

User Profile Correlation-Based Similarity Algorithm in Movie Recommendation System

Triyanna Widiyaningtyas

Universitas Gadjah Mada

Indriana Hidayah

Universitas Gadjah Mada

Teguh Bharata Adji (✉ adji@ugm.ac.id)

Universitas Gadjah Mada <https://orcid.org/0000-0001-7856-1498>

Research

Keywords: collaborative filtering, user rating value, user behavior value, UPCSim

Posted Date: November 18th, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-107401/v1>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published at Journal of Big Data on March 29th, 2021.
See the published version at <https://doi.org/10.1186/s40537-021-00425-x>.

RESEARCH

User profile correlation-based similarity algorithm in movie recommendation system

Triyanna Widiyaningtyas^{1,2,7,8}, Indriana Hidayah¹ and Teguh B. Adji^{1*}

*Correspondence: adji@ugm.ac.id

¹Department of Information Technology and Electrical Engineering, Universitas Gadjah Mada, Yogyakarta, Indonesia
Full list of author information is available at the end of the article

Abstract

A recommendation system is a software used in the e-commerce field that provides recommendations for customers to choose the items they like. Several recommendation systems have been proposed; however, collaborative filtering is the most widely used approach. The main issue in collaborative filtering is how to implement a similarity algorithm that can improve performance in the recommendation system. Several similarity algorithms based on user rating value have been developed, and recently a similarity algorithm has been developed that combines the user rating value and the user behavior value. However, the existing research is still based only on a single user behavior value, which is the genre data. Therefore, we propose a new similarity algorithm that considers not only the genre data but also the user profile data (namely age, gender, occupation, and location). The new similarity we are proposing is called User Profile Correlation-based Similarity (UPCSim). The user profile correlation similarity was obtained by calculating the correlation coefficient between the user profile data and the user rating or behavior values. An experiment was done to compare the accuracy of the UPCSim algorithm with that of the previous algorithm. The experiment results showed that the UPCSim algorithm can improve the recommendation performance MAE by 1.64% and RMSE by 1.4% compared to the previous algorithm.

Keywords: collaborative filtering; user rating value; user behavior value; UPCSim

Introduction

The exponential growth of information on the internet causes users can get more information resources to dig up and collect. However, users will get lost in the sea of information and will have difficulty in processing that information [1, 2]. Users have to spend more time and more energy finding the information they want, but users may not necessarily get satisfactory results. Fortunately, user behavior on e-commerce sites and other social networks can be recorded and be tracked, making it easier to analyze user interests [3, 4]. One of the tools used to solve this problem in analyzing user interest is known as a recommendation system.

The recommendation system is software that helps users to get relevant items from millions of items in the database [5, 6]. The recommendation system's main task is to offer users personalized item recommendations through information filtering. This system has become a commercial platform that assists the user in providing suggestions for items to be selected. The provided suggestions are useful to support users in various decision-making processes, such as what books to read, which locations to visit, what news to read, and more [7].

¹ Based on the utilized data source and computation method, the recommendation¹
² system is divided into three approaches, namely: collaborative filtering, content-²
³ based filtering, and hybrid filtering [6, 8]. The collaborative filtering approach uses³
⁴ the collaborative power of ratings given by all users to make recommendations. The⁴
⁵ content-based filtering approach uses descriptive attributes of items to make rec-⁵
⁶ ommendations. Meanwhile, the hybrid filtering approach combines several filtering⁶
⁷ methods to get a list of items according to user preferences [9].⁷

⁸ Of the three recommendation system approaches, collaborative filtering is one⁸
⁹ of the most popular, successful, and widely used methods in recommendation sys-⁹
¹⁰ tems [5, 10, 11]. This recommendation method is widely used because it is simple,¹⁰
¹¹ efficient, and has an acceptable accuracy level. Based on the advantages of this¹¹
¹² collaborative filtering, several real-world systems have used this method, such as¹²
¹³ Amazon, MovieLens, and Netflix [8, 10, 12].¹³

¹⁴ The collaborative filtering approach is categorized into two approaches, viz¹⁴
¹⁵ ranking-oriented collaborative filtering and rating-oriented collaborative filtering¹⁵
¹⁶ [13, 14]. Ranking-oriented collaborative filtering directly provides a preference or-¹⁶
¹⁷ der for items from users without predicting ratings for items that have not been¹⁷
¹⁸ rated. In contrast, rating-oriented collaborative filtering predicts ratings for items¹⁸
¹⁹ that have not been rated by users based on rating information from other users [14].¹⁹
²⁰ The rating-oriented collaborative filtering approach is more widely used because it²⁰
²¹ is faster in generating recommendations than the ranking-oriented collaborative²¹
²² filtering approach.²²

²³ The rating-oriented collaborative filtering approach is categorized into two meth-²³
²⁴ ods, viz the model-based method and the memory-based method [8, 15, 16]. The²⁴
²⁵ model-based method uses a rating database to build a model and uses that model to²⁵
²⁶ predict ratings for unrated items [17]. Some of the techniques that are often used to²⁶
²⁷ build models include clustering [18, 19], Bayesian network [20], Markovian factor-²⁷
²⁸ ization [21], and Singular Value Decomposition (SVD) [22, 23]. This model-based²⁸
²⁹ method has a drawback; namely, the computational complexity is very dependent²⁹
³⁰ on the model being built [3]. Therefore, the algorithm complexity value cannot be³⁰
³¹ ascertained.³¹

³² Meanwhile, the memory-based method uses the rating database to calculate the³²
³³ similarity between users or similarity between items [24]. In its implementation, this³³
³⁴ method is divided into two techniques, namely User-Based Collaborative Filtering³⁴
³⁵ (UBCF) and Item-Based Collaborative Filtering (IBCF) [1, 2, 25]. The UBCF pre-³⁵
³⁶ dicts ratings for all items that have not been rated based on the similarity of users,³⁶
³⁷ while the IBCF predicts ratings based on the similarity of items [26]. Some of the³⁷
³⁸ frequently used traditional similarities are the Cosine (COS), Pearson Correlation³⁸
³⁹ Coefficient (PCC), and Jaccard [2, 8].³⁹

⁴⁰ In recent years, the majority of researchers have focused more on developing a sim-⁴⁰
⁴¹ ilarity algorithm between users/items, because of its simplicity in computation. For⁴¹
⁴² example, Patra *et al.* [16] proposed a similarity algorithm known as Bhattacharyya⁴²
⁴³ similarity. Furthermore, Polatidis *et al.* [27] proposed a similarity algorithm that⁴³
⁴⁴ is an increase in PCC similarity, known as multi-level collaborative filtering. Their⁴⁴
⁴⁵ proposed similarity is an increase in the PCC similarity by taking into account the⁴⁵
⁴⁶ number of items co-rated on several levels. Sun *et al.* [28] proposed a similarity⁴⁶

¹algorithm by integrating the similarity of Triangle and Jaccard, known as TMJ¹
²similarity. Feng et al. [2] proposed a new similarity algorithm by integrating three²
³factors of similarity impact, namely S_1, S_2 , and S_3 . With S_1 expressing similarity³
⁴between users, S_2 is used to calculate the number of items co-rated less than the⁴
⁵specified threshold, and S_3 explains the weight of each user's rating preference. The⁵
⁶four similarity algorithms that have been developed by previous researchers perform⁶
⁷similarity calculations based only on the user's rating value data. ⁷

⁸ Furthermore, Wu et al. [29] proposed a new similarity by combining similarity⁸
⁹based on the user's rating value and similarity based on user behavior value. Wu⁹
¹⁰et al. [29] view that the traditional recommendation algorithm only pays attention¹⁰
¹¹to the rating value given by users explicitly, but ignores implicit information on¹¹
¹²user behavior in assessing the item type (genre) that will also affect the accuracy¹²
¹³level. In their research, similarity based on the user's rating value (S_r) was calculated¹³
¹⁴using the UBCF method. Meanwhile, similarity based on user behavior value (S_b) is¹⁴
¹⁵calculated from the probability of the user score in giving a rating to the genre. The¹⁵
¹⁶proposed similarity algorithm is known as the User Score Probability Collaborative¹⁶
¹⁷Filtering (UPCF). ¹⁷

¹⁸ By combining the two similarities, the similarity calculation between users not¹⁸
¹⁹only require user rating value data but also genre data to calculate the user be-¹⁹
²⁰havior value. The research results can provide an increase in accuracy compared²⁰
²¹to the application of the traditional similarity algorithm, which only considers the²¹
²²user rating value. However, there is still a need to improve the recommendations²²
²³accuracy by exploring other user behavior affecting user interests. ²³

²⁴ Therefore, we propose a new similarity algorithm – that is called User Profile²⁴
²⁵Correlation-based Similarity (UPCSim) – which not only pays attention to genre²⁵
²⁶data, but also adds other user behavior data in term of user profile data (namely age,²⁶
²⁷gender, occupation, and location) on giving weight to similarity. In this UPCSim²⁷
²⁸algorithm, the similarity weighting technique is obtained by calculating the corre-²⁸
²⁹lation coefficient between the user profile data (namely age, gender, occupation,²⁹
³⁰and location) and the user rating value for the similarity weight of S_r , while the³⁰
³¹similarity weight of S_b is obtained by calculating the correlation coefficient between³¹
³²the user profile data and the user behavior value. ³²

³³ The structure of this paper is as follows. In "Similarity Algorithm" section de-³³
³⁴scribes the several similarity algorithm which developed in the collaborative filter-³⁴
³⁵ing approach. Then, "Research Method" section explains the proposed similarity³⁵
³⁶algorithm in detail. Subsequently, "Experiment" section present the experiment's³⁶
³⁷results using the MovieLens dataset and its discussion. Finally, "Conclusion and³⁷
³⁸Future Work" section provides some conclusions and suggestions for further research³⁸
³⁹development". ³⁹

⁴⁰ **Similarity Algorithm** ⁴¹

⁴²In the previous section, it has been explained that the frequently used traditional⁴²
⁴³similarity algorithms in recommendation systems are the Cosine similarity (COS)⁴³
⁴⁴and the Pearson Correlation Coefficient (PCC) [30]. To explain this similarity cal-⁴⁴
⁴⁵culatation, we assume that the user and item sets are defined as $U = \{u_1, u_2, \dots, u_m\}$ ⁴⁵
⁴⁶and $I = \{i_1, i_2, \dots, i_n\}$. The user rating matrix for the item is denoted as⁴⁶

$^1R = [r_{ui}]^{m \times n}$, where m and n are the number of users and the number of items,
 2 respectively, and r_{ui} is the rating given by user u on item i
 3 Cosine similarity measures the angle between two rating vectors (user or item)
 4 [2, 3, 16]. The Cosine similarity formula between user u_1 and user u_2 is stated in
 5 (1).
 6

$$^7 \quad \text{Sim}(u_1, u_2)^{COS} = \frac{\vec{r}_{u_1} \cdot \vec{r}_{u_2}}{\|\vec{r}_{u_1}\| \cdot \|\vec{r}_{u_2}\|} = \frac{\sum_{i \in I_{u_1} \cap I_{u_2}} r_{u_1 i} \cdot r_{u_2 i}}{\sqrt{\sum_{i \in I_{u_1} \cap I_{u_2}} r_{u_1 i}^2} \cdot \sqrt{\sum_{i \in I_{u_1} \cap I_{u_2}} r_{u_2 i}^2}} \quad (1)^8$$

