

Movie Popularity and Target Audience Prediction Using the Content-Based Recommender System

Mr.Mabubasha¹, Thorlikonda Siva Prasad², Vishnumolakala Sarvan Manikanta³,
Sabavath Mohan Naik⁴

¹Assistant Professor, Dept. of Computer Science and Engineering, R.V.R. & J.C College of Engineering, Chowdavaram, Guntur, Andhra Pradesh, India

^{2,3,4}B. Tech Student, Dept. of Computer Science and Engineering, R.V.R. & J.C College of Engineering, Chowdavaram, Guntur, Andhra Pradesh, India

Abstract: The movie is one of the integral components of our everyday entertainment. The worldwide movie industry is one of the most growing and significant industries and seizing the attention of people of all ages. It has been observed in a recent study that only a few movies achieve success. Uncertainty in the sector has created immense pressure on the film production stakeholder. Moviemakers and researchers continuously feel it necessary to have some expert systems predicting the movie success probability preceding its production with reasonable accuracy. A maximum of the research work has been conducted to predict the movie's popularity in the post-production stage. To help the movie maker estimate the upcoming film and make necessary changes, we need to conduct the prediction at the early stage of movie production and provide specific observations about the upcoming movie.

Keywords: Movie, audience, prediction, moviemakers, film production.

I. INTRODUCTION

The worldwide movie industry is a fast-moving, revenue-generating industry, and billions of dollars are involved in it. A large number of people are associated with this industry, and massive investments—both qualitative and quantitative—are required. In 2019, the total box office revenue of the United States and Canada was \$11.32 billion. However, in reality, only a few movies achieve significant success. Film producers and researchers consistently feel the need for expert systems that can predict a movie's chance of success with appropriate accuracy during its production. The movie industry is massive and diverse. A significant number of parameters from various dimensions are involved in creating a movie. Representing an upcoming movie's success or its degree of success is a highly complex task. Research has been conducted to predict movie popularity. Earlier works have focused on post-production or post-release forecasting. However, such predictions are not beneficial, as investors have already committed their funds to the movie's production. Predictions made during the early production or pre-production stage with satisfying accuracy are more beneficial in securing investment. Forecasts made soon after the cast, director, and storyline have been finalized can help investors make informed decisions. After thorough study, it has been found that significant research has been conducted on predicting movie success before the official release. Predictions performed shortly before or after the official release may utilize additional data and produce more precise outcomes. However, they are too late for investors to make any critical decisions. Early-stage (production phase) forecasting of movie success is the most valuable.

Very little work has been done to forecast movie success at this stage, and the accuracy of existing models is not significantly high. Most of the previous works focus only on determining the success probability of upcoming movies. Some classify the problem as binary (hit/flop), while others use a multiclass approach. Movie makers often target specific audience groups when creating a new movie. Audience age is one of the essential criteria. Some movies are designed for younger audiences, some for teenagers, others for middle-aged or senior audiences, and some aim to appeal to all age groups. If we could predict whether an upcoming movie would be popular among the target audience groups, it would benefit the movie makers. They could also measure the influence of the movie on different age groups during the early production stage and make necessary changes if needed. Forecasting a movie's success and predicting its impact on target audiences during early production are interrelated and meaningful. The outcome of such work could reduce the risks involved in the movie industry.

II. LITERATURE REVIEW

[1] L. Sharma and A. Gera, "A survey of recommendation system: Research challenges," *Int. J. Eng. Trends Technol.*, vol. 4, no. 5, pp. 1989_1992, 2013.

A recommender system is a set of tools for information retrieval. It improves access and proactively recommends items and services that match users' tastes by considering their explicit and implicit preferences and behaviors. Recommender systems have become very popular in the e-commerce field. Today, the internet is flooded with diverse information that makes it very difficult for the end-users to reach out for what they need. Recommender systems provide tailored views to users who are constantly adapted to the users' changing tastes. Although many recommendation techniques have been developed in multiple domains, recommender systems still face problems and challenges that hinder their precision. This paper provides a comprehensive summary of the key challenges and problems when developing recommender systems and summarizes the latest research achievements and directions to resolve them. In addition to that, we go beyond this by presenting the evaluation techniques used to judge the performance of recommender systems.

