

Research on Parameter Optimization in Collaborative Filtering Algorithm

Zijiang Zhu

South China Business College, Guangdong University of Foreign Studies, Guangzhou, China

Email: zzjdwh2002@163.com

How to cite this paper: Zhu, Z.J. (2018) Research on Parameter Optimization in Collaborative Filtering Algorithm. *Communications and Network*, 10, 105-116. <https://doi.org/10.4236/cn.2018.103009>

Received: March 10, 2018

Accepted: August 24, 2018

Published: August 27, 2018

Copyright © 2018 by author and Scientific Research Publishing Inc.

This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).

<http://creativecommons.org/licenses/by/4.0/>



Open Access

Abstract

Collaborative filtering algorithm is the most widely used and recommended algorithm in major e-commerce recommendation systems nowadays. Concerning the problems such as poor adaptability and cold start of traditional collaborative filtering algorithms, this paper is going to come up with improvements and construct a hybrid collaborative filtering algorithm model which will possess excellent scalability. Meanwhile, this paper will also optimize the process based on the parameter selection of genetic algorithm and demonstrate its pseudocode reference so as to provide new ideas and methods for the study of parameter combination optimization in hybrid collaborative filtering algorithm.

Keywords

Collaborative Filtering Algorithm, Genetic Algorithm, Parameter Combination Optimization

1. Introduction

User-based collaborative filtering and item-based collaborative filtering are deemed as two classic methods in collaborative filtering algorithms while traditional collaborative filtering algorithms have problems such as cold start and data sparsity. A single filtering algorithm cannot give the best recommendation. According to the advantages of parallel and distributed cloud platform, multiple algorithms can be executed simultaneously based on applying multiple nodes [1]. What's more, user-item-based collaborative filtering can be implemented on the cloud platform while the advantages of its parallelization can be adopted to process big data efficiently [2]. On the basis of this, by optimizing the nearest neighbor's selection strategy and scoring strategy, the recommendation system's

effect and execution efficiency can be greatly improved. Therefore, the running of algorithm fusion research is of pivotal significance in both theory and reality [3].

A well-established recommendation system is usually composed of three parts: the recording plate, recording the user's historical behavior; analysis plate; algorithm plate. As the core of the whole recommendation system, the related research of the recommendation algorithm has become the hot research direction, because the recommendation result is closely related to the performance of the recommended algorithm. According to different needs of users, recommendations will differ from each other [4]. At present, there here are four main kinds of recommendation algorithms: content-based recommendation algorithms, collaborative filtering recommendation algorithm, hybrid recommendation algorithm and network recommendation algorithm. Collaborative filtering recommendation technology is one of the successful technologies in the application of personalized recommendation technology [5]. Thanks to collective wisdom, it will unearth out a small part of "neighbors" with similar hobbies among a large number of users, according to the analysis and records of these neighbors' favorite contents, and generate a sorted catalog, which is called recommended results, then push it to the "neighbors", which will reduce the user's workload of picking process. The traditional collaborative filtering algorithm performs the similarity calculations with the user scores without considering the behavior time or the same label between the items, as a result, a lot of failures are exposed, such as the cold start problem, sparse matrix problems, recommended accuracy issues, etc. The accuracy of recommendation results is not satisfying, and will hardly meet the actual needs of users [6]. My research is based on the item-based collaborative filtering algorithm, and adds the users' behavior time and the label attributes also other information into the similarity calculation, to avoid the cold start problem, thereby improving the quality of the recommend results. By adjusting the parameters of the collaborative filtering algorithm to meet the actual needs of different users as much as possible, the personalized recommendation service is realized.

2. Thoughts to Improve Collaborative Filtering Algorithm

However, the recommended algorithm fusion based on the cloud platform has shortages in terms of parameter selection and combination optimization. Among the proposed algorithms, whether the K-Nearest Neighbor (KNN) [7] is used to generate the nearest neighbor set or the threshold method is used to generate the nearest neighbor set, the optimal parameter selection cannot be avoided. This paper aims to solve the cold start problem of the traditional collaborative filtering algorithm, and further alleviate the data sparsity. In addition, it fully analyzes and takes advantage of user and item characteristics, and combines people's craziness for products and their similarity based on weight method while the weight directly determines the proportion of the nearest neighbor

that will be affected. Those two are complementary and then the recommendation of user-to user based on nearest neighbor effect will be generated. At the same time, the similarity between items and characteristics will also be combined based on weight method and the weight directly determines the proportion of the nearest neighbor that will be affected. Those two are complementary and then the recommendation of item-to-item based on nearest neighbor effect will be generated [8]. Finally, the two recommendations are combined by using appropriate weights and the final recommendation will be generated. The flow chart of this hybrid collaborative filtering algorithm is shown in **Figure 1**.

