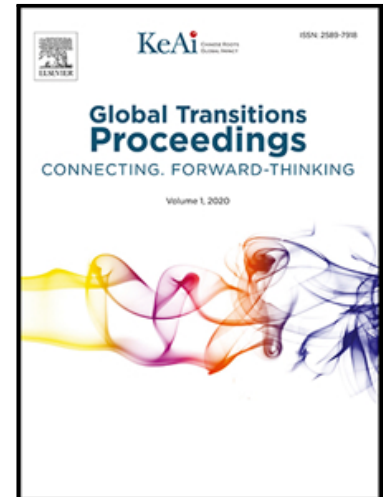


Movie Recommendation and Sentiment Analysis Using Machine Learning

N Pavitha , Vithika Pungliya , Ankur Raut , Roshita Bhonsle ,
Atharva Purohit , Aayushi Patel , R Shashidhar

PII: S2666-285X(22)00017-6
DOI: <https://doi.org/10.1016/j.gltp.2022.03.012>
Reference: GLTP 112



To appear in: *Global Transitions Proceedings*

Please cite this article as: N Pavitha , Vithika Pungliya , Ankur Raut , Roshita Bhonsle ,
Atharva Purohit , Aayushi Patel , R Shashidhar , Movie Recommendation and Sentiment
Analysis Using Machine Learning, *Global Transitions Proceedings* (2022), doi:
<https://doi.org/10.1016/j.gltp.2022.03.012>

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2022 The Authors. Publishing Services by Elsevier B.V. on behalf of KeAi Communications Co. Ltd.

This is an open access article under the CC BY-NC-ND license
(<http://creativecommons.org/licenses/by-nc-nd/4.0/>)

Pavitha N^a, Vithika Pungliya^{a,*} vithika.pungliya20@vit.edu, Ankur Raut^a, Roshita Bhonsle^a, Atharva Purohit^a, Aayushi Patel^a, Shashidhar R^b

^aDepartment of Artificial Intelligence and Data Science, Vishwakarma Institute of Technology, Pune (411037), Maharashtra, India.

^bDepartment of Electronics and Communication, Sri Jayachamarajendra College of Engineering JSS Science & Technology University, Mysuru, Karnataka, India.

*Corresponding author.

Abstract

In the modern world, where technology is at the forefront of every industry, there has been an overload of information and data. Thus, a recommendation system comes in handy to deal with this large volume of data and filter out the useful information which is fast and relevant to the user's choice. This paper describes an approach to a movie recommendation system using Cosine Similarity to recommend similar movies based on the one chosen by the user. Although the existing recommendation systems get the job done, it does not justify if the movie is worth spending time on. To enhance the user experience, this system performs sentiment analysis on the reviews of the movie chosen using machine learning. Two of the supervised machine learning algorithms Naïve Bayes (NB) Classifier and Support Vector Machine (SVM) Classifier are used to increase the accuracy and efficiency. This paper also gives a comparison between NB and SVM on the basis of parameters like Accuracy, Precision, Recall and F1 Score. The accuracy score of SVM came out to be 98.63% whereas accuracy score of NB is 97.33%. Thus, SVM outweighs NB and proves to be a better fit for Sentiment Analysis.

Keywords

Cosine similarity; Movie recommendation; Naïve Bayes; Sentiment analysis; Support Vector Machine.

1. Introduction

Since its invention, the Internet has grown rapidly and continues to grow each day. The abundance of information available online, has made it a strenuous task to access the right information quickly and easily [1]. Fortunately, this problem can be solved with the help of recommendation systems.

Recommendation systems are used extensively today and have found applications in multiple industries such as e-commerce, retail, banking, entertainment etc. [2] These systems collect and auto-analyse the user data to generate personalised recommendations for the users [3,4]. The most common approaches to implement recommendation systems are Content-based Filtering (CBF), Collaborative Filtering (CF) and Hybrid Filtering [5]. CBF is an approach that is used to analyse the content of each item and recommend other items that have similar characteristics [6]. CF addresses some of the limitations of CBF and provides recommendations by comparing the similarities between the users and the items. It uses the knowledge of the user's previous preferences as well as the preferences of other similar users to generate a recommendation. Many recommendation systems are also known to use the Hybrid-filtering technique combining the features of both CBF and CF methods [7,8].

