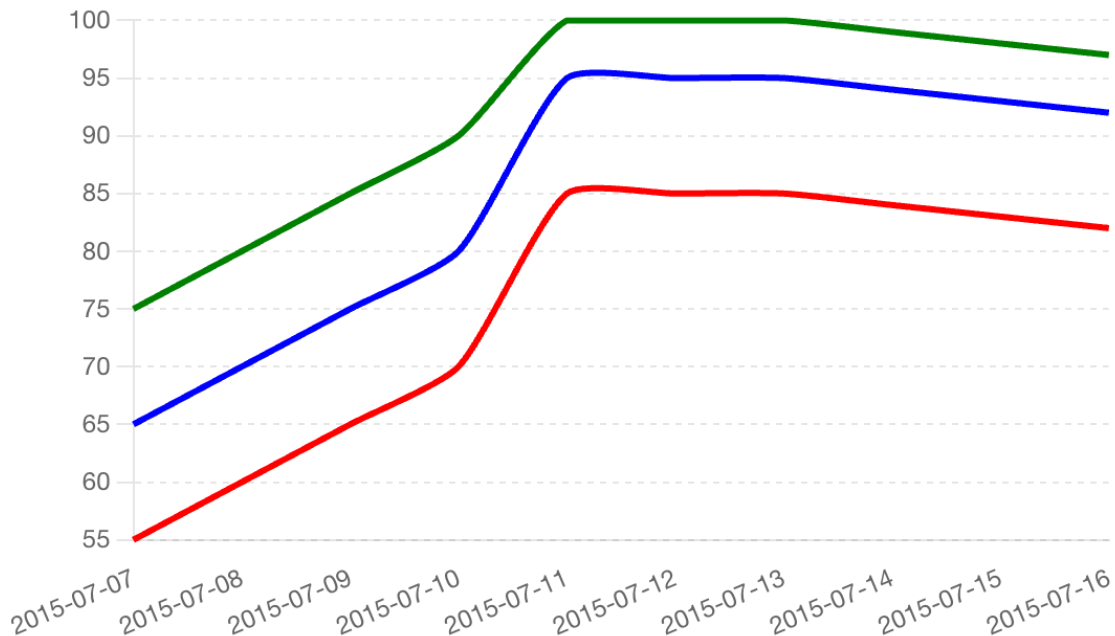
- 1) **Tools:** Tableau, Power BI, Grafana, or custom dashboards using D3.js and Flask.
- 2) **Capabilities:**
 - **Real - Time Monitoring:** Dashboards displaying key metrics like Regional Popularity Index (RPI), Cultural Acceptance Indicator (CAI), Engagement Index (EI), social media sentiment, box office performance, and streaming viewership.
 - **Trend Analysis:** Visualizations showing trends over time for various metrics, helping to identify patterns and shifts in audience behavior.

- **Geographical Insights:** Maps highlighting regional variations in audience sentiment and engagement, providing a clear view of market penetration across different areas.
- **Interactive Filters:** Ability to filter data by time period, region, and other dimensions for deeper analysis and targeted insights.
- **Alerts and Notifications:** Set up alerts for significant changes in metrics, with notifications via email or SMS for timely responses to trends and anomalies.

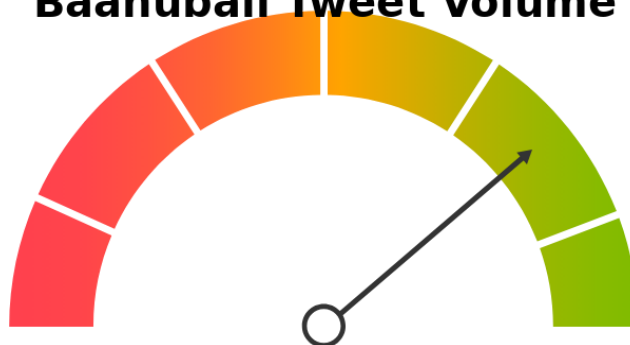
3) **Example Dashboard Components:**

a) **Overview Dashboard:**

- **Metrics:** Displays key metrics such as RPI, CAI, EI, overall sentiment score, box office revenue, and streaming viewership.
- **Visualizations:** Includes line charts for trend analysis, bar charts for comparative analysis, and pie charts for distribution.