[2] N. Das, S. Borra, N. Dey, and S. Borah, "Social networking in web based movie recommendation system," in *Social Networks Science: Design, Implementation, Security, and Challenges*. Cham, Switzerland: Springer, 2018, pp. 25_45.

Movie Recommendations Systems are a common practice by most of the online stores today. The web based movie recommendation systems makes predictions about the responses of the users based on their search history or known preferences. Recommendation of items is usually done based on the properties or content of the item or collaboration of the user's ratings, and by using intelligent algorithms that include classification or clustering techniques. Accurate prediction of what the customer may likely to buy or the user my visit is of utmost important, as it benefits both the service providers and customers. This chapter provides the evolution, fundamental concepts, classification, traditional and novel models, requirements, similarity measures, evaluation approaches, issues, challenges, impacts due to social networking, and future of movie recommendation systems.

[3] P. Nagarnaik and A. Thomas, "Survey on recommendation system methods," in *Proc. 2nd Int. Conf. Electron. Commun. Syst. (ICECS)*, Feb. 2015, pp. 1603_1608.

In recent years recommendation systems have changed the way of communication between both websites and users. Recommendation system sorts through massive amounts of data to identify interest of users and makes the information search easier. For that purpose many methods have been used. Collaborative Filtering (CF) is a method of making automatic predictions about the interests of customers by collecting information from number of other customers, for that purpose many collaborative base algorithms are used. CHARM algorithm is one of the frequent patterns finding algorithm which is capable to handle huge dataset, unlike all previous association mining algorithms which do not support huge dataset. This paper covers different techniques which are used in recommendation system and also proposes a new system for efficient web page recommendation based on hybrid collaborative filtering i.e. using collaborative technique and CHARM algorithm which are coupled with the pattern discovery algorithms such as clustering and association rule mining.

[4] M. A. Hameed, O. Al Jadaan, and S. Ramachandram, "Collaborative Filtering based recommendation system: A survey," *Int. J. Comput. Sci. Eng.*, vol. 4, no. 5, p. 859, 2012.

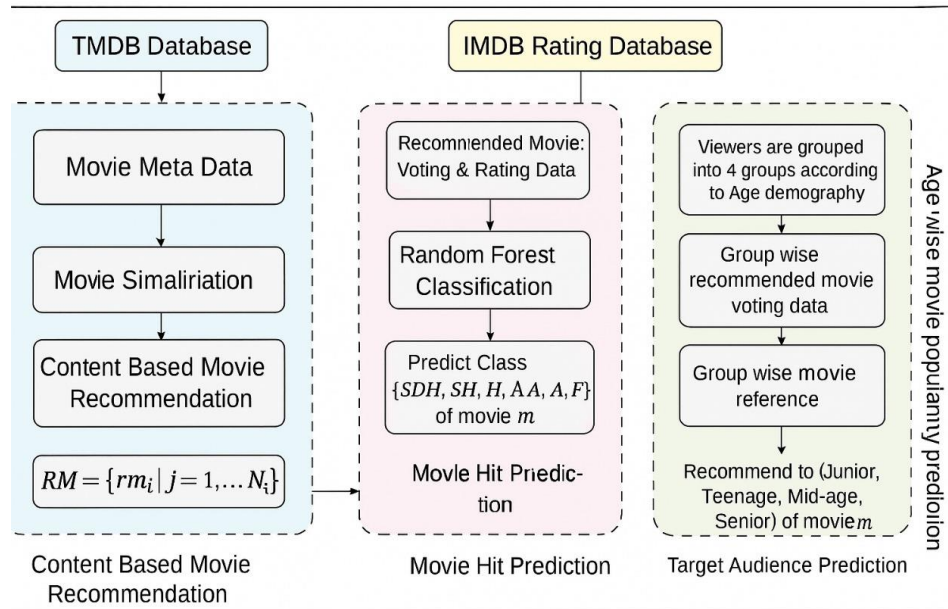
In the era of big data, recommender system (RS) has become an effective information filtering tool that alleviates information overload for Web users. Collaborative filtering (CF), as one of the most successful recommendation techniques, has been widely studied by various research institutions and industries and has been applied in practice. CF makes recommendations for the current active user using lots of users' historical rating information without analyzing the content of the information resource. However, in recent years, data sparsity and high dimensionality brought by big data have negatively affected the efficiency of the traditional CF-based recommendation approaches. In CF, the context information, such as time information and trust relationships among the friends, is introduced into RS to construct a training model to further improve the recommendation accuracy and user's satisfaction, and therefore, a variety of hybrid CF-based recommendation algorithms have emerged. In this paper, we mainly review and summarize the traditional CF-based approaches and techniques used in RS and study some recent hybrid CF-based recommendation approaches and techniques, including the latest hybrid memory-based and model-based CF recommendation algorithms. Finally, we discuss the potential impact that may improve the RS and future direction. In this paper, we aim at introducing the recent hybrid CF-based recommendation techniques fusing social networks to solve data sparsity and high dimensionality and provide a novel point of view to improve the performance of RS, thereby presenting a useful resource in the state-of-the-art research result for future researchers.