A movie's popularity is based on the type of reviews it gets from the audience [9]. These reviews are also responsible for affecting the choice of other users. Users are more likely to choose a movie that was preferred by most people rather than a movie that was largely disliked [10]. Analysing these reviews, ignoring the reviews that contain misleading information also adds to the difficulty of decision-making [11]. Sentiment Analysis provides a solution to this problem.

Sentiment Analysis facilitates a way to use NLP (natural language processing) to extract information from a textual source and classify the statement or word or document as

positive or negative [12]. It is very useful to understand the opinion of the author and indicate the user experience [13].

Opinion mining uses the concepts of data mining to extract and classify the opinions expressed in various online forums or platforms. This enables better understanding of the user's sentiment or feeling towards a particular subject matter [14].

The paper presents a system that not only recommends movies to the users but also analyses and classifies the reviews into positive or negative. The movie recommendation part is performed using Cosine Similarity and a comparison is drawn between SVM and the NB algorithm to perform the Sentiment Analysis of the reviews.

The objective of the study is to deal with the large volume of data and filter useful information, recommend similar movies based on user's choice and perform Sentimental Analysis on the reviews of the movie chosen.

The paper follows the given structure; Section II covers the Literature Review. Section III, discusses the Methodology which includes the dataset, the pre-processing of data, mining of data for movie recommendation, machine learning for sentiment analysis and finally, the performance report. Section IV comprises the results and discussions. Finally, the conclusion can be found in Section V.

2. Related Works

In this section, the various existing methods and the drawbacks of the existing work are discussed in detail.

In this paper, the authors propose a hybrid approach that combines a content-based approach with genre correlation to implement a recommendation system. This system takes into

genres while making recommendations to the user [15].

In this paper, a new similarity algorithm is introduced. This is called User Profile Correlation-based Similarity (UPCSim) and it allows other user behavioural data to influence the accuracy of the recommendation. It calculates the weights of similarity which depend on the user's rating and the user's behaviour value and classifies the preferences of the user using K-nearest Neighbours algorithm. While this approach shows a decrease in the Mean Absolute Error (1.64%) and the Root Mean Square Error (1.4%), it requires more computation time [16].

The authors discuss the implementation of a movie recommendation system using two algorithms, Cosine Similarity and K-Nearest Neighbours. The movie recommendation is done using the cosine similarity algorithm. A normalised popular score is used to obtain the function for computing distance and the K-Nearest Neighbours algorithm is applied to enhance the accuracy [17].

In this paper, a hybrid recommendation system is proposed that uses sentiment analysis of user tweets for movies to obtain a sentiment score to improve the recommendation made to the users using a weighted fusion score method [18].

This paper implements five machine learning classifiers – Multinomial Naïve Bayes, SVM, Decision Tree, Bernoulli Naïve Bayes, Maximum Entropy are applied on the pre-processed data containing feature vectors to classify the movie reviews data [19].

3. Methodology

Under this section, the methods used for the execution of the study and implementation of the algorithms have been discussed. The diagram below shows the flowchart of the methodology.

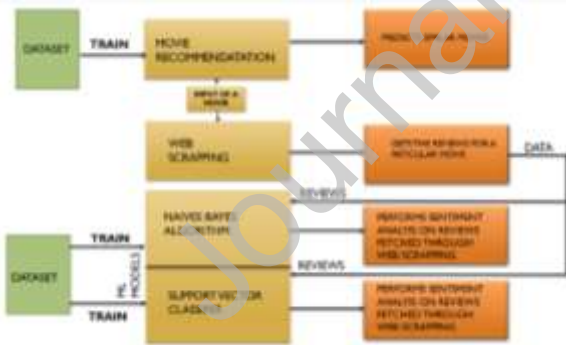


Fig 1. Flowchart of proposed Method

Fig 1 explains the methodology that has been used in the project. The study has used a dataset to train the cosine similarity model which is used for recommending movies. Then using another dataset, Naïve Bayes (NB) and Support Vector Machine (SVM) Classifier for Sentiment Analysis has been trained. Now a movie name is taken as an input and sent to the movie recommendation model to predict similar movies. Through web scraping from the IMDB site, the reviews of that movie are obtained and sent to the Sentiment Analysis model for classifying the reviews as positive or negative.

