
Creation of Movie Recommender

Danielė Ačaitė, Evita Petkevičiūtė, Rūta Okulič-Kazarinaitė

Abstract

1 The World Wide Web continues to grow at a massive rate and all the sites with
2 all the information and their complexity continues to grow along with it. Thus, it
3 becomes time consuming and burdensome for people to find similar information
4 themselves. To enhance the search for relevant things and make it easier to use
5 it, websites can be personalized. This is where recommender systems become
6 valuable because they can dynamically add hyperlinks which results in making it
7 easier for users to find interesting things and in the meantime making the interaction
8 between the system and the user better. These personalized recommender systems
9 are used in various movie sites.

10 Movie recommendation systems are increasingly popular and widely used in the
11 entertainment industry. They aim to predict the ratings of preferences a user would
12 give to a movie, based on the user's past ratings and the ratings of other users. By
13 recommending movies that are likely to be highly rated by a user, the system can
14 help users discover new films to watch and improve their experience. However,
15 it requires expensive computations in order to find the essential properties as the
16 number of users as well as movies increase exponentially.

17 According to M.Sridevi and Dr. R.Rajeswara Rao [1], there are several approaches
18 to building recommender systems, including Collaborative filtering, Content based
19 filtering, Demographic filtering and Hybrid Recommenders system. Collaborative
20 filtering is based on the idea of predicting users' ratings for a movie based on the
21 ratings given by other users with similar tastes. Content-based filtering uses the
22 characteristics of the movies themselves (for example: genre, director, actors) to
23 recommend similar movies to users. Hybrid methods combine both collaborative
24 filtering and content-based filtering to make recommendations.

25 In this research paper, we will explore the use of machine learning techniques to
26 build a movie recommendation system based on I. Ahmed method [2], . We will
27 discuss the challenges and limitations of building movie recommendation systems.

1 Project goals and Outcomes

29 Our objective is to create a movie recommendation system based on 4 types of recommender systems
30 which is acknowledged by literature review.

31 The principal steps of our study are: 1)Reviewing the types of movie recommender systems in
32 literature 2)Finding the dataset required for the project 3)Performing the modeling 4)Report results
33 and explain all the steps of the modeling process thoroughly

34 Completing this project will require familiarising with various types of movie recommenders: Collab-
35 orative filtering, Content based filtering, Demographic filtering and Hybrid Recommenders system.
36 Gaining practice and expertise in distinguishing these types and learn reproducible experimental
37 design.

38 2 Literature Research

39 2.1 COLLABORATIVE FILTERING

40 Collaborative filtering is a method of making recommendations for a user by collecting preferences
41 or taste information from many other users. This information is then used to identify patterns and
42 suggest items that are likely to be of interest to the user. In the case of a movie recommender system,
43 collaborative filtering involves collecting ratings or reviews from many users and using those ratings
44 to suggest movies to a particular user. The assumption behind collaborative filtering is that if a user
45 A has similar preferences to a user B, and B likes a particular movie, then A is more likely to like
46 that movie as well. M.Sridevi and Dr .R.Rajeswara Rao address 3 main issues of this type of movie
47 recommender system [3]:

- 48 • Most users do not rate the movies.
- 49 • The recommender fails to recommend a movie when new users come to the environment
50 since they do not rate movies and the system lacks data.
- 51 • As the number of users increases exponentially it becomes expensive and computationally
52 hard to compute an accurate system.

53 2.2 CONTENT-BASED RECOMMENDER SYSTEM

54 Content-based filtering systems choose items based on the correlation between items' content and
55 the user's preferences as opposed to the collaborative filtering system meaning that they suggest
56 items based on a particular topic [4]. Recommendations are made by comparing a profile with the
57 matter of each file in the correlation where there are a lot of wordings. This system has item metadata
58 where there is information about the director, genre, description, actors, etc for films to to create these
59 recommendations. Some words in the plot description like tags and stop words are removed because
60 they occur pretty often and cannot be used as discriminators. The remaining usually are minimized
61 by removing prefixes and suffixes. In the end, content-based recommender systems look for items
62 that a person has liked and then search for similar ones. Indicating that the movies who were the most
63 liked will likely be similar and recommended.

64 Steps that we taken:

- 65 • Overview columns indicate the plot's description, therefore, a conversion into matrix by
66 computing TF-IDF vectors was needed.
- 67 • The cosine similarity was used to denote the similarity between two different movies.
- 68 • Define recommendation system.

69 2.3 DEMOGRAPHIC FILTERING

70 Demographic filtering is a method of making recommendations for a user based on demographic
71 information about the user. The demographic properties may be things like the user's age, gender,
72 location, or occupation. This information can be used to suggest movies that are likely to be of interest
73 to the user based on the preferences or characteristics of people with similar demographics. For
74 example, if a recommender system knows that a user is a 20-year-old woman, it might recommend
75 romantic comedies or dramas, which are known to be popular among young women. It is more
76 advantageous than collaborative filtering since it addresses the issue of new users with limited data
77 and is easier to implement. If new users have not rated any movies the system can recommend movies
78 based on one's demographic information. On the other hand, it lacks uniqueness. For instance, let's
79 say there is a new young male user and the system recommends action movies as it is the most
80 popular among males. However, it does not capture that this male user might not enjoy action movies
81 since he is more of a romantic comedies fan.