10 Pearson similarity measures how two users or items are linearly related to each
 11 other [3, 16]. After identifying the items rated jointly between user u_1 and user u_2 ,
 12 Pearson similarity calculates the linear correlation between the two users using the
 13 formula specified in (2). Pearson similarity values ranged in the range $[-1, +1]$. A
 14 value of $+1$ indicates a very high correlation and -1 indicates a negative correlation.
 15

$$^{17} \quad \text{Sim}(u_1, u_2)^{PCC} = \frac{\sum_{i \in I_{u_1} \cap I_{u_2}} (r_{u_1 i} - \bar{r}_{u_1}) \cdot (r_{u_2 i} - \bar{r}_{u_2})}{\sqrt{\sum_{i \in I_{u_1} \cap I_{u_2}} (r_{u_1 i} - \bar{r}_{u_1})^2} \cdot \sqrt{\sum_{i \in I_{u_1} \cap I_{u_2}} (r_{u_2 i} - \bar{r}_{u_2})^2}} \quad (2)^{18}$$

21 Several similarity improvements are continuously being developed to increase the
 22 recommendations' accuracy. Among them are Bhattacharyya's similarity [16], the
 23 multi-level collaborative filtering similarity [27], the TMJ similarity [28], and the
 24 similarity integrating three impact factors, namely S_1 , S_2 , dan S_3 [2]. These four
 25 similarity algorithms only consider user rating value data explicitly to calculate the
 26 similarity between users. These similarity algorithms assume that users who give a
 27 high rating on an item indicate that they like the item, while users who give a low
 28 rating indicate that they do not like the item.
 29

29 In the current era of online shopping, users assess items based on their quality,
 30 delivery speed, service attitude to customers, and other factors. If the user is not
 31 interested in an item, the user will not select/buy the item. The illustration of
 32 user behavior in assessing items is shown in Figure 1. The User u likes to watch
 33 genre animation and does not like to watch drama and adventure genres. The User
 34 u selects the movie title "Toy Story" from a collection of animated films. After
 35 watching the movie, the User u gave a low rating to "Toy Story" because the "Toy
 36 Story" movie gave unclear sound quality.
 37

37 In the similarity algorithm, which is only based on the user rating value, if a
 38 user gives a low rating, it means that the user does not like the movie title as well
 39 as the movie genre. This indicates that the user's preference for similar items will
 40 decrease and will affect the resulting recommendations. Therefore, it is necessary
 41 to involve user behavior in rating items to calculate similarity as a form of invisible
 42 user preferences.
 43

43 To adopt this user behavior, Wu et al. [29] developed a new similarity algorithm
 44 that involves not only explicit rating data but also user behavior data in giving an
 45 implicit rating. Wu et al. [29] assumed that users who gave a low rating to an item
 46 did not necessarily dislike the item. The similarity algorithm combines similarity

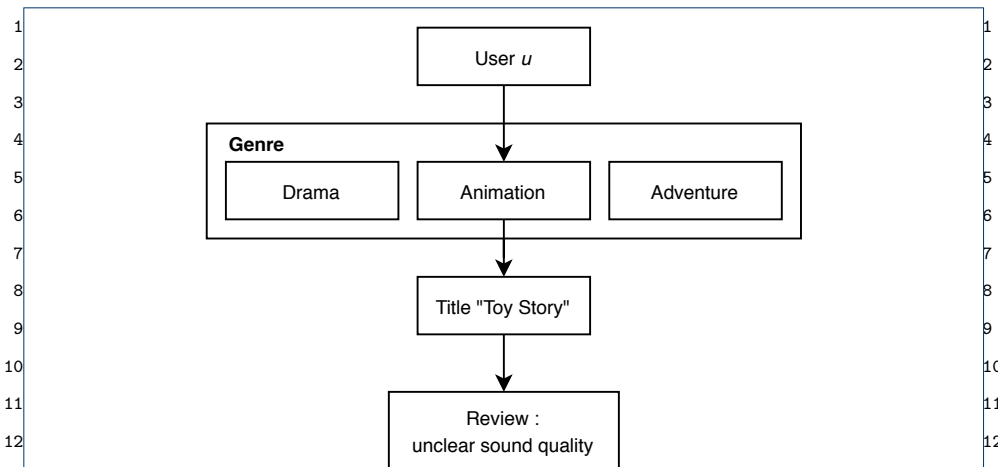


Figure 1 Example of user behavior in rating the item

based on user rating value (S_r) and similarity based on user behavior value (S_b). The formula of this similarity is shown in (3) below.

$$Sim(u_1, u_2)^{UPCF} = \beta S_r(u_1, u_2) + (1 - \beta) S_b(u_1, u_2) \tag{3}$$

$S_r(u_1, u_2)$ is the similarity between user u_1 and user u_2 calculated by the UBCF method, β is the threshold (which is between 0 to 1) that can be set to the average value of the similarity of all the users who are similar to the active user u_1 , and S_b is defined in (4).

$$S_b(u_1, u_2) = \frac{\sum_{g \in G_{u_1} \cap G_{u_2}} (P_{u_1g} - \bar{P}_{u_1}) \cdot (P_{u_2g} - \bar{P}_{u_2})}{\sqrt{\sum_{g \in G_{u_1} \cap G_{u_2}} (P_{u_1g} - \bar{P}_{u_1})^2} \cdot \sqrt{\sum_{g \in G_{u_1} \cap G_{u_2}} (P_{u_2g} - \bar{P}_{u_2})^2}} \tag{4}$$

G_{u_1} and G_{u_2} are the set of item types (genres) rated by users u_1 and users u_2 respectively. P_{u_1g} and P_{u_2g} are the probability scores for the type of item g from users u_1 and users u_2 respectively. \bar{P}_{u_1} and \bar{P}_{u_2} are the average probability scores for all item types from users u_1 and users u_2 , respectively, and g is one type of items rated by both users.

The combination of the two similarities in the research conducted by Wu et al. [29] has several limitations, viz the calculation of similarity based on user behavior value only considers the genre data of the item so that it does not guarantee the resulted recommendations accuracy. Based on the problem, this paper focuses on increasing accuracy by considering other user behavior data that is the user profile data (namely age, gender, occupation, and location) that will influence user behavior in determining the selected item.

We assumed that age, gender, occupation, and location would influence the user's interest in the item. As an illustration, young users would have different preferences from older users. Female users would have different preferences from male users. Users with a technician job would have different preferences from users with a lawyer job, and users living on the coast would have different preferences from

¹users living in cities. Based on this assumption, the main contribution of our study¹
²is the proposed new similarity algorithm that considers not only genre but also age,²
³gender, occupation, and location data to improve accuracy in dealing with items³
⁴that have not been rated by users (data sparsity). The detail of this algorithm is⁴
⁵be described in the next section.⁵

⁷Research Method ⁷

⁸The conducted research was more focused on the proposed UPCSim algorithm. This⁸
⁹UPCSim algorithm would be compared with the UPCF algorithm [29], because the⁹
¹⁰UPCF already accommodates one user behavior, viz the genre data. Meanwhile, the¹⁰
¹¹UPCSim algorithm adds other user behavior data, namely the user profile data. For¹¹
¹²this comparison, an MBCF system applying the UPCSim algorithm was created. It¹²
¹³should be noted that all the pre-processing and processing of MBCF we did were¹³
¹⁴the same as that done by Wu et al. [29]. So that the comparison between the two¹⁴
¹⁵studies is equal.¹⁵

¹⁶ The MBCF system that we have developed is a development of the MBCF system¹⁶
¹⁷that was previously developed by Wu et al. [29]. A detailed illustration of the¹⁷
¹⁸proposed MBCF system is shown in Figure 2.¹⁸

¹⁹ The system is divided into four blocks, namely input, data preparation, the MBCF¹⁹
²⁰process, and output. The input block is the input dataset used in the MBCF system.²⁰
²¹The data preparation block consists of the data pre-processing stage, which results²¹
²²in data source ready for use in the MBCF process. The data source includes the²²
²³rating data and the behavior data. While the behavior data used in Wu et al.²³
²⁴[29] only employs the genre data (the green component in data preparation block),²⁴
²⁵our research also accommodated the user profile data (the red component in data²⁵
²⁶preparation block). The MBCF process block is the development of the MBCF²⁶
²⁷method using the weighting similarity. The similarity weighting carried out by Wu²⁷
²⁸et al. used a threshold value ranging from 0 to 1 (the green component in the²⁸
²⁹MBCF process). Whilst, the weighting similarity in our research used the coefficient²⁹
³⁰correlation between the user profile data and the user rating or behavior values (the³⁰
³¹red component in MBCF process). Furthermore, the final block of the developed³¹
³²system is the output block that provides an evaluation of the UPCSim algorithm³²
³³application in the MBCF method.³³

³⁴ The detail of our proposed UPCSim algorithm is explained as a component of³⁴
³⁵similarity calculation, which is presented in "Similarity Calculation" section. Mean-³⁵
³⁶while the detail of the developed MBCF system is explained in "Developed MBCF³⁶
³⁷System" section.³⁷

³⁹Similarity Calculation ³⁹

⁴⁰In this study, the similarity calculation between users was divided into three com-⁴⁰
⁴¹ponents. The first component is the S_r similarity calculation component (shown in⁴¹
⁴²the dashed blue box). The second block is the S_b similarity calculation component⁴²
⁴³(shown in the dashed green box). Finally, the UPCSim block (shown in the dashed⁴³
⁴⁴red box) gives weight to both similarities. Details of the three component in our⁴⁴
⁴⁵similarity calculation can be seen in Figure 3.⁴⁵