A) Dataset

Three datasets have been used for study. 2 of them are for Movie Recommendation and 1 is for Sentiment Analysis. The

'tbmd_5000_credits.csv' [20] and the one used for sentiment analysis is 'reviews.txt'. The 2 datasets used in movie recommendation are then merged to form a single data set and the columns kept under it are 'movie_id', 'title' and 'tags' [21].

The reviews data set has only 2 columns, one for the 'reviews' and other for the 'comments'. The positive comments have been labelled as 1 and the negative ones have been labelled as 0. There are 3943 positive comments and 2975 negative comments [23].

B) Data Pre- Processing

After merging the 2 datasets into a single dataset, only the essential columns such as 'movie_id', 'title', 'overview', 'genres', 'keywords', 'cast' and 'crew' are kept, rest are removed from the dataset [24].

Then using Abstract Syntax trees, the columns of 'genres', 'keywords', 'cast', 'crew' have been refined.

Furthermore, these columns have been combined under 'tags'. Then using the count vectorizer, the column 'tags' is tokenised. Tokenising means to divide the sentences into words. Here the pre-processing for Movie Recommendation comes to an end.

The pre-processing for Sentiment Analysis requires Natural Language Tool Kit (NLTK). The NLTK is a leading platform for building Python programs to work with human language data.

This is a standard python library used for natural language processing and computational linguistics. Using this library, the stop words are downloaded. Stop Words are the usually used words in any language. Examples of such words include 'a', 'an', 'the', 'if', 'or'. They are used in Text Mining and Natural Language Processing (NLP) to eliminate such words as they carry very little useful information. Then using the TfidfVectorizer, the column of Comments is tokenised. Now the Data Pre-processing of the datasets ends here.

	movie_id	title	tags
0	1995	Avatar	In the 22nd century, a paraplegic Marine is di...
1	285	Pirates of the Caribbean: At World's End	Captain Barbossa, long believed to be dead, ha...
2	206647	Spectre	A cryptic message from Bond's past sends him o...
3	49026	The Dark Knight Rises	Following the death of District Attorney Harve...
4	48525	John Carter	John Carter is a war-weary, former military ca...

Fig 2. Final dataset being used for Movie Recommendation.

The Fig 2 shows the dataset used for Movie Recommendation, which has the 'movie_id', 'title' and 'tags'.

	Reviews	Comments
0	1	The Da Vinci Code book is just awesome.
1	1	this was the first clive cussler i've ever rea...
2	1	i liked the Da Vinci Code a lot.
3	1	i liked the Da Vinci Code a lot.
4	1	I liked the Da Vinci Code but it ultimately did...

Fig 3. Dataset for Sentiment Analysis.

The Fig 3 shows the dataset for sentiment analysis, has 2 columns, one for the 'comments' and other for the 'reviews'.

C) Data Mining for Movie Recommendation

(Proposed)				
Bernoulli's Naive Bayes (Existing) [11]	0.875	0.884	0.8633	0.8735
Multinomial NB (Existing) [11]	0.885	0.9294	0.8333	0.8787
SVM (Existing) [11]	0.8733	0.859	0.8933	0.8753
NB (Existing)[5]	0.8183	0.84	0.79	0.82
SVM (Existing)[5]	0.8745	0.87	0.88	0.88
Random Forest (Existing)[5]	0.9601	0.93	1.00	0.96
Stacked-LSTM (Existing)[14]	0.9365	0.94	0.94	-
Minimal-RNN(Existing) [14]	0.8564	0.86	0.86	-
CNN(Existing) [15]	0.8915	0.8259	0.8246	0.8253
LSTM (Existing)[15]	0.9550	0.9087	0.8228	0.8636

Table 1. Comparison of Models

The Table 1 shows the comparison between different Accuracy, Precision, Recall and AUC scores of the 2 proposed models and 3 existing models (Bernoulli's Naïve Bayes, Multinomial NB, SVM) studied in [11], 3 existing models (SVM, NB, Random Forest) studied in [5] 2 existing models (Stacked-LSTM, Minimal-RNN) studied in [14] and 2 existing models (CNN, LSTM) studied in [15]. Proposed SVM model is better than NB in all parameters.