82 2.4 HYBRID FILTERING

83 Hybrid filtering is a technique used in movie recommendation systems that combines the properties
84 of collaborative filtering and content-based filtering. Collaborative filtering relies on the preferences

85 of similar users to make recommendations, while content-based filtering uses the characteristics of
86 the movie itself to make recommendations. By combining these two approaches, hybrid filtering
87 is able to make more accurate and personalized recommendations. For example, a hybrid system
88 might consider a user's past ratings of comedies, as well as the genre and actors of the movie being
89 recommended, to provide a more tailored suggestion. Hybrid filtering can help to improve the
90 effectiveness of movie recommendation systems by combining the two methods.

91 **3 Methodology**

92 **3.1 Datasets**

93 All provided datasets were extracted from kaggle [5]. The features of the datasets were displayed in
94 the coding environment.

95 **4 Demographic filtering modification**

96 This paper examines the demographic filtering methodology based on the TMDb (The Movie
97 Database) website. This is one of the major databases about movies and television shows. TMDb
98 allows users to access metadata for various tv shows, movies. There users are able to find the status
99 of the movies, get a list of trending films and get the features. One of the simplest metrics that could
100 be used is movie rating. However, this method is not as advantageous since, firstly, it does not take
101 into account the popularity of the movie. If there is a movie with rating of 9 by 50 000 users will
102 have a lower place than a movie which was rated 9.5 by only three users, hence, the list would not be
103 reliable. Therefore, one of the main metrics that are used when building a movie recommender is the
104 weighted rate used in the TMDb [6].

$$105 \text{WR} = (v/(v+m))*R + (m/(v+m))*C$$

- 106 • v - number of votes for the movie -> already given
- 107 • m - minimum number of votes that keep movie in the list (we take 80
- 108 • R - mean rating of the movie -> already given
- 109 • C - mean rating of all the movies in the list (calculated in the code, mean is 6.092)

110 Steps:

- 111 • filtering of movies that qualify for the list → 96
- 112 • define function of tmdb weighted rating and create new feature score and calculate the value
113 by applying a function to our DataFrame of qualified movies
- 114 • sorting movies based on score:

115 These are the results of basic Demographic Filtering Movie Recommender. The table of 10 movies
116 displays the list of the movies which can be treated as “Trending now” movies. That proves that
117 demographic filtering does not take into account the individual preferences as tastes of a particular
118 user.

119 **5 CONTENT-BASED FILTERING**

120 This recommendation system is based on the pairwise similarity index. The index is used according
121 to the plot descriptions and movies are recommended based on the similarity score. Each word vector
122 needs to be converted by computing TF-IDF vectors. This provides a matrix whose each column is a
123 word in the overview vocabulary and each row shows the original movie name. Multiplied TF-IDF
124 is used to minimize the importance of repeating words that are found in overview. After using the
125 scikit-learn function it showed that there are more than 20 000 different words that describe more
126 than 4800 films in the given dataset. Cosine similarity was used to calculate a numeric quantity that
127 denotes the likeness between two different movies. This index is more independent of magnitude
128 and is in comparison not difficult and faster to calculate. In order the system would recommend, the

	title	vote_count	score
1881	The Shawshank Redemption	8205	8.248353
662	Fight Club	9413	8.096134
3337	The Godfather	5893	8.077404
3232	Pulp Fiction	8428	8.074738
65	The Dark Knight	12002	8.044250
809	Forrest Gump	7927	7.972814
96	Inception	13752	7.969290
95	Interstellar	10867	7.937399
1990	The Empire Strikes Back	5879	7.904757
1818	Schindler's List	4329	7.900080

Figure 1: List representing movies based on weighted rating

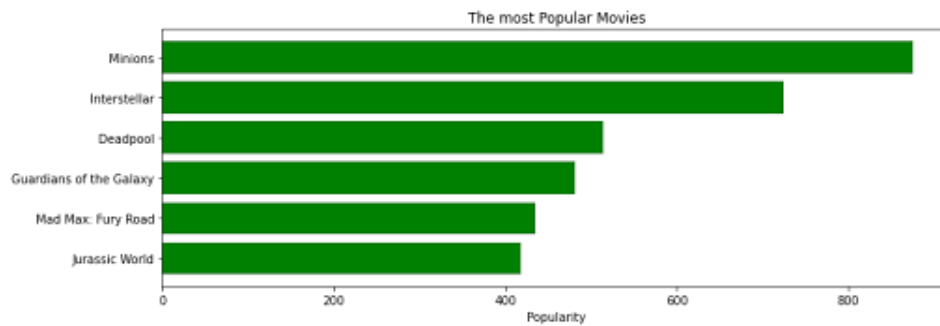


Figure 2: List representing movies based on popularity

129 output is set to a movies name and the input is set to the movies names and their description of other
130 liked 10 similar movies. Later for this, a reverse napping of movie titles and DataFrame indices are
131 needed.

```

▶ get_recommendations("Grown Ups")

4420          Black Rock
1816          The Best of Me
1057          Coach Carter
3382          The Mighty Macs
3478          College
4541          The Slaughter Rule
2562    When the Game Stands Tall
507          Independence Day
2228          An American Carol
3900          Air Bud
Name: title, dtype: object