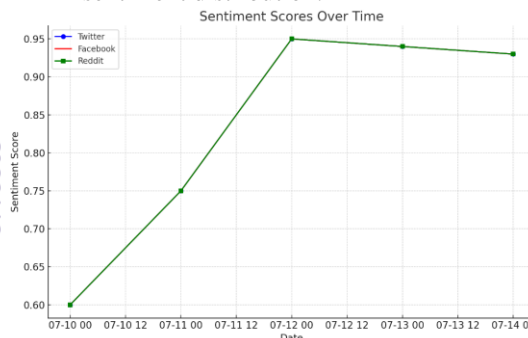
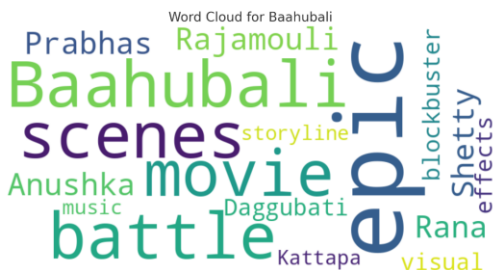
Baahubali Tweet Volume



Sentiment Analysis Dashboard:

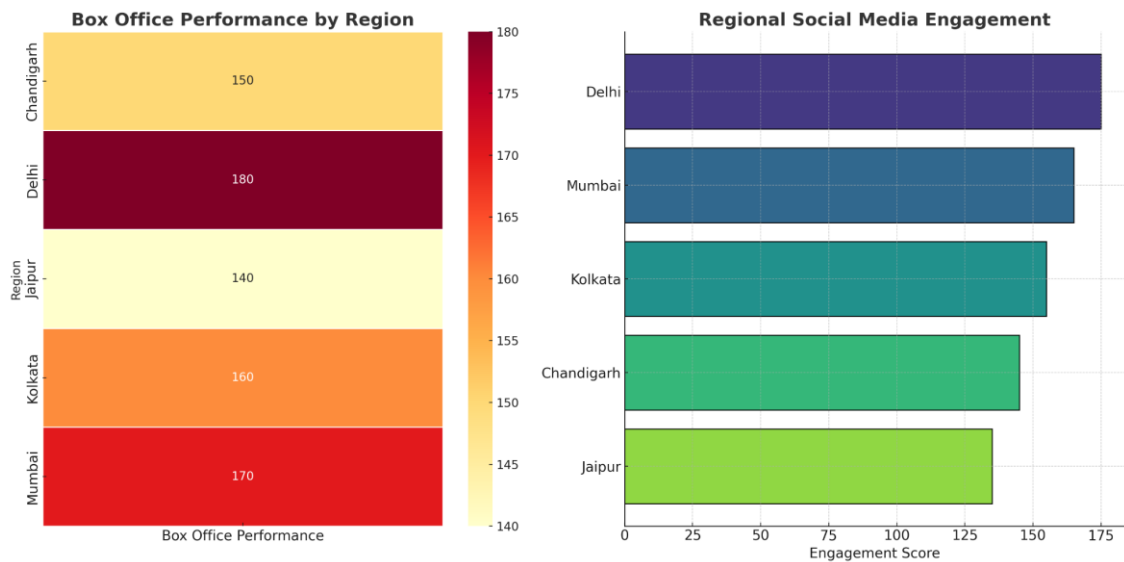
- **Metrics:** Tracks sentiment scores from social media, traditional media, and forums.

- **Visualizations:** Word clouds for common phrases, line charts for sentiment trends over time, and bar charts for sentiment distribution.



Regional Performance Dashboard:

- **Metrics:** Shows regional variations in box office performance, social media engagement, and streaming viewership.
- **Visualizations:** Heatmaps for geographical insights, bar charts for regional comparisons, and line charts for regional trends.

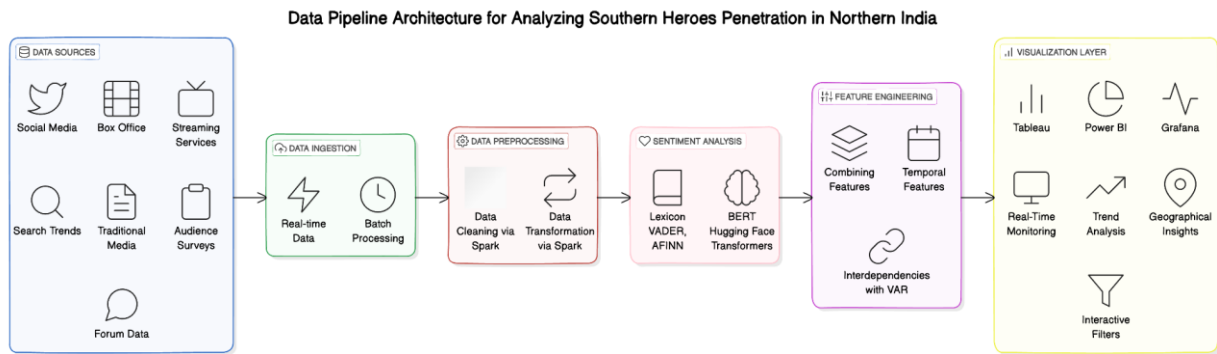


Engagement and Reach Dashboard:

- **Metrics:** Monitors engagement metrics from social media platforms, streaming services, and traditional media.
- **Visualizations:** Scatter plots for engagement vs. reach, line charts for engagement trends, and bar charts for platform - specific engagement.