⁴⁶ Each component in Figure 3 is described as follows.⁴⁶

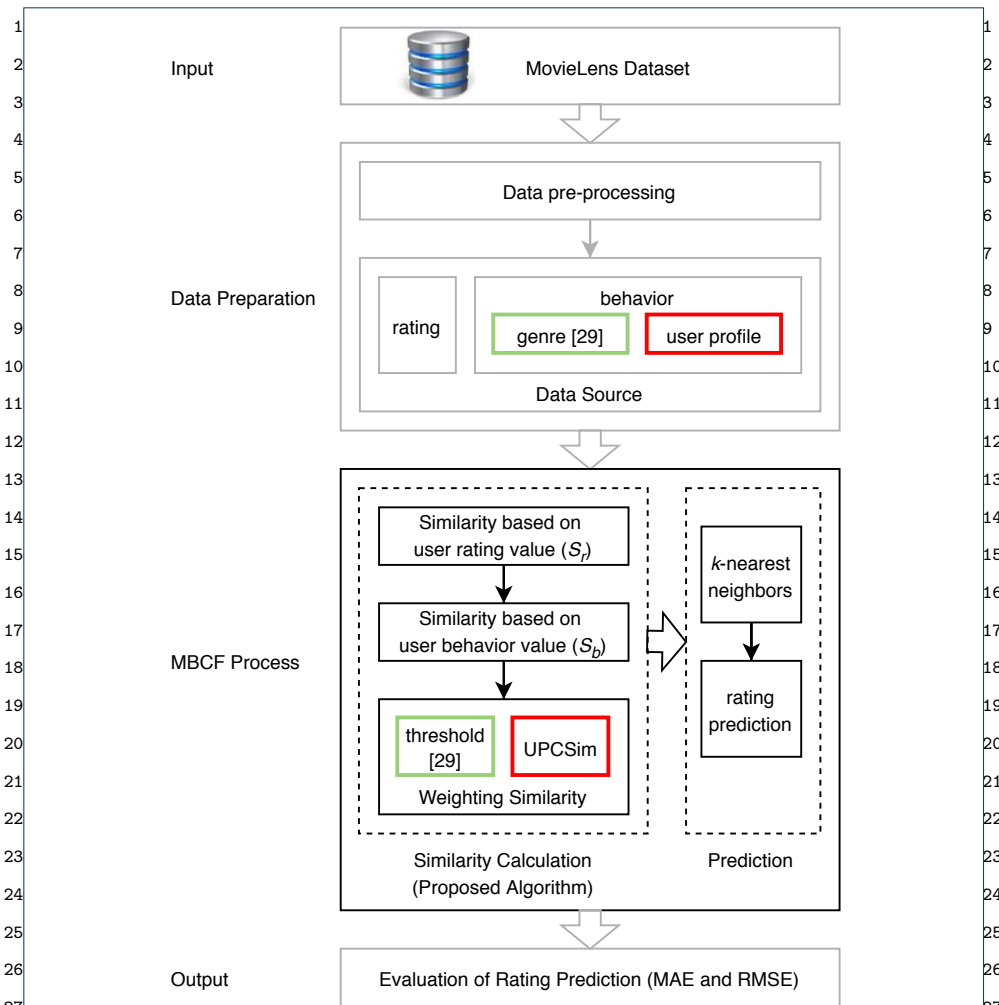


Figure 2 The developed MBCF system

S_r Similarity

The S_r similarity component is the similarity calculation based on the user rating value. As an example of the similarity calculation, we used the MovieLens 100K dataset. The initial stage in calculating the S_r similarity was done by reading the rating data from the resulted pre-processing data. Based on the rating data, we obtained a rating matrix of 943×1682 . The value of 943 represented the number of users, and the value of 1682 represented the number of movies in the dataset. The rating matrix illustration is shown below.

$$R = \begin{bmatrix} R_{11} & R_{12} & R_{13} & R_{14} & \dots & R_{1_1682} \\ R_{21} & R_{22} & R_{23} & R_{24} & \dots & R_{2_1682} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ R_{943.1} & R_{943.2} & R_{943.3} & R_{943.4} & \dots & R_{943_1682} \end{bmatrix}$$

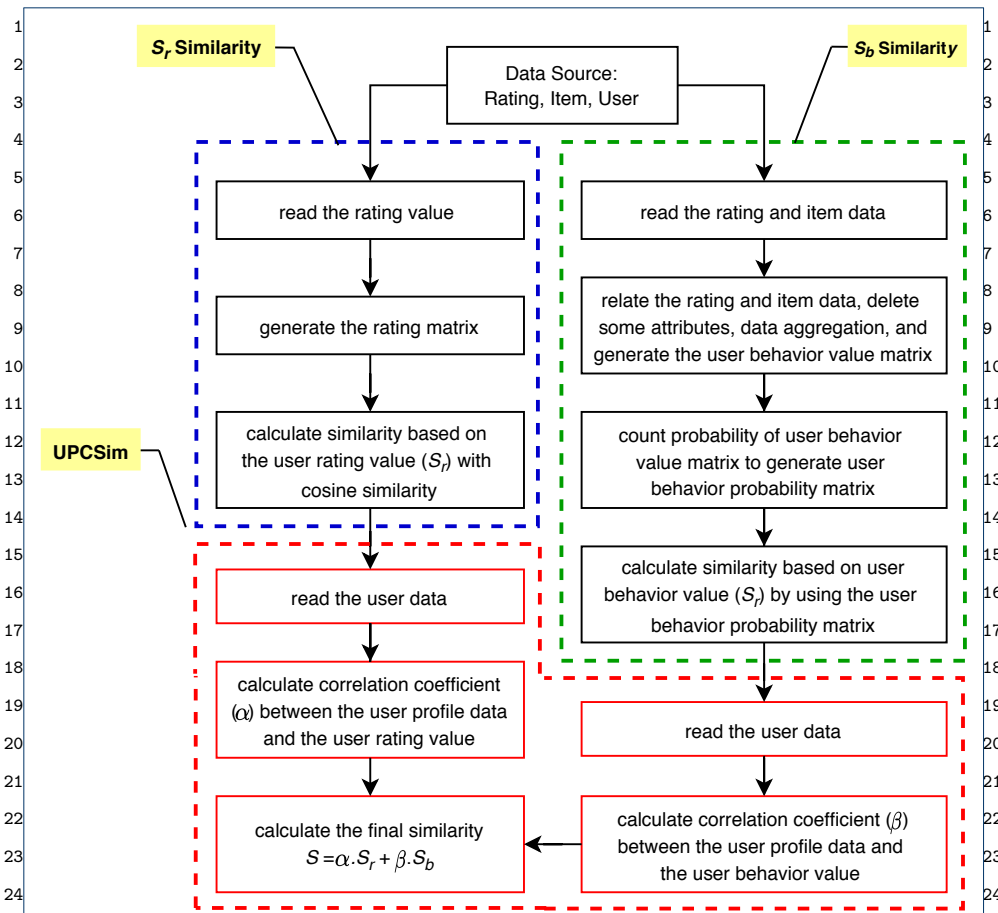


Figure 3 Flowchart of similarity calculation

$R_{943.1682}$ is the rating value given by the 943rd user for the 1682nd item. Values R_{11} to $R_{943.1682}$ range from 0 to 5, with a value of 0 indicating that the user unrated the item.

After the rating matrix is formed, the next step was to calculate the S_r similarity using the Cosine similarity formula referring to (1). The final result of the S_r similarity calculation would form the S_r similarity matrix with the order 943×943 . Illustration of the similarity matrix S_r is shown as follows.

$$S_r = \begin{bmatrix} S_{11} & S_{12} & S_{13} & \dots & S_{1.943} \\ S_{21} & S_{22} & S_{23} & \dots & S_{2.943} \\ \dots & \dots & \dots & \dots & \dots \\ S_{943.1} & S_{943.2} & S_{943.3} & \dots & S_{943.943} \end{bmatrix}$$

$S_{1.943}$ is the similarity value based on the rating between the 1st user and the 943rd user.

S_b Similarity

The S_b similarity component is the similarity calculation based on the user behavior value. The initial stage in calculating the S_b similarity was done by reading the

rating and item data. In the MovieLens 100K dataset, the rating data stated the user's rating value for the watched movies, and the item data stated the movie title data containing the genre information of each movie. Each movie title can include several genres. For example, the movie "Toy Story" has the genre of animation, children, and comedy. The rating and item data are used to calculate the user behavior value.

The user behavior value could be calculated by relating the rating data and the item data, removing some unused attributes from the results of the relationship between the rating data and the item data, and performing data aggregation using the sum function grouped by user. This data aggregation results are illustrated in the user behavior value matrix with the order matrix of 943×19 , where the value 943 represents the number of users, and the value 19 represents the number of genres that exist. The illustration of user behavior value matrix is shown as follows.

$$B = \begin{bmatrix} B_{11} & B_{12} & B_{13} & \dots & B_{1,19} \\ B_{21} & B_{22} & B_{23} & \dots & B_{2,19} \\ \dots & \dots & \dots & \dots & \dots \\ B_{943,1} & B_{943,2} & B_{943,3} & \dots & B_{943,19} \end{bmatrix}$$

$B_{943,19}$ is the 943rd user behavior value for the 19th genre, representing the total number of 19th genre watched by 943rd user. After the user behavior value matrix was formed, the next stage was to calculate the probability of genres occurrence from the user behavior value matrix to produce a probability matrix of user behavior value using (5) below.

$$P = \frac{B(g)}{N} \quad (5)$$

$B(g)$ is the value of user behavior for the target genre g , and N is the total number of users who gave rate to the target genre g . The illustration of the probability matrix of user behavior value is shown as follows.

$$B = \begin{bmatrix} P_{11} & P_{12} & P_{13} & \dots & P_{1,19} \\ P_{21} & P_{22} & P_{23} & \dots & P_{2,19} \\ \dots & \dots & \dots & \dots & \dots \\ P_{943,1} & P_{943,2} & P_{943,3} & \dots & P_{943,19} \end{bmatrix}$$

$P_{943,19}$ is the probability value of the 943rd user behavior for the 19th genre. The probability matrix of user behavior value was used as the basis for calculating the S_b similarity referring to (4). The results of the S_b similarity calculation would form a matrix with the order 943×943 . The illustration of the S_b similarity matrix is shown as follows.

$$S_b = \begin{bmatrix} S_{11} & S_{12} & S_{13} & \dots & S_{1,943} \\ S_{21} & S_{22} & S_{23} & \dots & S_{2,943} \\ \dots & \dots & \dots & \dots & \dots \\ S_{943,1} & S_{943,2} & S_{943,3} & \dots & S_{943,943} \end{bmatrix}$$

¹ $S_{1,943}$ is the similarity value based on the user behavior value the 1st user and the
²943rd user.

⁴*UPCSim*

⁵The UPCSIm component is the component of the similarity calculation using the
⁶UPCSIm algorithm, which calculates the weights of both similarities (S_r and S_b)
⁷based on user profile. The initial stage in calculating the weights was done by read-
⁸ing the user data. In the MovieLens 100K dataset, the user data stated the user
⁹profile data that consisted of age, gender, occupation, and location. The user pro-
¹⁰file data was used to calculate the weights of S_r and S_b similarities. The weights of
¹¹these two similarities were calculated based on the correlation coefficient (R) using
¹²multiple linear regression. The weight of S_r similarity was obtained by calculating
¹³the correlation coefficient between user profile data (age, gender, occupation, and
¹⁴location) and the user rating value, which was then symbolized by α . While the
¹⁵weight of S_b similarity was obtained by calculating the correlation coefficient be-
¹⁶tween user profile data (age, sex, occupation, and location) and the user behavior
¹⁷value, which was then symbolized by β .

¹⁸ After weighting the two similarities, the next stage was to calculate the final
¹⁹similarity matrix by combining the weighted S_r and S_b similarities. By combining
²⁰these two similarities, we could obtain the final similarity matrix S with the order
²¹matrix 943×943 . The formula of the final similarity between user u and user v is
²²defined in (6) below.

$$S(u, v) = \alpha S_r(u, v) + \beta S_b(u, v) \quad (6)$$

²⁶ $S(u, v)$ is the final similarity between user u and user v . $S_r(u, v)$ is the similarity
²⁷based on user rating value between user u and user v . $S_b(u, v)$ is the similarity based
²⁸on user behavior value between users u and user v . α is the weight of the similarity
²⁹ S_r , and β is the weight of the similarity S_b .

