

Personalized Recommendations of Products to Users

Spoorthi Chinivar

Masters in Data Science, FAU

Machine Learning and Data Analytics for Industry 4.0 Seminar

Erlangen, Germany

spoorthi.chinivar@fau.de

Abstract—Many organizations utilize recommendation systems to increase their profitability and win over their customers, including Facebook, which suggests friends, LinkedIn, which promotes employment, Spotify, which recommends music, Netflix, which recommends movies, and Amazon, which recommends purchases. When it comes to movie recommendation system, suggestions are made based on user similarities (collaborative filtering) or by considering a specific user's behavior (content-based filtering) that he or she wishes to interact with. Using TF-IDF, cosine similarity method for content-based filtering, and deep learning for a collaborative approach, this study compares two movie recommendation system.

The proposed systems are evaluated by calculating the precision and recall values. On a small dataset, a content-based filtering methodology had a precision of 5.6% whereas a collaborative approach had a precision of 57%. Collaborative filtering clearly worked better than content-based filtering. Future improvements involve creating a single hybrid recommendation system that combines a collaborative and content-based approach to improve the outcomes.

Keywords—movie recommendation system, content-based filtering, collaborative-based approach, deep learning, TF-IDF, cosine similarity.

I. INTRODUCTION

Personalized product recommendations are when a website provides a series of unique product recommendations based on each visitor's behavior and profile. Recommender systems are utilized in a wide range of contexts, including e-commerce and online advertising. To filter through millions of items of information and provide its customers with personalised suggestions, online businesses like YouTube, Netflix, Amazon, Pinterest, and a long list of others rely on recommender systems. Recommender systems have been researched and shown to provide both Internet firms and their customers significant advantages [12].

With this in mind, the goal is to formulate machine learning and deep learning strategies for movie recommender systems. There are two primary categories of recommender systems: collaborative filtering, which may be further split into memory-based and model-based techniques, and content-based filtering. Content-based filtering methods are based on item descriptions and user-preferred profiles. For simplicity, it suggests similar movies based on a specified movie (movie name would be an input) or based on all movies watched by the user (user is an input). It extracts features of items and can also look at the user's history to make recommendations. Collaborative filtering is a

technique for predicting a user's interests by integrating preferences from several other users. It suggests unwatched movies to a user from a collection of movies that have already been seen by users who share similar interests.

There are many ways to implement a recommendation engine. Some of the limitations of today's systems include the lack of data, cold start problems, and changes in user preferences. Some examples of highly efficient recommendation systems include engines used by multinational technology giants such as YouTube, Netflix, Amazon and Pandora. These systems process huge amounts of data at various stages, including the training stages [2]. The main aim of this work is to build an efficient movie recommendation system based on 'MovieLens' dataset using Content-based and Collaborative filtering approach and compare the results of both.

II. RELATED WORK

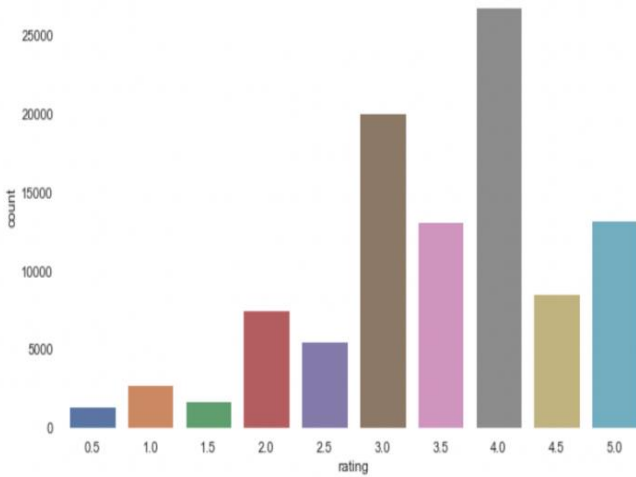
The most common way to perform content-based filtering is with a similarity metric. A similarity metric is a mathematical measure used to determine how similar a vector is to a given vector. They are computed from the item's feature vector and the user's preference feature vector from past records. Then some of the most similar items are recommended. Similarity metrics are mainly cosine similarity (cosine angle between vectors), dot product (product of cosine angle and magnitude of vectors), Euclidean distance (element-wise squared distance between two vectors), and Pearson similarity (ratio between the covariance and the standard deviation of both vectors). B.S. Oladapo [1] applied content-based technique in paper recommendation system. The author used Jaccard similarity coefficient also known as Jaccard index to compute similarity between users' query (users' attributes) and the attributes of the papers. The recommendations suggested by the system were sent via emails to the intended users. N. Muthurasu et. al. [2] created a movie Recommendation System using Term Frequency-Inverse Document Frequency (TF-IDF) and Cosine Similarity Method. Donghui Wang [3] proposed a recommender system on computer science publications referred to as the Publication Recommender System (PRS). This system is based on a new content-based filtering (CBF) recommendation model using chi-square and SoftMax regression, which are combined to construct a real-time online system.