In the improved hybrid collaborative filtering algorithm described above, multiple weights and thresholds are used. If the exhaustive method is used to select the optimal combination, the efficiency of the algorithm will be greatly affected. This paper will use the Genetic Algorithm [9] to search for the optimal parameter combinations that meet the conditions so as to improve the accuracy of the recommendation algorithm, and take advantage of the distributed cloud platform. Hybrid collaborative filtering combines multiple algorithms to ensure execution efficiency which can significantly improve the recommended effect.

3. Hybrid Collaborative Filtering Algorithm

3.1. Description of User and Item Characteristics

In order to perform mathematical operations on user features and item features,

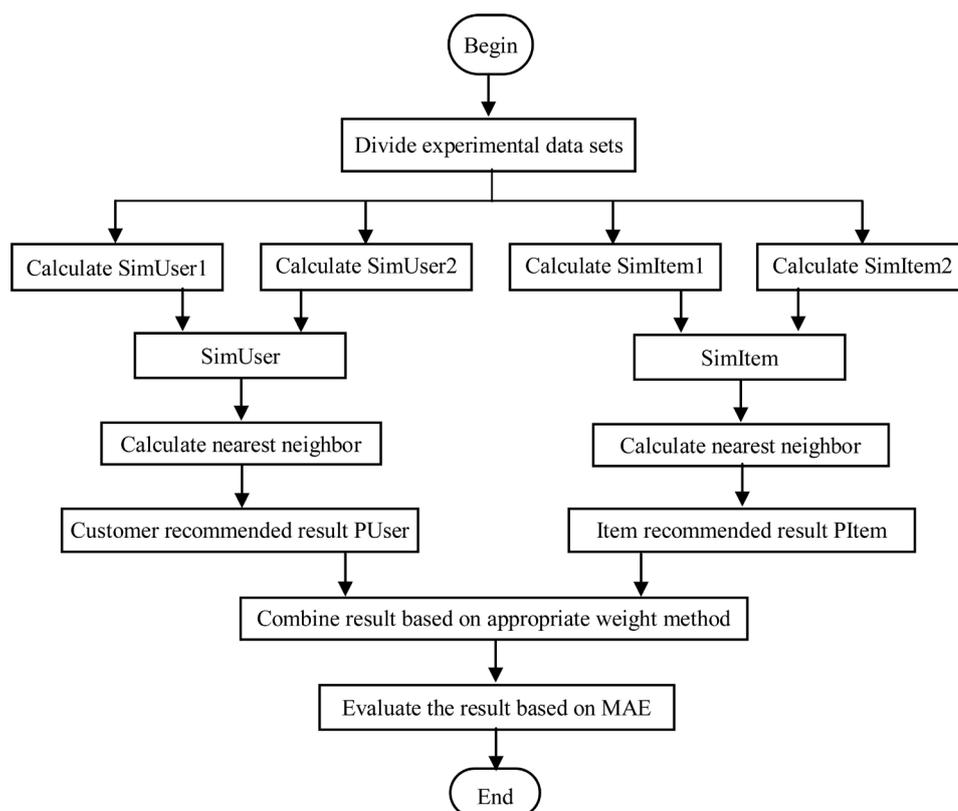


Figure 1. Flow chart of hybrid collaborative filtering algorithm.

encoding is necessary. This paper selects the MovieLens dataset [10] to evaluate the recommended results of the improved algorithm. The MovieLens dataset provides historical information on more than 1600 movies from more than 900 users. In this set of data, the user's rating for the movie is divided into five levels, with 1 being the worst and 5 being the best. The MovieLens dataset also has detailed user and movie feature information. This article has coded the following information for these feature messages. The original user feature data of the MovieLens data set is described as follows:

The user's age is divided into: 0 - 16 for children, 17 - 39 for young people, 40 - 60 for middle-aged people, 60 or older for the elderly, and the integers for these four types of people are respectively demonstrated as 1, 2, 3, 4 code representation. The user gender is coded as 1 for men and 0 for women. In the MovieLens dataset, the movie project attributes include 19 movie tags such as romance films and suspense films. The movie can have multiple types and the encoding rule of this article is: the target movie belongs to the current type, and is represented by 1, and if it is not, it is represented by 0.