³¹Developed MBCF System

³²Based on the illustration shown in Figure 2, this section describes each block of the
³³developed MBCF system.

³⁵*Input*

³⁶The first block of the developed MBCF system was the input dataset. In this paper,
³⁷we used the MovieLens dataset (<https://grouplens.org/datasets/movielens/>). This
³⁸dataset was collected by the “GroupLens Study Group of the University of Min-
³⁹nesota” [31]. There were several versions of the dataset, including ml-100K, ml-1M,
⁴⁰ml-10M, ml-20M, etc. In this experiment, we chose the dataset used in previous
⁴¹studies [29], namely ml-100K. This ml-100K dataset contained several data files.
⁴²Our study used 3 data files, namely rating data, item data, and user data..

⁴³ The rating data consisted of 100,000 ratings as rated by 943 users on 1682 movies.
⁴⁴Each user has rated at least 20 movies. The rating value given by the user ranged
⁴⁵from 1 to 5. A score of 1 stated that the user did not like the watched movie. A
⁴⁶score of 5 stated that the user liked the watched movie. This rating data had a

¹sparsity of 93.7% and a density of 6.3%. This rating data structure consisted of¹
²user-id, movie-id, rating, and timestamp. 2

³ Item data contains information about items (movies). This item's data structure³
⁴consisted of 24 attributes, namely, movie-id, movie title, release date, video release⁴
⁵date, IMDb URL, and 19 attributes of movie type (genre). Each item (movies) can⁵
⁶have several genres. 6

⁷ User data contains information about the user's profile. This user data structure⁷
⁸consisted of 5 attributes, namely user-id, age, gender, occupation, and zip code⁸
⁹(which states the user's location). 9

¹¹*Data Preparation* 11

¹²The second block was data preparation. In this section, data pre-processing was¹²
¹³carried out. The purpose of this process was to prepare data to obtain a quality¹³
¹⁴dataset. In this study, data pre-processing was done by reducing attributes that¹⁴
¹⁵were not relevant to data processing. 15

¹⁶ Based on the existing file data structure, several attributes were not needed in data¹⁶
¹⁷processing. These attributes included the timestamp on the rating data, movie title,¹⁷
¹⁸release date, video release date, and IMDb URL on the item data. At the dataset¹⁸
¹⁹preparation stage, these attributes would be deleted so that the data processing¹⁹
²⁰stage ran more effectively. 20

²²*MBCF Process* 22

²³The third block of the MBCF system was the MBCF process. The MBCF pro-²³
²⁴cess was divided into two sub blocks, namely the similarity calculation and the²⁴
²⁵prediction. 25

²⁶ The similarity calculation was the initial process used in the information filter-²⁶
²⁷ing process using the MBCF approach. In this study, the similarity calculation was²⁷
²⁸divided into three components, namely S_r similarity calculation, S_b similarity cal-²⁸
²⁹ulation, and the UPCSim algorithm. In detail, these three components have been²⁹
³⁰described in "Similarity Calculation" section. 30

³¹ The prediction was carried out to provide a predicted rating for items that had not³¹
³²been rated by active users. The initial stage taken in this prediction was to determine³²
³³the number (k) of the nearest neighbors of an active users. k was an integer number³³
³⁴representing the number of neighbors, ranging from 10 to 100 [2, 27–29]. After the³⁴
³⁵ k value was determined, the next stage was to determine the rating prediction for³⁵
³⁶some unrated items. 36

³⁷ The formula used to determine the predicted rating for an item (i) unrated by³⁷
³⁸active users (u) is shown in (7) below [2, 16]. 38

$$p_{ui} = \bar{r}_u + \frac{\sum_{v \in NN_u} S(u, v) \cdot (r_{vi} - \bar{r}_v)}{\sum_{v \in NN_u} |S(u, v)|}, v \neq u \quad (7)$$

⁴³ p_{ui} represents the predicted rating value of user u to item i . \bar{r}_u dan \bar{r}_v is the average⁴³
⁴⁴rating of user u and user v respectively. r_{vi} is the rating value given by user v to⁴⁴
⁴⁵item i . $S(u, v)$ is the final similarity between user u and user v . NN_u is the set of⁴⁵
⁴⁶nearest neighbors to user u . 46

¹*Output* 1

²The fourth block of the MBCF system was the output block. This block was used to ²
³evaluate the performance of the UPCSim Algorithm in predicting ratings for items ³
⁴that had not been rated by an active user. 4

⁵ To measure the performance of the recommendation system, mean absolute error ⁵
⁶(MAE), rooted mean squared error (RMSE), precision, and recall were the most ⁶
⁷popular measures. According to Jalili et al. [32], the metrics for evaluating recom- ⁷
⁸mendation systems can be classified into two categories, namely prediction metrics ⁸
⁹and classification metrics. The MAE and RMSE are primarily used to evaluate ⁹
¹⁰prediction metrics [33, 34], whereas precision and recall are used to evaluate classi- ¹⁰
¹¹fication metrics, namely evaluation of the quality of top-N recommendations [3]. 11

¹² In this study, we adopted the MAE and RMSE metrics to measure the prediction ¹²
¹³metrics of the UPCSim Algorithm. The MAE is the most widely used metric in ¹³
¹⁴recommendation systems with a collaborative filtering approach. It is used to esti- ¹⁴
¹⁵mate the average absolute deviation between the actual and the predicted rating ¹⁵
¹⁶values. A lower MAE provides good recommendation quality [35]. The formula for ¹⁶
¹⁷calculating MAE is defined in (8). 17

$$\sup_{19} MAE = \frac{1}{TN} \sum_{u \in U, i \in I} |p_{ui} - r_{ui}| \quad (8) \quad \sup_{20}$$

²²Meanwhile, RMSE reflects the degree of deviation between the predicted rating and ²²
²³the actual rating. A lower RSME is associated with higher prediction metrics [36]. ²³
²⁴The RMSE formula is defined in (9). 24

$$\sup_{26} RMSE = \sqrt{\frac{1}{TN} \sum_{u \in U, i \in I} (p_{ui} - r_{ui})^2} \quad (9) \quad \sup_{27}$$

²⁹ TN is the total number of predicted items. p_{ui} and r_{ui} represent the predicted ²⁹
³⁰rating and actual rating of the user u to item i , respectively. 30

³²Experiment 32

³³This section begins with the experiment design described in "Experiment Design" ³³
³⁴section. Next, the "Experiment Result and Analysis" section explain the comparison ³⁴
³⁵between the proposed UPCSim algorithm and the previous similarity algorithms, ³⁵
³⁶namely the Cosine similarity algorithm and the UPCF similarity algorithm. The ³⁶
³⁷comparison utilized the MAE and RMSE values of the MBCF system that had ³⁷
³⁸been made. Finally, "Discussion" section provides conclusions from the experiment ³⁸
³⁹results. 39

⁴¹Experiment Design 41

⁴²To evaluate the performance of the proposed UPCSim algorithm, the experiment ⁴²
⁴³design in this study implemented the following four steps: 43

⁴⁴ The first step was to divide the dataset. The dataset would be divided into two ⁴⁴
⁴⁵parts, viz training data, and testing data. The k -fold cross-validation method was ⁴⁵
⁴⁶used in dividing the dataset. In this experiment, the chosen k was 5. It meant ⁴⁶

¹that 80% of the dataset was used for training data and the remaining 20% was for¹
²testing data. The training data were train1, train2, train3, train4, and train5, and²
³the testing data are called test1, test2, test3, test4, and test5.³

⁴ The second step was to calculate the similarity matrix between users. The S_r ⁴
⁵similarity was obtained based on the user rating value matrix, and the S_b similarity⁵
⁶was obtained based on the user behavior value matrix. Then, we calculated the⁶
⁷weights of the two similarities using the correlation coefficient (R), with multiple⁷
⁸linear regression analysis. The final similarity matrix was obtained by combining⁸
⁹the two weighted similarities.⁹

¹⁰ The third step was to calculate the predicted ratings for the testing data. The¹⁰
¹¹nearest neighboring k was selected based on the final similarity matrix. In this¹¹
¹²experiment, the k values varied from 10 to 100, with an increase in the k value by¹²
¹³10.¹³

¹⁴ The fourth step was to measure the proposed UPCSim algorithm’s performance¹⁴
¹⁵using the MAE and RMSE prediction metrics.¹⁵

¹⁶
¹⁷Experiment Results and Analysis¹⁷

¹⁸This section aims to compare the proposed UPCSim algorithm’s performance re-¹⁸
¹⁹sults with the traditional Cosine similarity and the UPCF similarity. The MAE and¹⁹
²⁰RMSE values were obtained by comparing the three algorithms using variations in²⁰
²¹the number of different neighbors. Experiments were carried out in five iterations.²¹
²²The first iteration was done using the train1 and test1 dataset. The second iteration²²
²³was done using the train2 and test2 datasets, and so on. The number of nearest²³
²⁴neighbors used in this experiment ranged from 10 to 100, which were used in each²⁴
²⁵train1, train2, train3, train4, and train5 data. The three similarity algorithms (Co-²⁵
²⁶sine, UPCF, and UPCSim) were applied to the five-training data. The average MAE²⁶
²⁷performance for the three algorithms is shown in Table 1.²⁷

²⁸In Table 1, MAE_c is the average MAE value of the UBCF experiment using²⁸
²⁹Cosine similarity. The MAE_p is the average MAE value of the UBCF experiment²⁹
³⁰using UPCF similarity. The MAE_ps is the average MAE value of the UBCF ex-³⁰
³¹periment using the proposed UPCSim similarity. The MAE_ps-c is the difference³¹
³²between the MAE value using the proposed UPCSim similarity and the MAE value³²
³³using the Cosine similarity. And the MAE_ps-p is the difference between the MAE³³
³⁴values using the proposed UPCSim similarity and the MAE values using the UPCF³⁴
³⁵similarity.³⁵

³⁶**Table 1** Comparison of average MAE values of the three algorithms³⁶

Number of Neighbors	MAE_c	MAE_p	MAE_ps	MAE_ps-c	MAE_ps-p
10	0.8227	0.7792	0.7669	0.0558	0.0123
20	0.8099	0.7631	0.7483	0.0616	0.0148
30	0.8051	0.7565	0.7410	0.0641	0.0155
40	0.8048	0.7551	0.7387	0.0661	0.0164
50	0.8043	0.7544	0.7369	0.0674	0.0175
60	0.8041	0.7535	0.7364	0.0677	0.0171
70	0.8049	0.7529	0.7359	0.0690	0.0170
80	0.8051	0.7526	0.7355	0.0696	0.0171
90	0.8056	0.7525	0.7347	0.0709	0.0178
100	0.8072	0.7521	0.7337	0.0735	0.0184
average	0.8074	0.7572	0.7408	0.0666	0.0164