[5] B. Sarwar, G. Karypis, J. Konstan, and J. Reidl, "Item-based collaborative Filtering recommendation algorithms," in *Proc. 10th Int. Conf. WorldWide Web (WWW)*, 2001, pp. 285_295.

In order to overcome the limitations of data sparsity and inaccurate similarity in personalized recommendation systems, a new collaborative filtering recommendation algorithm by using items categories similarity and interestingness measure is proposed. In this algorithm, first the items categories similarity matrix is constructed by calculating the item-item category distance, and then analyzes the correlation degree of different items by using interestingness measure, last an improved collaborative filtering algorithm is proposed by combining the information of items categories with item-item interestingness and utilizing improved conditional probability method as the standard item-item similarity measure. Experimental results show this algorithm

can effectively alleviate the dataset sparsity problem and achieve better prediction accuracy compared to other well-performing collaborative filtering algorithms.

III. SYSTEM ARCHITECTURE



Algorithms:

- TF-IDF
- KNN
- RF

TF-IDF Vectorizer

TF-IDF (Term Frequency–Inverse Document Frequency) is primarily used within the Content-Based Recommendation module to compute the similarity between movies based on their textual metadata — such as genres, overview (plot), keywords, director names, cast, and more. TF-IDF transforms these textual attributes into meaningful numerical vectors that reflect how important a word is to a particular movie description in relation to the entire dataset. For example, common terms like “love” or “fight” that appear in many movies will have lower weights, while more unique and distinguishing words like “time-travel” or “heist” will carry higher weights.

Once each movie’s metadata is converted into TF-IDF vectors, cosine similarity is used to compare the vectors and identify which movies are most alike. This similarity scoring helps form the initial list of content-based similar movies for a given input movie. These similar movies are then passed on to other modules (Random Forest for popularity prediction and MLP Classifier for audience prediction). Thus, TF-IDF helps ensure that the recommendations are more accurate and contextually relevant by focusing on the distinguishing features of each movie’s metadata.

K- NEAREST NEIGHBOUR ALGORITHM (KNN):

The K-Nearest Neighbour (KNN) algorithm plays a crucial role in refining the results of the Content-Based Recommendation System. After identifying a broad set of similar movies based

on metadata such as genre, director, and cast, KNN is applied to select the top-K most similar movies from this set. By calculating the distance (or similarity score) between the input movie and others in the dataset, KNN ensures that only the most relevant and closely matched movies are chosen. This filtering step enhances the accuracy of the subsequent modules by providing more meaningful data. The selected top-K movies are then used to extract voting and rating patterns for the Random Forest Classifier to predict the movie's popularity, and for the MLP Classifier to determine age-wise audience preferences. Overall, KNN helps eliminate noise, improves efficiency, and ensures that only high-quality, relevant data flows through the system, ultimately boosting prediction performance.