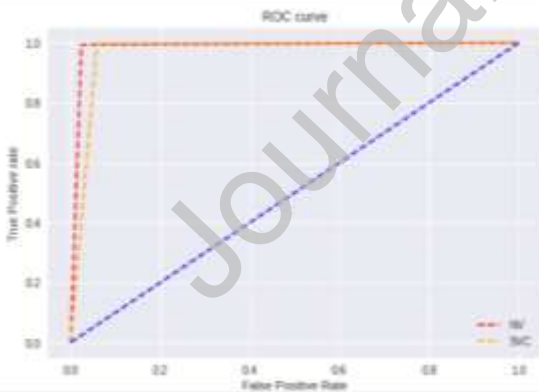


Fig 7. ROC Curves.

The Fig 7 shows the ROC Curves between the 2 algorithms.

One of the examples that taken for the study is the movie Spectre.

Through the Cosine Similarity algorithm, predictions of 5 other movies - Quantum of Solace, Never Say Never Again, Skyfall, Thunderball, From Russia with Love were made.

```
print("Enter a movie")
movie=input()
recommend(movie)
```

Enter a movie
Spectre
Quantum of Solace
Never Say Never Again
Skyfall
Thunderball
From Russia with Love

Fig 8. Prediction of Movies using the Cosine Similarity Algorithm.

Fig 8, shows the reviews about the movie - Spectre, entered by the user. Web scraping is used to get the taglines of the reviews. Web scraping is performed from the IMDB website and perform Sentiment Analysis on it using the NB and SVC algorithm.

```
The second installment of the Craig Bond [8]
A Bond film to the classic world [1]
could not get great [1]
Quality craftsmanship, basic McGary's credit missing [8]
Bond vs. surveillance [1]
Interesting and moving [1]
Solid and suspense movie [1] I could have used more character and space to make it more than 'a Bond movie' [8]
I am your best friend of staying alive [1]
Craig's list ... [1]
Just go and see this movie. Make your own mind up. [8]
I can really put you through it, James [1]
Quality craftsmanship, McGary's credit missing [8]
The James Bond franchise should have ended decades ago [8]
Spectre [1]
All the thrills one hopes for in a Bond film [1]
Enjoyable installment in Bond series with lots of noisy action, thrills, emotion and spectacular scenes [1]
It's time to shake up the franchise and add some FBI [8]
Too long at the tail, not [8]
At times Spectre is a lot of fun, other times its just... over-elaborate, over the top [1]
Spectre is perhaps the best of the Daniel Craig Bond movies yet [1]
Oh well [1]
Just follow the movie logic [8]
What else there is there to do, mate [8]
Oh, what a solid Bond film this is [1]
Now [1]
```

Fig 9. Sentiment Analysis on Reviews using NB Algorithm.

```
The second installment of the Craig Bond [8]
A Bond film to the classic world [1]
could not get great [1]
Quality craftsmanship, basic McGary's credit missing [8]
Bond vs. surveillance [1]
Interesting and moving [1]
Solid and suspense movie [1] I could have used more character and space to make it more than 'a Bond movie' [8]
I am your best friend of staying alive [1]
Craig's list ... [1]
Just go and see this movie. Make your own mind up. [1]
I can really put you through it, James [1]
Quality craftsmanship, McGary's credit missing [8]
The James Bond franchise should have ended decades ago [8]
Spectre [1]
All the thrills one hopes for in a Bond film [1]
Enjoyable installment in Bond series with lots of noisy action, thrills, emotion and spectacular scenes [1]
It's time to shake up the franchise and add some FBI [8]
Too long at the tail, not [8]
At times Spectre is a lot of fun, other times its just... over-elaborate, over the top [1]
Spectre is perhaps the best of the Daniel Craig Bond movies yet [1]
Oh well [1]
Just follow the movie logic [8]
What else there is there to do, mate [8]
Oh, what a solid Bond film this is [1]
Now [1]
```

Fig 10. Sentiment Analysis on Reviews using SVC.