The most general method of achieving collaborative filtering is to use a k-nearest-neighbor approach between users as proposed by W.Croft et. al. [4]. Clustering also helps in finding similar groups in user space. K-Means clustering is the most common clustering algorithm. An effective collaborative movie

recommendation system was proposed which uses a hybrid K-Means clustering with a Cuckoo search algorithm for optimization by Katarya R et. al. [5]. There is another popularly known technique called, matrix factorization for performing collaborative filtering as described by C.Volinsky et. al. [6]. Xiangnan He et. al. [7] presented a neural network architecture for modeling latent features of users and objects and developing a general framework for neural network-based collaborative filtering. Snehal R Chavare et. al. [8] described a recommendation system by translating the inner product of traditional matrix factorization into deep learning. A collaborative filtering deep learning technique based on an autoencoder accurately predicted movie ratings as presented by Jeffrey Lund et.al. [9].

III. DATASET DESCRIPTION

The proposed model uses the MovieLens dataset. The MovieLens dataset was collected by GroupLens research project at the University of Minnesota [13]. This dataset contains 3683 movie tags of 9742 movies with 100836 ratings ranging from 0.5 to 5 (see Graph 1). The dataset is preprocessed and contains only users who have rated 20 or more different movies. Data is contained in links.csv, movies.csv, ratings.csv, and tags.csv files. No demographic information is included. Each user is represented by his/her ID and no other information is provided.



Graph 1: Count of ratings of all movies in dataset.

IV. PROPOSED SYSTEM

Proposed systems include implementation of a content-based filtering approach using TF-IDF, a cosine similarity method as explained in paper [2] and a collaborative approach using deep learning by authors Ram Murti Rawat et. al. [10]. The main advantage of this system is that the algorithms are designed to work efficiently even on small datasets. Recommendations are based on movie plot information, history, user profile information, and user likes and dislikes.

A. TF-IDF and Cosine Similarity

This method is commonly used as part of content-based recommender systems. It consists of two terms, Term Frequency (TF) and Inverse Document Frequency (IDF). Term frequency

deals with the frequency of interests and favorites in user profiles. Inverse document frequency, on the other hand, deals with the inverse of term frequency in the total data provided by user profiles. These two concepts are combined to provide recommendations to users based on data provided by user profiles as explained in [14].

$$W_{i,j} = tf_{i,j} \cdot \log\left(\frac{N}{df_i}\right)$$

$tf_{i,j}$ = number of occurrences of i in j

df_i = number of documents

N = total number of documents

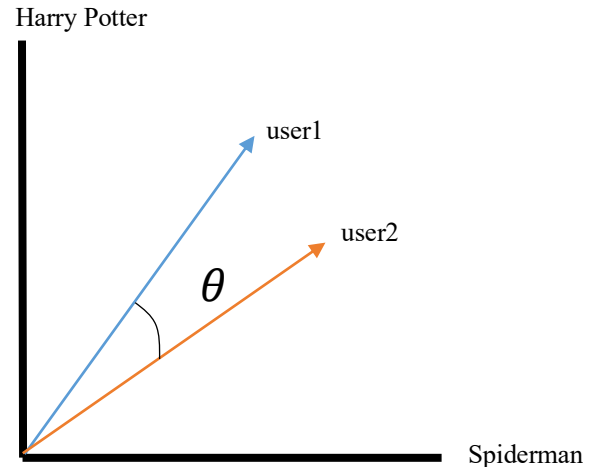
Here, we have the term frequency, or the number of times a certain term (such as a genre) appears in a document (such as a movie's genres), multiplied by the right-side term, which essentially scales the term frequency according to the number of times a particular word appears in all documents (movies).

The resultant weight is higher when there are fewer films that belong to a certain genre (df_i). The main purpose of the logarithm is to minimize large discrepancies caused by the right-hand component by smoothing the division result.