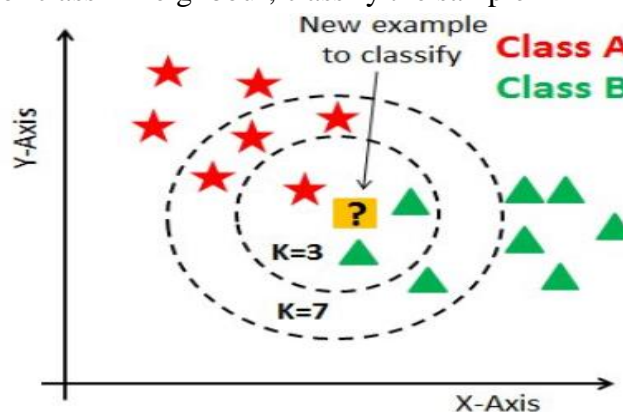
$$\begin{aligned} \text{Euclidian Distance} &= D(x, y) \\ &= (x_i - y_i)_{2k_i} = 1 \end{aligned} \quad (1)$$

K=number of cluster

x, y=co-ordinate sample spaces

The algorithm for KNN is defined in the steps given below:

1. D represents the samples used in the training and k denotes the number of nearest neighbour.
2. Create super class for each sample class.
3. Compute Euclidian distance for every training sample
4. Based on majority of class in neighbour, classify the sample



K- Nearest Neighbour

Random Forest

In our experiment, we use random forest as a classifier. The popularity of decision tree models in data mining is owed to their simplification in algorithm and flexibility in handling different data attribute types. However, single-tree model is possibly sensitive to specific training data and easy to overfit. Ensemble methods can solve these problems by combine a group of individual decisions in some way and are more accurate than single classifiers. Random forest, one of ensemble methods, is a combination of multiple tree predictors such that each tree depends on a random independent dataset and all trees in the forest are of the same distribution. The capacity of random forest not only depends on the strength of individual tree but also the correlation between different trees. The stronger the strength of single tree and the less the correlation of different trees, the better the performance of random forest. The variation of trees comes from their randomness which involves bootstrapped samples and randomly selects a subset of data attributes.

Below is the step by step Python implementation. ...

Step 1: Import and print the dataset.

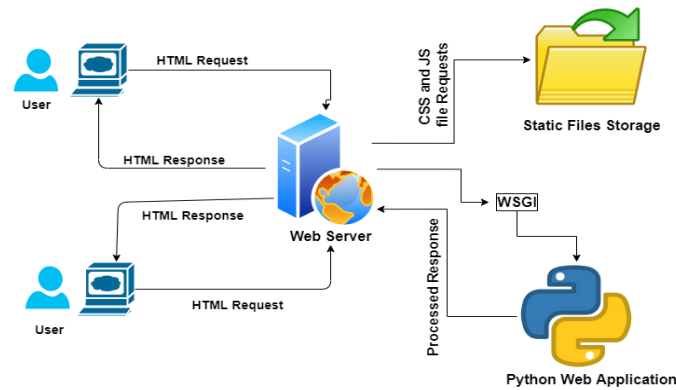
Step 2: Select all rows and column 1 from dataset to x and all rows and column 2 as y.

Step 3: Fit Random Forest regressor to the dataset.

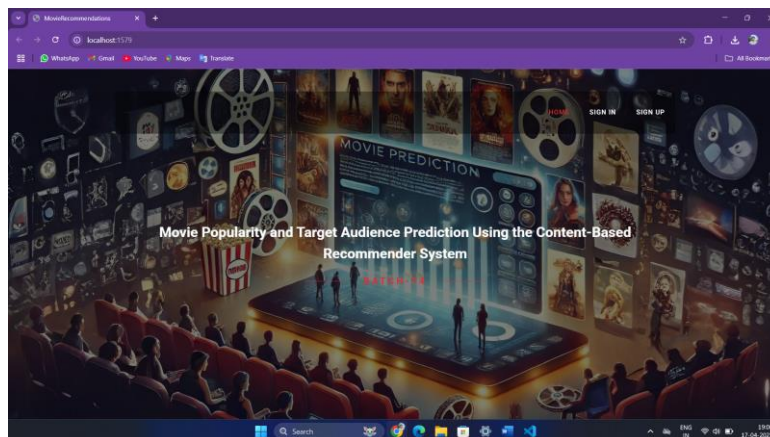
Step 4: Predicting a new result.