⁴⁵ Based on Table 1, an increase in the number of nearest neighbors results in a⁴⁵
⁴⁶decreased MAE value in the UPCF similarity algorithm and the UPCSim similarity⁴⁶

¹algorithm. For the Cosine algorithm, the MAE value decreases only to $k = 60$. It ¹
²shows that the number of nearest neighbors influenced MAE value; in other words, ²
³the number of nearest neighbors affects the algorithm's performance. The smallest ³
⁴MAE value is obtained in the UPCSim algorithm, which means that this algorithm's ⁴
⁵error prediction is the smallest. So, it can be said that UPCSim's algorithm is more ⁵
⁶reliable than others. Compared with the Cosine similarity, the UPCSim algorithm ⁶
⁷has an increase in MAE values ranging from 5.58% to 7.35%, with the average ⁷
⁸MAE of 6.66% for all k nearest neighbors. Compared with the UPCF similarity, the ⁸
⁹UPCSim algorithm has an increase in MAE values ranging from 1.23% to 1.84%, ⁹
¹⁰with the average MAE of 1.64% for all k nearest neighbors. By referring to the ¹⁰
¹¹data in Table 1, the average MAE value of the three algorithms can be illustrated ¹¹
¹²graphically, as in Figure 4. ¹²

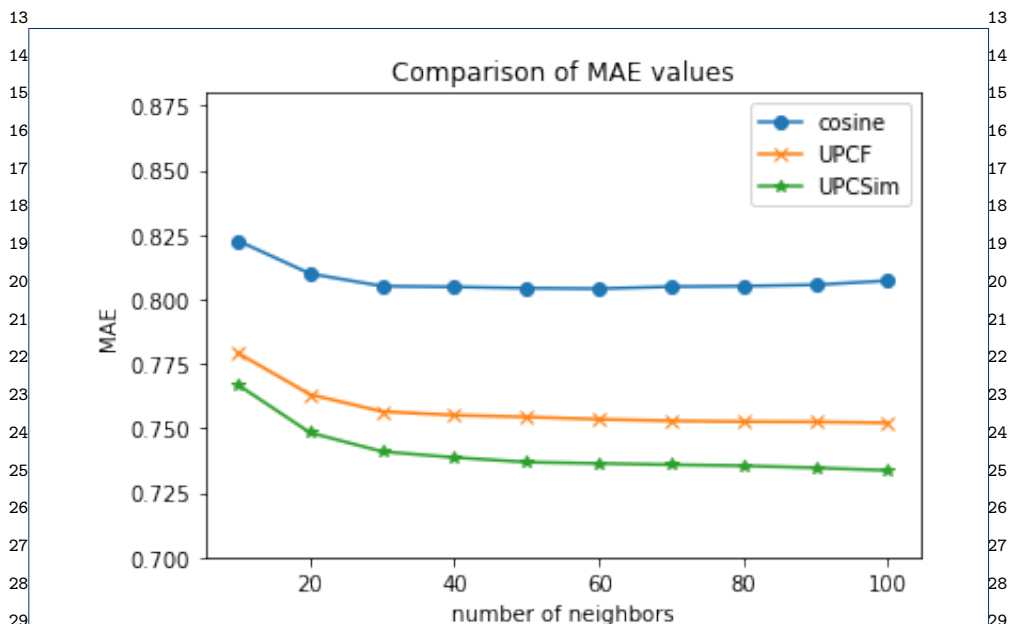


Figure 4 Comparison of the average MAE value of the three similarity algorithms using Movielens 100k dataset

³³ Figure 4 shows that the three algorithms' MAE value decreases with an increas- ³³
³⁴ing number of nearest neighbors. At the beginning of the curve, it can be seen that ³⁴
³⁵the decline in MAE value is very sharp with the increase in the number of nearest ³⁵
³⁶neighbors, while at the end of the curve, the greater the number of nearest neigh- ³⁶
³⁷bors, the MAE value tends to be stable. It can be said that the number of nearest ³⁷
³⁸neighbor variables affects the MAE value, where the greater the number of nearest ³⁸
³⁹neighbors, the smaller the MAE value. With the same number of nearest neighbors, ³⁹
⁴⁰the MAE value of the UPCSim algorithm is always smaller than that of the other ⁴⁰
⁴¹recommendation algorithms. In other words, the error between the actual rating ⁴¹
⁴²and the predicted rating of the proposed UPCSim algorithm is the smallest with a ⁴²
⁴³more accurate rating prediction. ⁴³

⁴⁴ Furthermore, the comparison of the average RMSE values of the three recommen- ⁴⁴
⁴⁵dation algorithms is shown in Table 2. The RMSE_c is the average RMSE value from ⁴⁵
⁴⁶the UBCF experiment using Cosine similarity. The RMSE_p is the average RMSE ⁴⁶

¹value of the UBCF experiment using UPCF similarity, and RMSE_{ps} is the average
²RMSE value of the UBCF experiment using the proposed UPCSim similarity. The
³RMSE_{ps-c} is the difference between the RMSE value based on the proposed UPC-
⁴Sim similarity and the RMSE value based on the Cosine similarity. The RMSE_{ps-p}
⁵is the difference between the RMSE value based on the proposed UPCSim similarity
⁶and RMSE value based on UPCF similarity.

⁷**Table 2** Comparison of average RMSE values of the three algorithms ⁷

⁸ Number of Neighbors	⁸ RMSE _c	⁸ RMSE _p	⁸ RMSE _{ps}	⁸ RMSE _{ps-c}	⁸ RMSE _{ps-p}
⁹ 10	1.0417	0.9835	0.9793	0.0624	0.0042
¹⁰ 20	1.0248	0.9643	0.9541	0.0707	0.0102
¹¹ 30	1.0189	0.9574	0.9453	0.0736	0.0121
¹² 40	1.0163	0.9558	0.9427	0.0736	0.0131
¹³ 50	1.0162	0.9556	0.9393	0.0769	0.0163
¹⁴ 60	1.0154	0.9547	0.9389	0.0765	0.0158
¹⁵ 70	1.0162	0.9544	0.9383	0.0779	0.0161
¹⁶ 80	1.0170	0.9543	0.9381	0.0789	0.0162
¹⁷ 90	1.0173	0.9542	0.9364	0.0809	0.0178
¹⁸ 100	1.0176	0.9538	0.9359	0.0817	0.0179
¹⁹ average	1.0201	0.9588	0.9448	0.0753	0.0140

¹⁷ As can be seen in Table 2, when the number of nearest neighbors is the same, the
¹⁸average RMSE value on the UPCSim similarity is always smaller than the other
¹⁹similarities. Meanwhile, an increase in the number of nearest neighbors results in
²⁰a decreased RMSE value in the three algorithms. This shows that the number of
²¹nearest neighbors influences the RMSE value. The UPCSim algorithm produces the
²²smallest RMSE value, which means that the proposed algorithm’s prediction error
²³is the smallest. So, it can be said that the UPCSim algorithm is superior. Compared
²⁴to the Cosine similarity, the UPCSim algorithm has an increase in RMSE values
²⁵ranging from 6.24% to 8.17%, with the average RMSE of 7.53% for all k nearest
²⁶neighbors. Compared with the UPCF similarity, the UPCSim algorithm has an
²⁷increase in the RMSE value ranging from 0.42% to 1.79%, with the average RMSE of
²⁸1.4% for all k nearest neighbors. Based on the data in Table 2, the three algorithms’
²⁹average RMSE value can be illustrated graphically, as shown in Figure 5.

³⁰ Figure 5 illustrates the effect of changes in the number of nearest neighbors on the
³¹RMSE value. Three algorithms show a decrease in the RMSE value first and tend
³²to be stable when the neighbors’ number is greater than 50. The RMSE value in
³³the UPCSim algorithm always shows the smallest value for each different number of
³⁴neighbors. It shows that the UPCSim algorithm has the lowest error rate compared
³⁵to the other two algorithms. So, it can be said that the UPCSim algorithm is
³⁶superior.

³⁷ The superiority of the UPCSim similarity algorithm can be obtained because the
³⁸algorithm considers the calculation of similarity weighting involving more complete
³⁹user behavior (genre and user profile) than that used in Wu et al. [29], which only
⁴⁰involved genre. So, the prediction metric is closer to the actual value. However, the
⁴¹UPCSim algorithm is considered computationally more complex than the previous
⁴²recommendation algorithms so that if it is implemented in a larger dataset, it will
⁴³take a longer time to produce recommendations.

⁴⁴Conclusion and Future Work ⁴⁴

⁴⁵Our experiment results on the MovieLens 100K dataset show that the UPCSim
⁴⁶algorithm can reduce MAE and RMSE values by 1.64% and 1.4%, respectively,

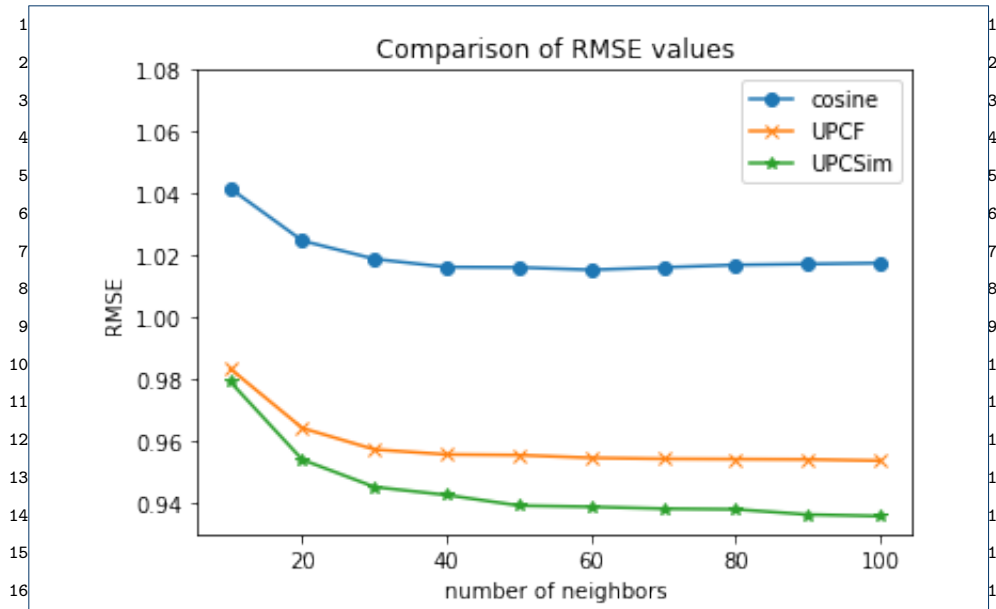


Figure 5 Comparison of the average RMSE value of the three similarity algorithms using MovieLens 100k dataset

compared to UPCF algorithm. The strength of our algorithm is the accommodation of the user profile for calculating the similarity weighting in order to capture the user interest more accurately. Although the UPCSim similarity algorithm can improve the accuracy of prediction metrics, this study still has some limitations. The UPCSim algorithm is more complex than the previous algorithms. Thus, it will be time consuming if it is applied in a larger dataset. Besides, the UPCSim similarity algorithm still needs to improve the resulting prediction metrics. Therefore, in future studies, a clustering method can be considered to overcome scalability problems due to the larger number of datasets and reduce computation time.