After performing Sentiment Analysis, [Fig 9, 10], one can tell whether it is a good movie or not.

So, the review 'Enjoyable installment in Bond series with lots of noisy action, thrills, emotion and spectacular scenes' gets assigned a value 1 which is the interpretation of a good review to the movie.

Another review 'the James Bond franchise should have ended decades ago' gets assigned a value 0 which is the interpretation of a bad review to the movie.

5. Conclusion

This paper is basically divided into two major parts. One of which focuses on Movie Recommendation system and the other on the Sentiment analysis. The study discusses both the systems in detail and has come to some important conclusions. For the Movie Recommendation System, the Cosine Similarity algorithm has been used to recommend the best movies that are related to the movie entered by the user based on different factors such as the genre of the movie, overview, the cast as well as the ratings given to the movie. Cosine Similarity has

been quite accurate at recommending the movies. Sentiment analysis also plays an important role in this study. It basically aims to classify the reviews into positive or negative. Two algorithms have been used for the same. One of which is NB and other is SVC. The main reason behind using two algorithms is to find out what which is the best algorithm to classify the reviews because the reviews have huge diversity in them, so it is very important to choose the right algorithm for classification. Finally, the experimental results show that SVM Algorithm has better accuracy than NB by a very small margin.

Some prospects of this study have been mentioned below:

1. Increasing the Accuracy of both Sentiment Analysis for better classification of sarcastic or ironic reviews.
2. Sentiment Analysis of the reviews in different languages other than English.
3. Movie recommendation according to users' preference (cast, genre, year of release, etc.).

Although the system is very accurate, it does have some limitations. One of which is, if the movie entered by the user isn't present in the dataset or if the user does not enter the name of the movie in the similar manner as that of in the dataset, then the system fails to recommend movies. One more limitation is the linguistic barrier while doing the sentimental analysis. As of now only reviews written in English can be analyzed. The Sentimental analysis also gives wrong classification if the reviews are sarcastic or ironic.