Step 5: Visualising the result.

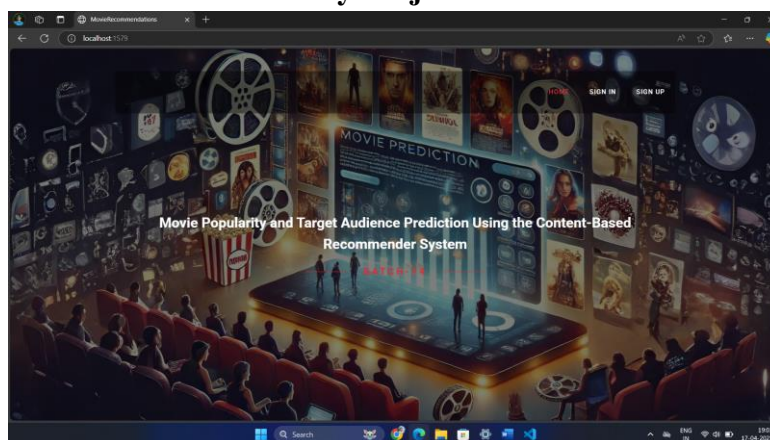
Deployment Model



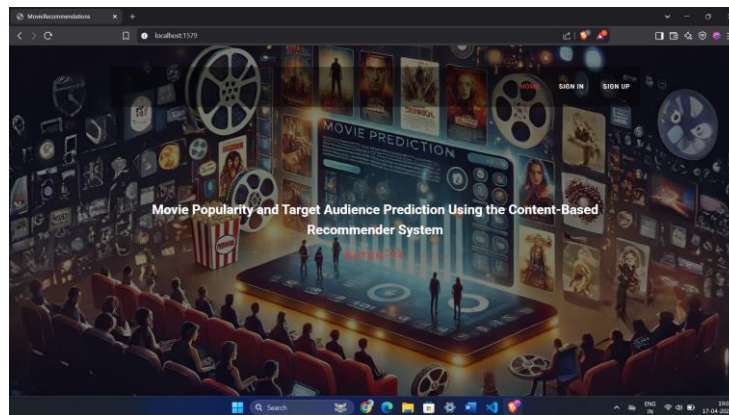
IV. RESULTS



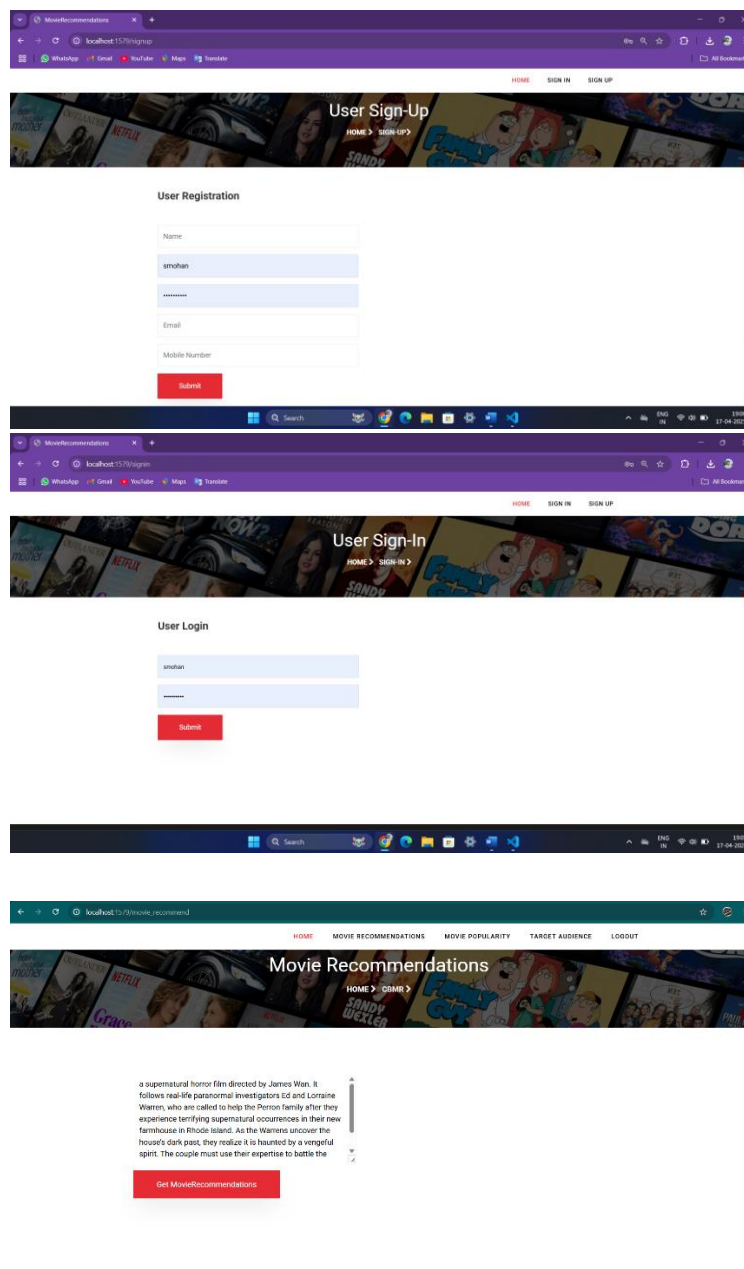
Result of my Project in chrome



Result of my Project in Microsoft edge



Result of my Project in Brave SCREENSHOTS



Content Based Movie Recommendations				
Movie name	Genres	Casts	Director	
The Conjuring	Horror,Thriller	Patrick Wilson,Vera Farmiga,Lili Taylor	James Wan	
Sinister	Horror,Thriller,Mystery	Ethan Hawke,Juliet Rylance,Vincent D'Onofrio	Scott Derrickson	
The Conjuring 2	Horror	Patrick Wilson,Vera Farmiga,Madison Wolfe	James Wan	
The Apparition	Horror,Thriller	Ashley Greene,Sebastian Stan,Tom Felton	Todd Lincoln	
Out of the Dark	Thriller,Horror	Julia Stiles,Scott Speedman,Stephen Rea	Luis奎lez	
An American Haunting	Horror,Thriller	Donald Sutherland,Sissy Spacek,James D'Arcy	Courtney Solomon	
Insidious: Chapter 2	Horror,Thriller	Patrick Wilson,Rose Byrne,Ty Simpkins	James Wan	