Acknowledgements

Not applicable

Funding

This research was funded by Indonesia Endowment Fund for Education (LPDP), Ministry of Finance of Republic of Indonesia: *Beasiswa Unggulan Dosen Indonesia - Dalam Negeri (BUDI – DN)* contract number 20200421211035.

Abbreviations

UPCSim: User Profile Correlation-based Similarity; MAE: Mean Absolute Error; RMSE: Root Mean Squared Error; SVD: Singular Value Decomposition; MBCF: Memory-Based Collaborative Filtering; UBCF: User-Based Collaborative Filtering; IBCF: Item-Based Collaborative Filtering; COS: Cosine similarity; PCC: Pearson Correlation Coefficient; TMJ: Triangle Multiplying Jaccard; UPCF: User Probability score Collaborative Filtering

Availability of data and materials

The datasets generated and analysed during the current study are available in the MovieLens dataset (<https://grouplens.org/datasets/movielens/>)

Ethics approval and consent to participate

Not applicable

Competing interests

The authors declare that they have no competing interests.

Consent for publication

Not applicable

1 Authors' contributions

2 All authors contributed both the concepts and contents of this study. TW provided the manuscript under supervised
3 by TBA and IH. All authors also performed discussion intensively for contents improvement. All authors have read
4 and approved the final manuscript.

5 Author details

6 ¹Department of Information Technology and Electrical Engineering, Universitas Gadjah Mada, Yogyakarta,
7 Indonesia. ²Departement of Electrical Engineering, Universitas Negeri Malang, Malang, Indonesia.

8 References

- 9 1. Xu, G., Tang, Z., Ma, C., Liu, Y., Daneshmand, M.: A collaborative filtering recommendation algorithm based
10 on user confidence and time context. *Journal of Electrical and Computer Engineering* **2019**, 1–12 (2019).
11 doi:[10.1155/2019/7070487](https://doi.org/10.1155/2019/7070487)
- 12 2. Feng, J., Fengs, X., Zhang, N., Peng, J.: An improved collaborative filtering method based on similarity. *Plos*
13 *One* **13**(9), 1–18 (2018). doi:[10.1371/journal.pone.0204003](https://doi.org/10.1371/journal.pone.0204003)
- 14 3. Liu, H., Hu, Z., Mian, A., Tian, H., Zhu, X.: A new user similarity model to improve the accuracy of
15 collaborative filtering. *Knowledge-Based Systems* **56**, 156–166 (2014). doi:[10.1016/j.knosys.2013.11.006](https://doi.org/10.1016/j.knosys.2013.11.006)
- 16 4. Camacho, L.A., Alves-Souza, S.N.: Social network data to alleviate cold-start in recommender system: A
17 systematic review. *Information Processing and Management* **54**, 529–544 (2018).
18 doi:[10.1016/j.ipm.2018.03.004](https://doi.org/10.1016/j.ipm.2018.03.004)
- 19 5. Sahu, A.K., Dwivedi, P.: User profile as a bridge in cross-domain recommender systems for sparsity reduction.
20 *Applied Intelligence* **49**, 2461–2481 (2019). doi:[10.1007/s10489-018-01402-3](https://doi.org/10.1007/s10489-018-01402-3)
- 21 6. Kumar, P., Kumar, V., Thakur, R.S.: A new approach for rating prediction system using collaborative filtering.
22 *Iran Journal of Computer Science* **2**, 81–87 (2019). doi:[10.1007/s42044-018-00028-5](https://doi.org/10.1007/s42044-018-00028-5)
- 23 7. Alonso, S., Bobadilla, J., Ortega, F., Moya, R.: Robust model-based reliability approach to tackle shilling
24 attacks in collaborative filtering recommender systems. *IEEE Access* **7**, 41782–41798 (2019).
25 doi:[10.1109/ACCESS.2019.2905862](https://doi.org/10.1109/ACCESS.2019.2905862)
- 26 8. Salah, A., Rogovschi, N., Nadif, M.: A dynamic collaborative filtering system via a weighted clustering
27 approach. *Neurocomputing* **175**, 206–215 (2015). doi:[10.1016/j.neucom.2015.10.050](https://doi.org/10.1016/j.neucom.2015.10.050)
- 28 9. Aggarwal, C.C.: *Recommender Systems*. Springer, New York (2016). doi:[10.1007/978-3-319-29659-3](https://doi.org/10.1007/978-3-319-29659-3)
- 29 10. Zhang, J., Lin, Y., Lin, M., Liu, J.: An effective collaborative filtering algorithm based on user preference
30 clustering. *Applied Intelligence* **45**, 230–240 (2016). doi:[10.1007/s10489-015-0756-9](https://doi.org/10.1007/s10489-015-0756-9)
- 31 11. Laishram, A., Padmanabhan, V., Lal, R.P.: Analysis of similarity measures in user-item subgroup based
32 collaborative filtering via genetic algorithm. *International Journal of Information Technology* **10**(4), 523–527
33 (2018). doi:[10.1007/s41870-018-0195-z](https://doi.org/10.1007/s41870-018-0195-z)
- 34 12. Bagher, R., Cami Hassanpour, H., Mashayekhi, H.: User trends modeling for a content-based recommender
35 system. *Expert Systems with Applications* **87**, 209–219 (2017). doi:[10.1016/j.eswa.2017.06.020](https://doi.org/10.1016/j.eswa.2017.06.020)
- 36 13. Li, G., Zhang, Z., Wang, L., Chen, Q., Pan, J.: One-class collaborative filtering based on rating prediction and
37 ranking prediction. *Knowledge-Based Systems* **124**, 46–54 (2017). doi:[10.1016/j.knosys.2017.02.034](https://doi.org/10.1016/j.knosys.2017.02.034)
- 38 14. Wang, S., Huang, S., Liu, T.-Y., Ma, J., Chen, Z., Veijalainen, J.: Ranking-oriented collaborative filtering: A
39 listwise approach. *ACM Transactions on Information Systems* **35**(2), 1–28 (2016). doi:[10.1145/2960408](https://doi.org/10.1145/2960408)
- 40 15. Karabadjji, N.E.I., Beldjoudi, S., Seridi, H., Aridhi, S., Dhifli, W.: Improving memory-based user collaborative
41 filtering with evolutionary multi-objective optimization. *Expert Systems with Applications* **98**, 153–165 (2018).
42 doi:[10.1016/j.eswa.2018.01.015](https://doi.org/10.1016/j.eswa.2018.01.015)
- 43 16. Patra, B.K., Launonen, R., Ollikainen, V., Nandi, S.: A new similarity measure using bhattacharyya coefficient
44 for collaborative filtering in sparse data. *Knowledge-Based Systems* **82**, 163–177 (2015).
45 doi:[10.1016/j.knosys.2015.03.001](https://doi.org/10.1016/j.knosys.2015.03.001)
- 46 17. Ocepeka, U., Rugelj, J., Bosnića, Z.: Improving matrix factorization recommendations for examples in cold
start. *Expert Systems with Applications* **42**(19), 6784–6794 (2015). doi:[10.1016/j.eswa.2015.04.071](https://doi.org/10.1016/j.eswa.2015.04.071)
18. Tran, C., Kim, J.Y., Shin, W.Y., Kim, S.W.: Clustering-based collaborative filtering using an
incentivized/penalized user model. *IEEE Access* **7**, 62115–62125 (2019). doi:[10.1109/ACCESS.2019.2914556](https://doi.org/10.1109/ACCESS.2019.2914556)
19. Bobadilla, J., Bojorque, R., Esteban, A.H., Hurtado, R.: Recommender systems clustering using bayesian non
negative matrix factorization. *IEEE Access* **6**, 3549–3564 (2018). doi:[10.1109/ACCESS.2017.2788138](https://doi.org/10.1109/ACCESS.2017.2788138)
20. Vander Aa, T., Chakroun, I., Haber, T.: Distributed bayesian probabilistic matrix factorization. In: *Procedia of*
21 *International Conference on Computational Science, ICCS, 12-14 June 2017, Zurich, Switzerland*, pp.
22 1030–1039 (2017). doi:[10.1016/j.procs.2017.05.009](https://doi.org/10.1016/j.procs.2017.05.009)
23. Zhang, R., Mao, Y.: Movie recommendation via markovian factorization of matrix processes. *IEEE Access* **7**,
24 13189–13199 (2019). doi:[10.1109/ACCESS.2019.2892289](https://doi.org/10.1109/ACCESS.2019.2892289)
25. Xian, Z., Li, Q., Li, G., Li, L.: New collaborative filtering algorithms based on svd++ and differential privacy.
26 *Mathematical Problems in Engineering* **2017**, 1–14 (2017). doi:[10.1155/2017/1975719](https://doi.org/10.1155/2017/1975719)
27. Guan, X., Li, C.T., Guan, Y.: Matrix factorization with rating completion: An enhanced svd model for
28 collaborative filtering recommender systems. *IEEE Access* **5**, 27668–27678 (2017).
29 doi:[10.1109/ACCESS.2017.2772226](https://doi.org/10.1109/ACCESS.2017.2772226)
30. Yue, L., Sun, X.X., Gao, W.Z., Feng, G.Z., Zhang, B.Z.: Multiple auxiliary information based deep model for
31 collaborative filtering. *Journal of Computer Science and Technology* **33**(4), 668–681 (2018).
32 doi:[10.1007/s11390-018-1848-x](https://doi.org/10.1007/s11390-018-1848-x)
33. Shams, B., Haratizadeh, S.: Item-based collaborative ranking. *Knowledge-Based Systems journal* **152**,
34 172–185 (2018). doi:[10.1016/j.knosys.2018.04.012](https://doi.org/10.1016/j.knosys.2018.04.012)
34. Park, Y., Park, S., Jung, W., Lee, S.G.: Reversed cf: A fast collaborative filtering algorithm using a k-nearest
35 neighbor graph. *Expert Systems with Applications* **42**(8), 4022–4028 (2015). doi:[10.1016/j.eswa.2015.01.001](https://doi.org/10.1016/j.eswa.2015.01.001)
36. Polatidis, N., Georgiadis, C.K.: A multi-level collaborative filtering method that improves recommendations.
37 *Expert Systems with Applications* **48**, 100–110 (2016). doi:[10.1016/j.eswa.2015.11.023](https://doi.org/10.1016/j.eswa.2015.11.023)