The movie genres and titles are transformed as vectors in a geometric space. Therefore, the angle between two vectors represents the closeness of those two vectors. Cosine similarity is a complex concept that has been widely discussed in information retrieval. This algorithm converts text documents into concept vectors. This model can find similarities between two data sets by determining the cosine value between the two vectors. The smaller the angle between two vectors, the larger the cosine and the higher the similarity.

The cosine similarity of two vectors A and B is represented using a dot product and magnitude as,

$$\cos(\theta) = \frac{A \cdot B}{||A|| ||B||}$$



Graph 3: Cosine similarity between movie 'Harry Potter' and 'Spiderman'

Graph 3 shows the cosine angle between two movies 'Harry Potter' watched by user1 and 'Spiderman' enjoyed by user2. This explains how relevant these two movies are for users user2 and user1, respectively.

B. Embedding Layer

In an embedding, discrete or categorical variables are transformed into a vector of continuous numerical values. Embeddings are multi-dimensional vectors in the deep learning perspective that represent discrete or categorical variables using learnt continuous values. The embeddings are used to simplify the representation of categorical variables and readily express categories in a modified space that can be fed into neural networks and other machine learning techniques.

C. Embedded-based Deep Learning

Collaborative filtering becomes very important as this technique is used in many recommender systems. Memory-based and model-based are two main types of collaborative filtering techniques. Memory-based collaborative filtering detects similar users in user space that match a given user and dynamically recommends movies to that user. This approach has his two problems: high computational complexity and small amount of data. A model-based collaborative filtering approach pre-builds a model based on a user's previous ratings of items, and the model provides item recommendations. Models are built using a variety of machine learning and neural network algorithms.

The proposed network predicts ratings for movies that users have not watched. Machine learning algorithms expect numbers as inputs, so the system uses embedding to convert the user ID and movie ID into fixed-size continuous vectors.

D. Network Architecture

User IDs and movie IDs are passed directly to the deep learning network as the model attempts to find relationships between non-existent IDs. From the embedding spaces of 610 users and 9742 movies, a user ID embedding vector and a movie ID embedding vector are created, respectively. The vectors are then concatenated together and passed to the dropout layer. A dropout layer is used to avoid overfitting the model. At each training epoch, Dropout randomly sets the layer's outgoing edge to 0 with a predefined probability. The received input is forwarded to the next layer, the hidden layer. Multiple fully connected hidden layers with dropout probabilities and ReLU activation features are added. Each hidden layer can have different values of output characteristics and dropout probabilities. The input features are determined based on the previous layer's output features. Finally, the network needs to send the predicted score as output for the input user id and model id. To do this, it uses an output layer with a sigmoid activation function and rescale it afterwards in the range 1-5. The neuron weights in each layer are updated at each training epoch using a backpropagation algorithm. The structure of the proposed

network is as shown in Fig 1 similar to network proposed by Ram Murti Rawat et. al. [10].