Diary of the Dead	Horror,Action,Science Fiction	Michelle Morgan,Joshua Close,Shawn Roberts	George A. Romero	

Movie Popularity				
Movie name	Genres	Casts	Director	Movie Popularity
The Conjuring	Horror,Thriller	Patrick Wilson,Vera Farmiga,Lili Taylor	James Wan	H
Sinister	Horror,Thriller,Mystery	Ethan Hawke,Juliet Rylance,Vincent D'Onofrio	Scott Derrickson	AA
The Conjuring 2	Horror	Patrick Wilson,Vera Farmiga,Madison Wolfe	James Wan	H
The Apparition	Horror,Thriller	Ashley Greene,Sebastian Stan,Tom Felton	Todd Lincoln	F
Out of the Dark	Thriller,Horror	Julia Stiles,Scott Speedman,Stephen Rea	Luis奎lez	F
An American Haunting	Horror,Thriller	Donald Sutherland,Sissy Spacek,James D'Arcy	Courtney Solomon	A
Insidious: Chapter 2	Horror,Thriller	Patrick Wilson,Rose Byrne,Ty Simpkins	James Wan	AA
Diary of the Dead	Horror,Action,Science Fiction	Michelle Morgan,Joshua Close,Shawn Roberts	George A. Romero	A