- 1 28. Sun, S.B., Zhang, Z.H., Dong, X.L., Zhang, H.R., Li, T.J., Zhang, L., Min, F.: Integrating triangle and jaccard 1
2 similarities for recommendation. *Plos One* **12**(8), 1–11 (2017). doi:[10.1371/journal.pone.0183570](https://doi.org/10.1371/journal.pone.0183570) 2
- 3 29. Wu, C., Wu, J., Luo, C., Wu, Q., Liu, C., Wu, Y., Yang, F.: Recommendation algorithm based on user score 3
4 probability and project type. *Eurasip Journal on Wireless Communications and Networking* **2019**(80), 1–13 4
5 (2019). doi:[10.1186/s13638-019-1385-5](https://doi.org/10.1186/s13638-019-1385-5) 5
- 6 30. Yu, P.: Collaborative filtering recommendation algorithm based on both user and item. In: Proceedings of 2015 6
7 4th International Conference on Computer Science and Network Technology, ICCSNT 2015, pp. 239–243 7
8 (2015). doi:[10.1109/ICCSNT.2015.7490744](https://doi.org/10.1109/ICCSNT.2015.7490744) 8
- 9 31. Harper, F.M., Konstan, J.A.: The movielens datasets: History and context. *ACM Transactions on Interactive 9
10 Intelligent Systems* **5**(4) (2015). doi:[10.1145/2827872](https://doi.org/10.1145/2827872) 10
- 11 32. Jalili, M., Ahmadian, S., Izadi, M., Moradi, P., Salehi, M.: Evaluating collaborative filtering recommender 11
12 algorithms: A survey. *IEEE Access* **6**, 74003–74024 (2018). doi:[10.1109/ACCESS.2018.2883742](https://doi.org/10.1109/ACCESS.2018.2883742) 12
- 13 33. Zhang, F., Gong, T., Lee, V.E., Zhao, G., Rong, C., Qu, G.: Fast algorithms to evaluate collaborative filtering 13
14 recommender systems. *Knowledge-Based Systems* **96**, 96–103 (2016). doi:[10.1016/j.knosys.2015.12.025](https://doi.org/10.1016/j.knosys.2015.12.025) 14
- 15 34. Zheng, M., Min, F., Zhang, H.R., Chen, W.B.: Fast recommendations with the m-distance. *IEEE Access* **4**, 15
16 1464–1468 (2016). doi:[10.1109/ACCESS.2016.2549182](https://doi.org/10.1109/ACCESS.2016.2549182) 16
- 17 35. Vellaichamy, V., Kalimuthu, V.: Hybrid collaborative movie recommender system using clustering and bat 17
18 optimization. *International Journal of Intelligent Engineering and Systems* **10**(5), 38–47 (2017). 18
19 doi:[10.22266/ijies2017.1031.05](https://doi.org/10.22266/ijies2017.1031.05) 19
- 20 36. Fan, X., Chen, Z., Zhu, L., Liao, Z., Fu, B.: A novel hybrid similarity calculation model. *Scientific Programming* 20
21 **2017**, 1–9 (2017). doi:[10.1155/2017/4379141](https://doi.org/10.1155/2017/4379141) 21
22 22
23 23
24 24
25 25
26 26
27 27
28 28
29 29
30 30
31 31
32 32
33 33
34 34
35 35
36 36
37 37
38 38
39 39
40 40
41 41
42 42
43 43
44 44
45 45
46 46