In this paper, the Mean Absolute Error (MAE) [11] index is used to evaluate the accuracy of the improved hybrid collaborative filtering algorithm and the traditional algorithm. The MovieLens data set is divided into two groups of different data, namely the training data set and Test the data set. The training data set is used to determine the recommended model and related parameters, and the test data set is used to verify the accuracy of the recommended results.

3.2. Establishment of User Scoring Matrix Data

In order to establish a scoring matrix suitable for the collaborative filtering algorithm, the initial data needs to be pre-processed, and the pre-processing results are shown in **Table 1**.

In **Table 1**, U_i represents the i -th user, I_j represents the j -th item, and r_{ij} represents the score of the i -th user on the j -th movie item. From this, a user-item scoring matrix $R = (r_{ij})$ can be constructed, where $i = [1, 2, \dots, m]$ and $j = [1, 2, \dots, n]$.

Table 1. User-Item rating data pre-processing table.

User	Item					
	I_1	I_2	...	I_j	...	I_n
U_1	r_{11}	r_{12}	...	r_{1j}	...	r_{1n}
U_2	r_{21}	r_{22}	...	r_{2j}	...	r_{2n}
...
U_j	r_{j1}	r_{j2}	...	r_{jj}	...	r_{jn}
...
U_m	r_{m1}	r_{m2}	...	r_{mj}	...	r_{mn}

3.3. Similarity Calculation

In the improved hybrid collaborative filtering algorithm, user score similarity and item score similarity are calculated by Pearson Correlation Coefficient (PCC) [12].

1) User Rating Similarity.

The similarity between user i and user j is as shown in Equation (1).

$$simUser1(i, j) = \frac{\sum_{s \in S} (r_{i,s} - \bar{r}_i)(r_{j,s} - \bar{r}_j)}{\sqrt{\sum_{s \in S} (r_{i,s} - \bar{r}_i)^2} \sqrt{\sum_{s \in S} (r_{j,s} - \bar{r}_j)^2}} \quad (1)$$

In formula (1), the \bar{r}_i and \bar{r}_j show the mean core of user i and user j , while $r_{i,s}$ and $r_{j,s}$ show the scores of the common item s from user i and user j .

The similarity of the user characteristic information is calculated by the formula (2) and the formula (3).

$$Dis(i, j) = \sqrt{\sum_{k=1}^p (i_k - j_k)^2} \quad (2)$$

$$simUser2(i, j) = \frac{1}{1 + Dis(i, j)} \quad (3)$$

In Equation (2), p represents the number of user characteristic attributes, i_k represents the value of the k -th feature of the i -th user, j_k represents the value of the k -th feature of the j -th user. Then, we calculate the value according to the characteristics of the user i and the user j . The Euclidean Metric (EM) [13] between users, and then the similarity of the characteristic information of user i and user j is calculated by formula (3).

Finally, the above two similarities are mixed together by the weight combination method. See Equation (4), and the comprehensive similarity between user i and user j is obtained.

$$simUser(i, j) = w_1 \times simUser1(i, j) + (1 - w_1) \times simUser2(i, j) \quad (4)$$

In the formula (4), w_1 is the weight.

2) Item Score Similarity.

The similarity calculation of the user i scores the item m and the item n is as shown in the formula (5).

$$simItem1(m, n) = \frac{\sum_{i \in I} (r_{i,m} - \bar{r}_m)(r_{i,n} - \bar{r}_n)}{\sqrt{\sum_{i \in I} (r_{i,m} - \bar{r}_m)^2} \sqrt{\sum_{i \in I} (r_{i,n} - \bar{r}_n)^2}} \quad (5)$$

In the formula (5), $r_{i,m}$ and $r_{i,n}$ respectively refer to the scores of the user i for the item m and the item n , \bar{r}_m and \bar{r}_n refer to the average scores of the items m and n , and I refers to the set of users who have evaluated both the item m and the item n .