References

1. Nassar, N., Jafar, A., & Rahhal, Y. (2020). A novel deep multi-criteria collaborative filtering model for recommendation system. *Knowledge-Based Systems*, 187, 104811.
2. Subramani, P., & BD, P. (2021). Prediction of muscular paralysis disease based on hybrid feature extraction with machine learning technique for COVID-19 and post-COVID-19 patients. *Personal and ubiquitous computing*, 1-14.
3. Beheshti, A., Yakhchi, S., Mousaeirad, S., Ghafari, S. M., Goluguri, S. R., & Edrisi, M. A. (2020). Towards cognitive recommender systems. *Algorithms*, 13(8), 176.
4. Tran, D. N., Nguyen, T. N., Khanh, P. C. P., & Trana, D. T. (2021). An iot-based design using accelerometers in animal behavior recognition systems. *IEEE Sensors Journal*.
5. Sharma, S., Rana, V., & Malhotra, M. (2021). Automatic recommendation system based on hybrid filtering algorithm. *Education and Information Technologies*, 1-16.
6. Yu, K., Lin, L., Alazab, M., Tan, L., & Gu, B. (2020). Deep learning-based traffic safety solution for a mixture of autonomous and manual vehicles in a 5G-enabled intelligent transportation system. *IEEE transactions on intelligent transportation systems*, 22(7), 4337-4347.
7. Reddy, S. R. S., Nalluri, S., Kuniseti, S., Ashok, S., & Venkatesh, B. (2019). Content-based movie recommendation system using genre correlation. In *Smart Intelligent Computing and Applications* (pp. 391-397). Springer, Singapore.
8. Parameshchari, B. D. Big Data Analytics on Weather Data: Predictive Analysis Using Multi Node Cluster Architecture. *International Journal of Computer Applications*, 0975-8887.
9. M. Yassen and S. Tedmori, "Movies Reviews Sentiment Analysis and Classification," 2019 IEEE Jordan International Joint Conference on Electrical Engineering and Information Technology (JEEIT), 10.1109/JEEIT.2019.8717422.
10. Le, N. T., Wang, J. W., Le, D. H., Wang, C. C., & Nguyen, T. N. (2020). Fingerprint enhancement based on tensor of wavelet subbands for classification. *IEEE Access*, 8, 6602-6615.
11. Rajput, Neha, and S. Chauhan. "Analysis of various sentiment analysis techniques." *International Journal of Computer Science and Mobile Computing* 8.2 (2019): 75-79.
12. Yu, K., Tan, L., Lin, L., Cheng, X., Yi, Z., & Sato, T. (2021). Deep-learning-empowered breast cancer auxiliary diagnosis for 5GB remote E-health. *IEEE Wireless Communications*, 28(3), 54-61.
13. Shaukat, Z., Zulfiqar, A. A., Xiao, C., Azeem, M., & Mahmood, T. (2020). Sentiment analysis on IMDB using lexicon and neural networks. *SN Applied Sciences*, 2(2), 1-10.
14. Rachana, B., Priyanka, T., Sahana, K. N., Supriya, T. R., Parameshchari, B. D., & Sunitha, R. (2021). Detection of polycystic ovarian syndrome using follicle recognition technique. *Global Transitions Proceedings*, 2(2), 304-308.
15. Widiyaningtyas, T., Hidayah, I. & Adji, T.B, 2021. User profile correlation-based similarity (UPCSim) algorithm in movie recommendation system. *J Big Data* 8, 52.
16. Vu, D. L., Nguyen, T. K., Nguyen, T. V., Nguyen, T. N., Massacci, F., & Phung, P. H. (2020). HIT4Mal: Hybrid image transformation for malware classification. *Transactions on Emerging Telecommunications Technologies*, 31(11), e3789.
17. Ramni Harbir Singh, Sargam Maurya, Tanisha Tripathi, Tushar Narula, Gaurav Srivastav. "Movie Recommendation System using Cosine Similarity and KNN" *International Journal of Engineering and Advanced Technology (IJEAT)* ISSN: 2249 – 8958, Volume-9 Issue-5, June 2020.
18. Zhang, J., Yu, K., Wen, Z., Qi, X., & Paul, A. K. (2021). 3D reconstruction for motion blurred images using deep learning-based intelligent systems. *CMC-computers Materials & Continua*, 66(2), 2087-2104.
19. Kumar, S., De, K., & Roy, P. P. (2020). Movie recommendation system using sentiment analysis from microblogging data. *IEEE Transactions on Computational Social Systems*, 7(4), 915-923.
20. Rahman, A., & Hossen, M. S. (2019, September). Sentiment analysis on movie review data using machine learning approach. In *2019 International Conference on Bangla Speech and Language Processing (ICBSLP)* (pp. 1-4). IEEE.
21. Uddin, S., Khan, A., Hossain, M. E., & Moni, M. A. (2019). Comparing different supervised machine learning algorithms for disease prediction. *BMC medical informatics and decision making*, 19(1), 1-16.
22. Ghosh, S., Dasgupta, A., & Swetapadma, A. (2019, February). A study on support vector machine based linear and non-linear pattern classification. In *2019 International Conference on Intelligent Sustainable Systems (ICISS)* (pp. 24-28). IEEE.
23. Kia Dashtipour, Mandar Gogate, Ahsan Adeel, Hadi Larijani and Amir Hussain. (2021) Sentiment Analysis of Persian Movie Reviews Using Deep Learning. *Entropy*, 23(5), 596.
24. Soubraylu, S., & Rajalakshmi, R. (2021). Hybrid convolutional bidirectional recurrent neural network based sentiment analysis on movie reviews. *Computational Intelligence*, 37(2), 735-757.