Figures

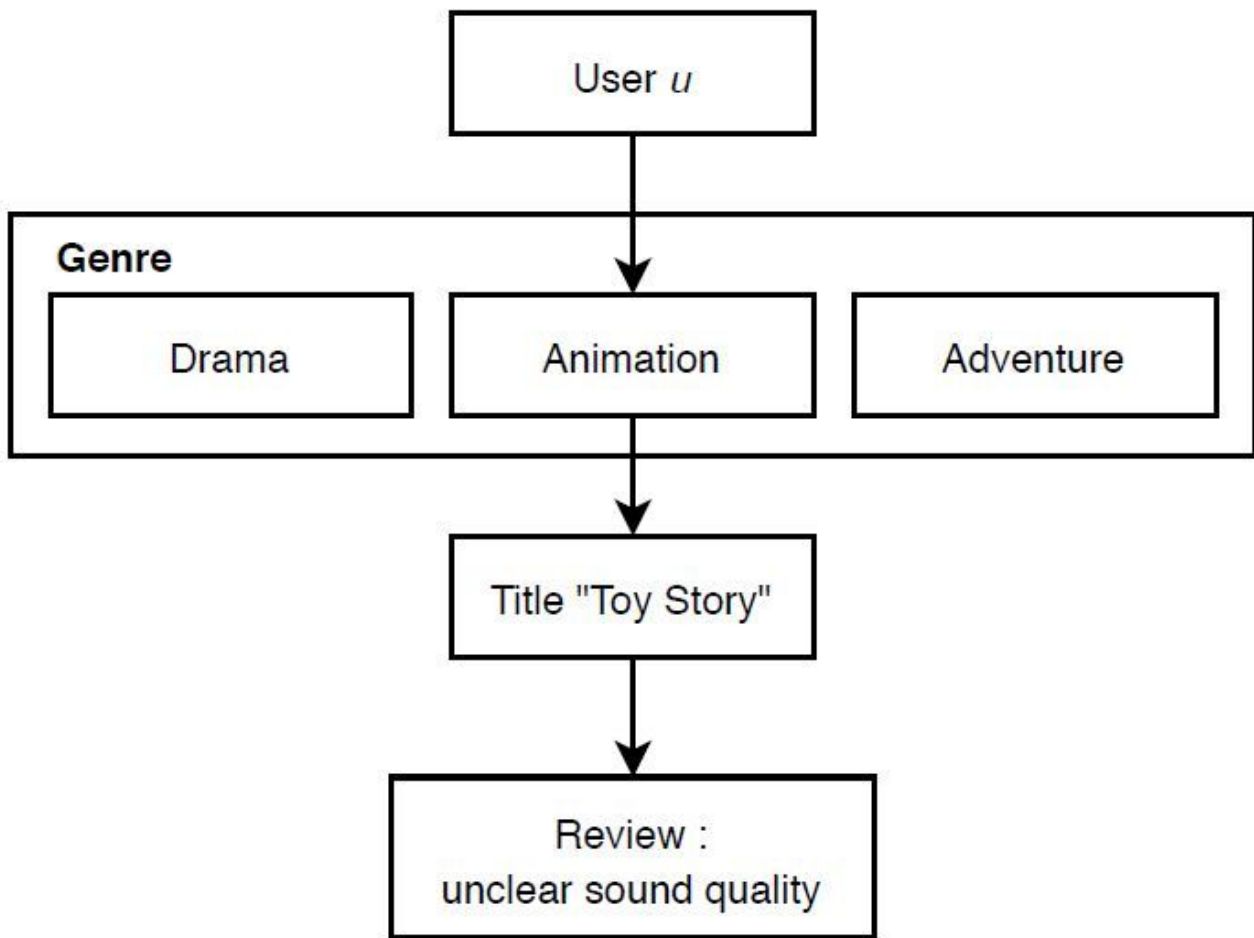


Figure 1

Example of user behavior in rating the item

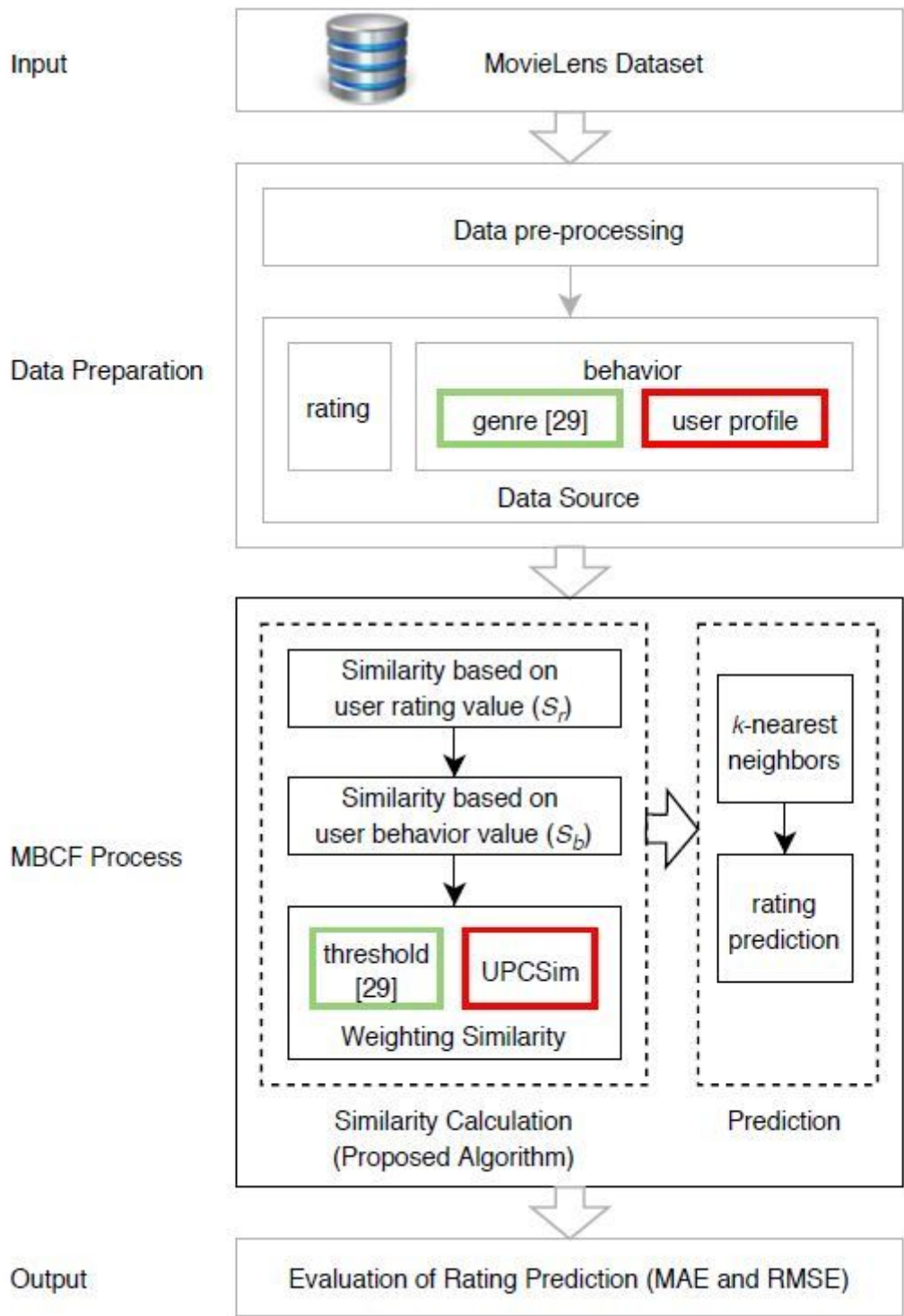


Figure 2

The developed MBCF system

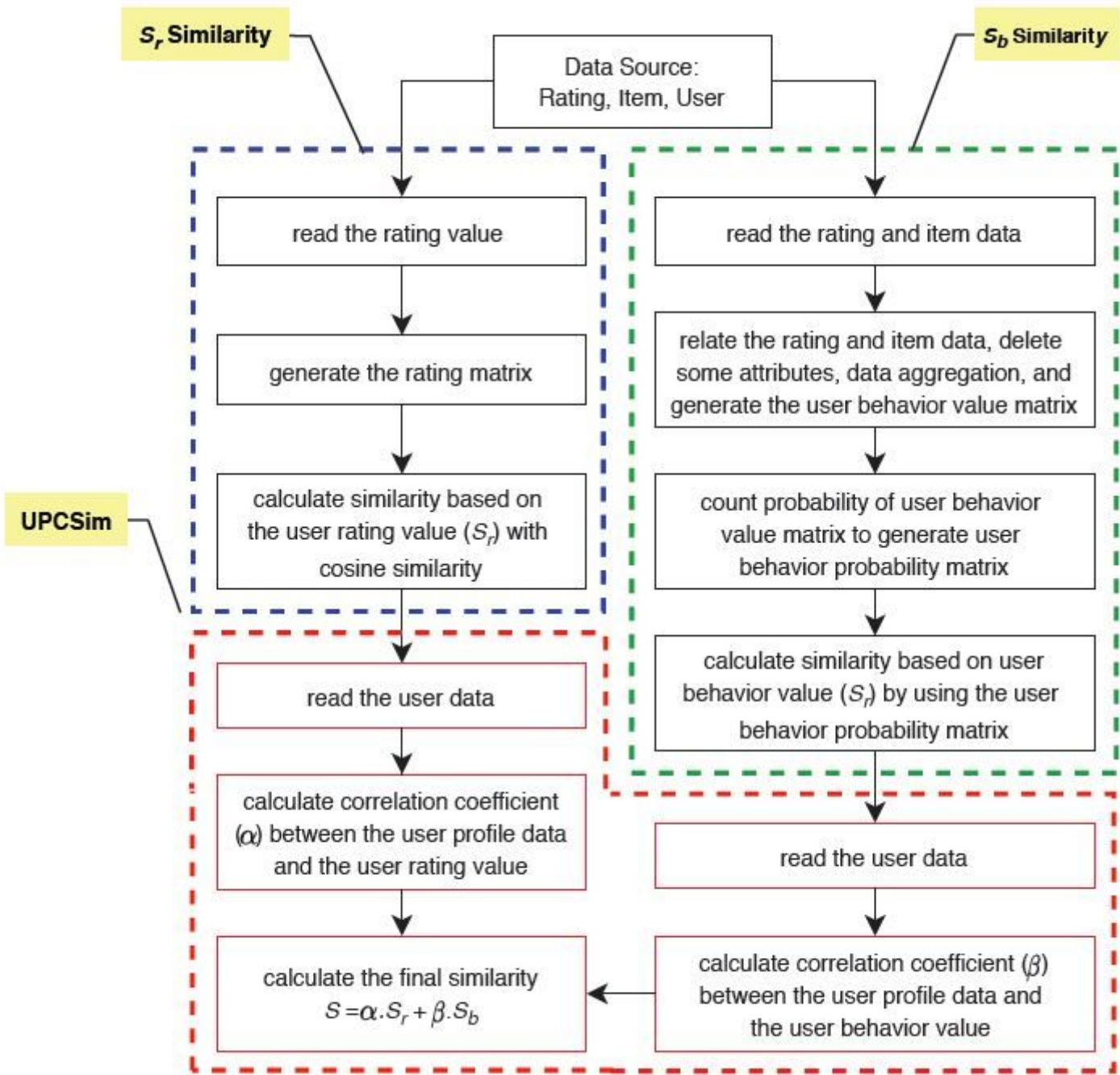


Figure 3

Flowchart of similarity calculation

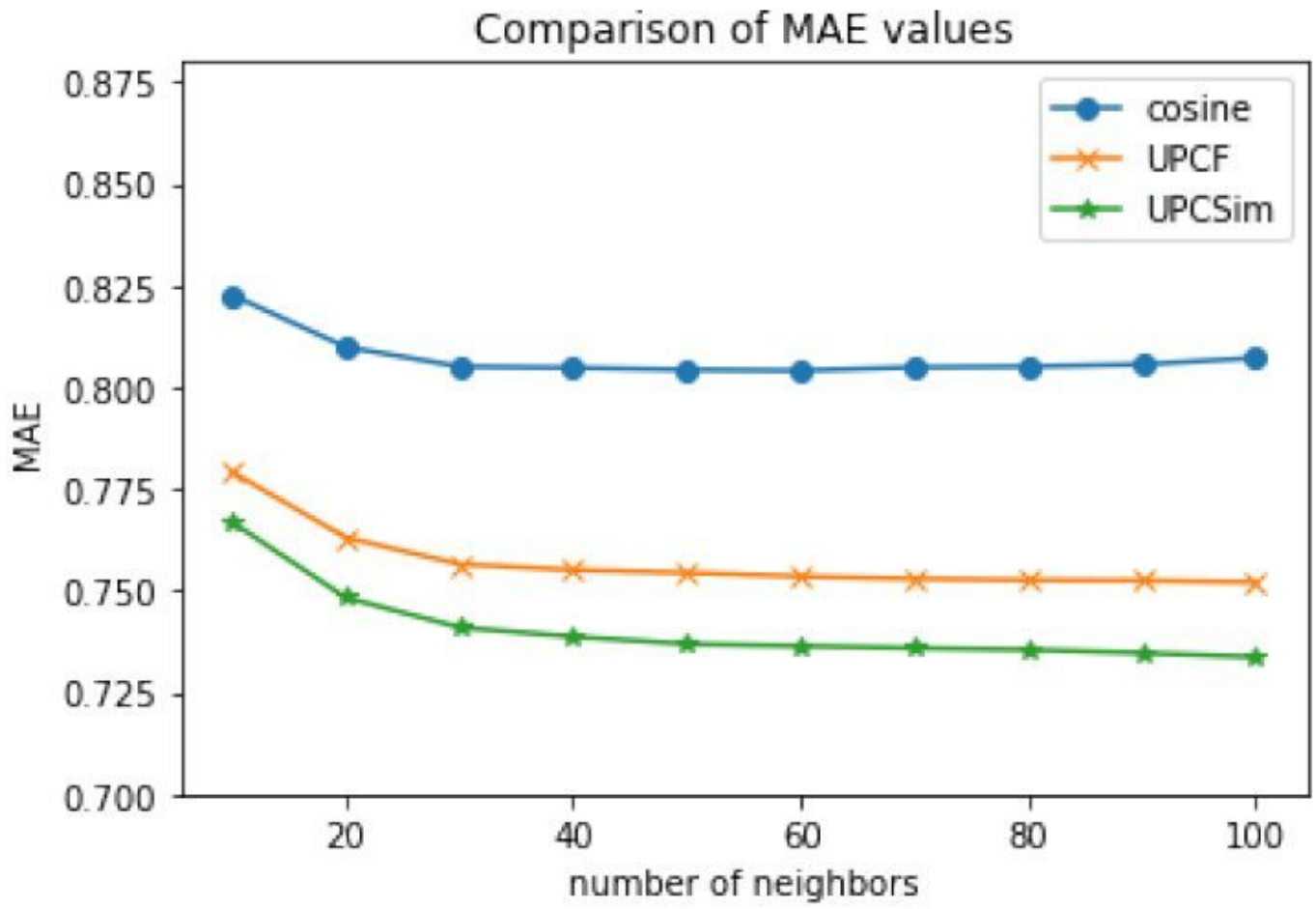


Figure 4

Comparison of the average MAE value of the three similarity algorithms using Movielens 100k dataset

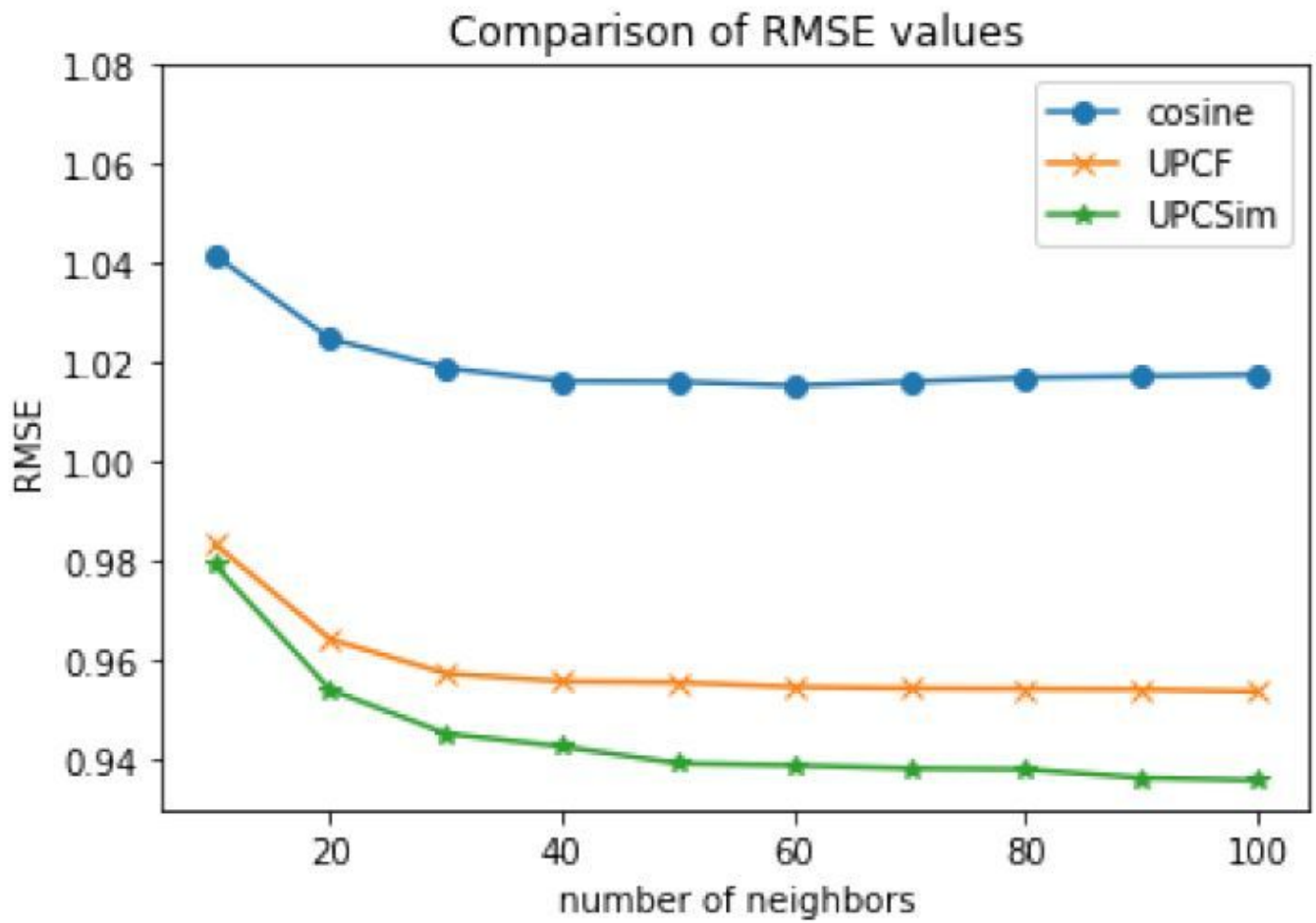


Figure 5

Comparison of the average RMSE value of the three similarity algorithms using Movielens 100k dataset