```

Figure 3: Recommendations made using TF-IDF

132 However, the given output quality can be improved because the recommendation system does not
133 take into account the cast, crew, keywords and genres. After converting the list into a safe and
134 usable structure, the functions are written to help to extract the needed information from each feature.

135 Later all the words are put into lowercases and the spaces are moved so that the system would
 136 calculate the name and surnames as one word. The string “metadata soup” is created where all the
 137 metadata is because it is needed to give to the vectorizer. Now the countvectorizer() is used instead
 138 of TF-IDF because there is no need to down-weight the presence of actors or directors. Finally, the
 recommendation system outputs similar movies.

```
| get_recommendations('Grown Ups', cosine_sim2)

436          Grown Ups 2
499          Jack and Jill
2989         Happy Gilmore
361    You Don't Mess with the Zohan
445          Just Go with It
1539         Big Daddy
796    The Ridiculous 6
802    That's My Boy
1126    Here Comes the Boom
1392    The Benchwarmers
Name: title, dtype: object
```

Figure 4: Recommendations made using countvectorizer()

139

140 6 COLLABORATIVE FILTERING

141 To provide an even more personalized recommendation system, it is worthy to implement collaborative
 142 filtering. In this case, we are looking at specific persons (userID) ratings for certain movies, and try
 143 to predict another person’s rating for the same movie. That can be used to recommend movies to
 144 people in the same interest groups.

145 To implement this model, we used Surprise Python Scikit due to its specific recommendation functions.
 146 We use a 5 fold cross-validation, to predict new data, together with RMSE and MAE to predict the
 147 quality of the predictions - we aim for a low RMSE for more precise results. We get the RSME equal
 148 to 0.89684716, which is decent. Then we train the model:

```
trainset = data.build_full_trainset()
svd.fit(trainset)
```

Figure 5: Training the data set

149 Now to check the results we pick a random user based on their id, and predict the rating for the movie
 (also based on the id), in this case: Here, we predict what user (10) would rate the movie (20). SVD -

```
[108] # predicting users rating for a movie
svd.predict(10, 20, 3)

Prediction(uid=10, iid=20, r_ui=3, est=3.0636579409030085, details={'was_impossible': False})
```

Figure 6: Predicting the ratings

150

151 Singular Value Decomposition algorithm is used to model the user and item biases from users and
 152 items, which uses stochastic gradient descent to optimize the parameters.

153 7 Summary

154 Current trends of movie recommender systems show that there is huge potential for them to improve to
 155 even more precise level. Our project summarizes and explains how the recommend systems work with

156 certain data and models. Digging deeper into recommender system issues, there are few main ones to
157 be named: one of them is a slow start for new users - their experience will not be pleasant and they
158 will not receive personalized recommendations until they have watched enough movies. Another issue
159 is that some recommender systems are vulnerable to feedback loops, when recommendations made
160 by the system influence the data that is used to generate future recommendations. Another problem is
161 that if the data that the models were trained on were biased, so will be the recommendations, and
162 the model will not understand that. This and many other issues do not have a quick or easy fix -
163 it requires attention and ethical decisions. To start with, the data set could be including metadata,
164 diverse (including different genders, age groups, races etc.), and updated often [7]. Besides these
165 issues, movie recommender systems can benefit not only the business owner by increasing profits and
166 user spent time on the platform, but it can benefit the user as well by saving them time and helping to
167 find good fitting movies faster.

168 **References**

- 169 [1] M. Sridevi. Table 1 from decors : A simple and efficient demographic collaborative recommender
170 system for movie recommendation: Semantic scholar, Jan 1970.
- 171 [2] Ibtesama. Getting started with a movie recommendation system, May 2020.
- 172 [3] Iateilang Ryngksai and L Chameikho. Recommender systems: types of filtering techniques.
173 *International Journal of Engineering Research & Technology, Gujarat*, 3(2278-0181):251–254,
174 2014.
- 175 [4]
- 176 [5] Rounak Banik. The movies dataset, Nov 2017.
- 177 [6]
- 178 [7] Jiawei Chen, Hande Dong, Xiang Wang, Fuli Feng, Meng Wang, and Xiangnan He. Bias and
179 debias in recommender system: A survey and future directions. *arXiv preprint arXiv:2010.03240*,
180 2020.