Target Audience				
Movie name	Genres	Casts	Director	Target Audience
The Conjuring	Horror,Thriller	Patrick Wilson,Vera Farmiga,Lili Taylor	James Wan	Teenage
Sinister	Horror,Thriller,Mystery	Ethan Hawke,Juliet Rylance,Vincent D'Onofrio	Scott Derrickson	Teenage
The Conjuring 2	Horror	Patrick Wilson,Vera Farmiga,Madison Wolfe	James Wan	Senior
The Apparition	Horror,Thriller	Ashley Greene,Sebastian Stan,Tom Felton	Todd Lincoln	Senior
Out of the Dark	Thriller,Horror	Julia Stiles,Scott Speedman,Stephen Rea	Luis奎lez	Teenage
An American Haunting	Horror,Thriller	Donald Sutherland,Sissy Spacek,James D'Arcy	Courtney Solomon	Senior
Insidious: Chapter 2	Horror,Thriller	Patrick Wilson,Rose Byrne,Ty Simpkins	James Wan	Senior
Diary of the Dead	Horror,Action,Science Fiction	Michelle Morgan,Joshua Close,Shawn Roberts	George A. Romero	Mid-age

V. CONCLUSION

A substantial amount of financing is consumed in every box-office movie. However, most movies fail to achieve success. Earlier, the most significant number of works have been done on post-production or post-release forecast. The estimate does not influence as the investor has already consumed their funds on the _lm production. The pre-production or early production stage forecast needs high accuracy and the best time to ensure investment. The objective of our study is to propose an expert system that could help the movie maker execute necessary changes if needed at the appropriate time. Our system can food cost the level of popularity of the upcoming movie before the production has started for the earliest stage of the production and with significant accuracy.

REFERENCES

- [1] L. Sharma and A. Gera, "A survey of recommendation system: Research challenges," *Int. J. Eng. Trends Technol.*, vol. 4, no. 5, pp. 1989_1992, 2013.
- [2] N. Das, S. Borra, N. Dey, and S. Borah, "Social networking in web based movie recommendation system," in *Social Networks Science: Design, Implementation, Security, and Challenges*. Cham, Switzerland: Springer, 2018, pp. 25_45.
- [3] P. Nagarnaik and A. Thomas, "Survey on recommendation system methods," in *Proc. 2nd Int. Conf. Electron. Commun. Syst. (ICECS)*, Feb. 2015, pp. 1603_1608.
- [4] M. A. Hameed, O. Al Jadaan, and S. Ramachandram, "Collaborative filtering based recommendation system: A survey," *Int. J. Comput. Sci. Eng.*, vol. 4, no. 5, p. 859, 2012.
- [5] B. Sarwar, G. Karypis, J. Konstan, and J. Reidl, "Item-based collaborative filtering recommendation algorithms," in *Proc. 10th Int. Conf. World Wide Web (WWW)*, 2001, pp. 285_295.
- [6] Y. Koren and R. Bell, "Advances in collaborative filtering," in *Recommender Systems Handbook*. Boston, MA, USA: Springer, 2015, pp. 77_118.
- [7] J. L. Herlocker, J. A. Konstan, and J. Riedl, "Explaining collaborative filtering recommendations," in *Proc. ACM Conf. Comput. supported Coop-erat. Work (CSCW)*, 2000, pp. 241_250.
- [8] M. J. Pazzani, "A framework for collaborative, content-based and demographic filtering," *Artif. Intell. Rev.*, vol. 13, no. 5, pp. 393_408, 1999.
- [9] R. Van Meteren and M. Van Someren, "Using content-based filtering for recommendation," in *Proc. Mach. Learn. Inf. Age, MLnet/ECML2000 Workshop*, vol. 30, 2000, pp. 47_56.
- [10] P. B. Thorat, R. M. Goudar, and S. Barve, "Survey on collaborative filtering, content-based filtering and hybrid recommendation system," *Int. J. Comput. Appl.*, vol. 110, no. 4, pp. 31_36, Jan. 2015.