By implementing this visualization layer, the hero backend team and production houses can monitor performance in real - time, analyze trends, and make informed decisions to optimize their marketing strategies and enhance audience engagement.

3. Conclusion



4. Summary of Data Flow

The data pipeline architecture for analyzing the penetration of Southern heroes in Northern India integrates multiple data sources into a cohesive system, facilitating comprehensive and real - time analysis. The pipeline begins with data collection from diverse sources, including social media, box office reports, streaming services, search trends, traditional media, audience surveys, and forum data. This data is ingested through real - time and batch processing methods into a unified data lake. Apache Spark is employed for data cleaning and transformation, ensuring the data is standardized and ready for analysis. Sentiment analysis is conducted using lexicon - based methods and advanced NLP models like BERT, providing nuanced insights into audience sentiment. Feature engineering combines various data points and incorporates temporal features and interdependencies using Vector Autoregression (VAR). Finally, the processed data is visualized through platforms like Tableau, Power BI, and Grafana, offering real - time monitoring, trend analysis, geographical insights, and interactive filters.

5. Interpretation and Implication

This robust data pipeline architecture allows for a detailed and dynamic analysis of how Southern heroes are received in Northern India. By leveraging real - time data from multiple sources, it provides a comprehensive view of audience engagement, sentiment, and regional preferences. The use of advanced sentiment analysis and feature engineering techniques enhances the accuracy and depth of insights, enabling stakeholders to make informed decisions. For movie production houses and the backend teams of heroes, this pipeline offers a valuable tool for optimizing marketing strategies, tailoring content to regional preferences, and tracking the effectiveness of promotional efforts. The ability to monitor trends and audience sentiment in real time also allows for agile responses to emerging trends and potential issues.

6. Limitations and Future Scope

While the data pipeline architecture provides extensive capabilities, it also has certain limitations. The accuracy of the analysis heavily relies on the quality and completeness of the

data collected. Incomplete or biased data can lead to skewed results. Additionally, the reliance on third - party APIs and scraping methods may introduce data latency and compliance issues.

Future enhancements could include the integration of additional data sources such as localized social media platforms, more sophisticated machine learning models for predictive analytics, and enhanced real - time processing capabilities. Expanding the geographical scope beyond Northern India to include other regions could provide a more holistic view of audience dynamics across the country. Moreover, incorporating feedback loops to continuously refine and improve the models based on new data and insights would further enhance the pipeline's effectiveness.

In conclusion, the proposed data pipeline architecture represents a significant advancement in understanding and analyzing the impact of Southern heroes in Northern India. By addressing current limitations and exploring future enhancements, this pipeline can continue to provide valuable insights and support data - driven decision - making in the movie industry.

References

- [1] M., Mujiya Ulkhaq, Finsaria Fidiyanti, Adyatama Arga, Zakia A. Maulani, Adi Nugroho. "Segmentation of Cinema Audiences: An Empirical Finding from Indonesia, " Proceedings of the 2nd International Conference on Computer Science, Information Technology, and Electrical Engineering (ICOMITEE), 2019, pp.3 - 8, doi: 10.1145/3354153.3354154.
- [2] Zhang, C., Tian, Y. - X., & Fan, Z. - P. "Forecasting the box offices of movies coming soon using social media analysis: A method based on improved Bass models, " Expert Systems With Applications, 2021, doi: 10.1016/J.ESWA.2021.116241.
- [3] Sinha, A., Gu, H., Kim, N., & Emile, R. "Signaling effects and the role of culture: movies in international auxiliary channels, " European Journal of Marketing, 2019, doi: 10.1108/EJM - 09 - 2017 - 0587.
- [4] Agarwal, P., Verma, R., & Majumdar, A. "Indian Regional Movie Dataset for Recommender Systems, " arXiv: Information Retrieval, 2018.
- [5] S., V., Srinivas., V., H., C., V., Megha, Shyam., Raghav, Nanduri., Vasundhara, Singhal., R, Vishnu, Dath. "From Single Screen to YouTube: Tracking the Regional Blockbuster, " Bioscope: South Asian Screen Studies, 2018, doi: 10.1177/0974927619841205.