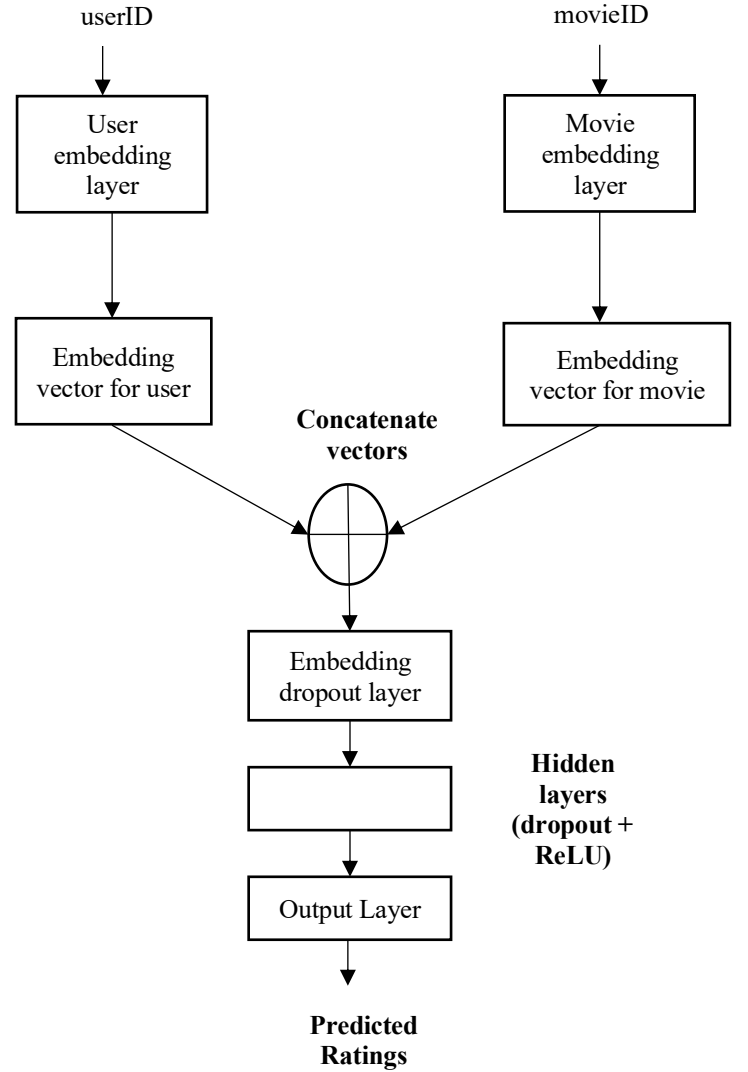


Fig1: Neural Network Architecture

E. Evaluation

First, some pre-processing of the data is done on. The dataset is divided into three parts in which 70% is for training data, 10% for validation, and the remaining 15% for testing. To check the performance of recommender system model, the RMSE is calculated during training, validation and testing.

Root mean squared error (RMSE) is a common metric used to calculate system accuracy. RMSE is computed by taking the square root of the mean squared error. The average is calculated by dividing the sum of the squared differences between the actual output and the predicted output by the number of samples. The equation below represents the RMSE.

$$RMSE = \sqrt{\sum_{k=0}^N \frac{(predicted_k - actual_k)^2}{N}}$$

N is the total number of instances. The smaller RMSE value means better performance for the system. Thus, a system which produces smaller RMSE value will provide better recommendations.

To compare the performance of the content-based and collaborative recommender system model, the precision and recall is calculated. Precision and recall are binary metrics used to evaluate models with binary output. Thus, there is a need to translate our numerical problem (ratings usually from 1 to 5) into a binary problem (relevant and not relevant items). To do the translation, assume that any true rating above 2.5 corresponds to a relevant item and any true rating below 2.5 is irrelevant. A relevant item for a specific user-item pair means that this item is a good recommendation for the user in question.

In the context of recommendation systems, recommending the user the top-N items is probably what we are most interested in doing. Therefore, it is more logical to compute precision and recall metrics for the first N items rather than the entire set of objects. As a result, there is the concept of precision and recall at k, where k is a user-definable integer that is specified by the user to correspond to the purpose of the top-N suggestions as explained in [15].

Precision at k is the proportion of recommended items in the top-k set that are relevant.

$$Precision@k = \frac{\# \text{ of recommended items @k that are relevant}}{\# \text{ of recommended items @k}}$$

Recall at k is the proportion of relevant items found in the top-k recommendations.

$$Recall @k = \frac{\# \text{ of recommended items @k that are relevant}}{\text{total \# of relevant items}}$$

V. RESULTS AND DISCUSSION

The precision and recall score obtained for content-based recommendation system (TF-IDF and Cosine Similarity) on test data is 5.6% and 4.3% respectively. Whereas the same evaluation metric provided 57% precision and 68% recall in collaborative-based approach (deep learning).

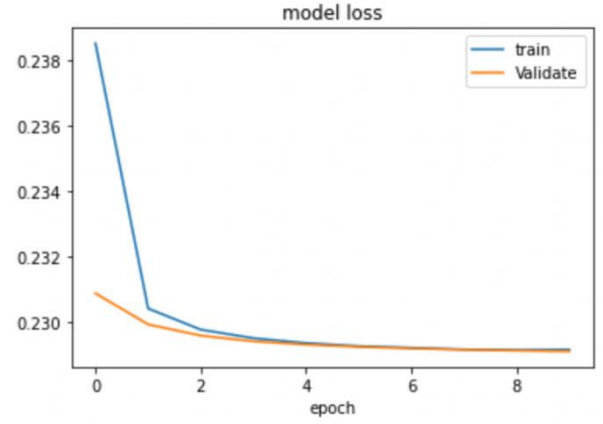


Fig 2 : Training and Validation Loss over Number of Epoch for Collaborative filtering approach

Fig. 2 shows how the training and validation errors are changing with increasing training epochs of the collaborative filtering model. The training error is decreasing continuously along with validation error till epoch 6 but almost stays the same after that. For deeper neural network at epoch 10 RMSE is 0.22. Therefore, the better prediction accuracy results can be explained by the fact that the use of users and items embeddings with deep neural network is very efficient.

The drawback of the used recommendation system is limited content analysis which leads to less accuracy of the recommendation system. Collaborative approach faces a cold-start problem where the profile of new user or item will be empty since he or she has not rated any item; hence, his or her taste is not known to the system. But, content-based filtering approach have the ability to recommend new items even if there are no ratings provided by users. So even if the database does not contain user preferences, recommendation accuracy is not affected. The deep learning technique with embedding layers has an additional advantage of scalability.

VI. CONCLUSION AND FUTURE WORK

In this work, a content-based recommender system is implemented using TF-IDF, cosine similarity, and a collaborative recommender system using an embedding-based deep learning approach. In a collaborative approach, user ratings of movies in the test data are predicted using deep learning methods, and an evaluation metric is calculated based on these predictions. Based on metrics, the collaborative model outperforms content-based on precision and recall scores.

Perhaps the biggest issue facing recommender systems is that they need a lot of data to effectively make recommendations. Vast amounts of data required to produce accurate recommendations open scope to other applications such as incorporation of big data tools and efficient data processing methodologies. The key to finding the perfect recommendation system for the user's needs lies in combination of content-based and collaborative approach into hybrid recommendation system. By implementing both in a single system, it can overcome drawbacks from both types of filtering techniques [11].

REFERENCES

- [1] [B.F. Oladapo . “A Research Proposal on Paper Recommender System”, unpublished, 2013.](#)
- [2] [N. Muthurasu, Nandhini Rengaraj, Kavitha Conjeevaram Mohan . “Movie Recommendation System using Term Frequency-Inverse Document Frequency and Cosine Similarity Method”, IJRTE, 2019.](#)
- [3] [Donghui Wang, Yanchun Liang, Dong Xu, Xiaoyue Feng, Renchu Guan. “A content-based recommender system for computer science publications”, Knowledge based Systems, 2018.](#)
- [4] [W.Croft, T.Strohman, and D.Metzler “SearchEngines: Information Retrieval in Practice”, Pearson Education, Inc., 2105.](#)
- [5] [Katarya R, Verma OP. “An effective collaborative movie recommender system with cuckoo Search”. Egyptian Informatics Journal, 2017.](#)
- [6] [C.Volinsky, R.Bell, Y.Koren. “MatrixFactorization Techniques for Recommender Systems”, IEEE, 2009.](#)
- [7] [Xiangnan He, Lizi Liao, Hanwang Zhang. “Neural Collaborative Filtering”, International World Wide Web Conference Committee, 2017.](#)
- [8] [Snehal R . Chavare, Chetan J. Awati , Dr. Suresh K. Shirgave. “Smart Recommender System Using Deep Learning”, Sixth International Conference on Inventive Computation Technologies, 2021.](#)
- [9] [Jeffrey Lund, Yiu-Kai Ng. “Movie Recommendation using deep learning approach”. IEEE Computer Society, 2018.](#)
- [10] [Ram Murti Rawat, Vikranth Tomar, Vinay Kumar. “An Embedding-based Deep Learning Approach for Movie Recommendation”, IEEE Xplore, 2020.](#)
- [11] [Robin Burke. “Hybrid Recommender Systems : Suvery and Experiments”, User Modeling and User-Adapted Interaction, 2002.](#)
- [12] [<https://towardsdatascience.com/prototyping-a-recommender-system-step-by-step-part-1-knn-item-based-collaborative-filtering-637969614ea>](#)
- [13] [F.Maxwell Harper and Joseph A. Konstan. 2015. “The MovieLens Datasets: History and Context.”, ACM Transactions on Interactive Intelligent Systems \(TiiS\), 2015.](#)
- [14] [<https://towardsdatascience.com/content-based-recommender-systems-28a1dbd858f5>](#)
- [15] [\[https://medium.com/@m_n_malaeb/recall-and-precision-at-k-for-recommender-systems-618483226c54\]\(https://medium.com/@m_n_malaeb/recall-and-precision-at-k-for-recommender-systems-618483226c54\)](#)