The similarity of the item characteristic information is calculated by the formula (6) and the formula (7).

$$Dis(m, n) = \sqrt{\sum_{k=1}^q (m_k - n_k)^2} \quad (6)$$

$$simItem2(m, n) = \frac{1}{1 + Dis(m, n)} \quad (7)$$

In formula (6), q represents the number of item feature information, m_k represents the code value of the k -th feature information of the m -th item, and n_k represents the code value of the k -th feature information of the n -th item. According to the encoded value of the characteristic information of given item m and item n , the Euclidean distance between them is calculated, and finally the characteristic information similarity of the item m and the item n is calculated according to the formula (7).

Finally, the above two similarities are mixed together by the weight combination method. See Equation (8), and the comprehensive similarity between the project m and the project n is obtained.

$$simItem(m, n) = w_2 \times simItem1(m, n) + (1 - w_2) \times simItem2(m, n) \quad (8)$$

In formula (8), C is the weight.

3) Selection of Nearest Neighbor

According to the analysis of the advantages and disadvantages of the nearest neighbor selection method, the threshold selection method is adopted as the nearest neighbor selection method. Set w_3 as the user similarity threshold, when the similarity between the user i and the user j is greater than or equal to w_3 , that is $simUser(i, j) \geq w_3$, and the user j is a member of the nearest neighbor group of the user i . Similarly, the item similarity threshold is set, when the similarity between the item m and the item n is greater than or equal to w_4 , that is $simItem(m, n) \geq w_4$, then the item n is a member of the nearest neighbor group of the item m .

3.4. Method of Selecting Nearest Neighbors

There are usually two methods for selecting nearest neighbors: one is the K-Nearest Neighbor (KNN), and the K users with the biggest similarity are selected. This method mainly performs sorting in descending order according to the calculated user similarity, and selects the user with the K closest neighbors. The second is the threshold selection method, and the user whose similarity is greater than the threshold is selected as the nearest neighbor. In this method, as long as a similarity threshold is preset, the user whose similarity is bigger than the threshold is used as the nearest neighbor set of the target user [14].

By comparing the two selection methods of nearest neighbors, the K-nearest neighbor method is simple and easy, but it has the following disadvantages: the number of nearest neighbors is artificially specified, and the K nearest neighbors of the target user do not necessarily belong to the nearest neighbor of the target user; in addition, it is difficult to determine the value of K, if the K value is too large, it may lead to an extensive coverage and the nearest neighbor selection will not be accurate and if the K value is too small there will be less targets, resulting in lower accuracy.

By using the threshold selection method, we can perfectly avoid the defect of KNN and this paper uses this method by improving collaborative filtering algorithm. However, the threshold selection method also has defects because it is difficult to ensure the threshold value. In this paper, the genetic algorithm is adopted to optimize the threshold selection while using collaborative so as to improve the accuracy of the algorithm.

3.5. Recommended Result

After a series of calculations of the nearest neighbor group of the user, the score of the target user for other unevaluated items is predicted by the recommendation formula according to the result. The definition of the user recommendation formula is shown in formula (9).

$$PItem(t, u) = \bar{r}_t + \frac{\sum_{k \in N_t} (r_{u,t} - \bar{r}_k) \times simItem(t, k)}{\sum_{k \in N_t} simItem(t, k)} \quad (9)$$

In formula (9), \bar{r}_t and \bar{r}_k mean the average score of item t and item k , N_t refers to the nearest neighbor set of item t , $r_{u,t}$ indicates the score of user u on item t .

Similarly, the weighting method can be used to mix the two results, see Equation (10).

$$P(t, u) = w_5 \times PUser(t, u) + (1 - w_5) \times PItem(t, u) \quad (10)$$

In formula (10), w_5 is weight value $PUser(t, u)$ is the recommendation result based on user-user relationship while $PItem(t, u)$ shows the recommendation result based on item-item relationship and $P(t, u)$ represents the final mixed score result of the user u on the item t .

4. Genetic Algorithm Optimizing Parameter Selection

4.1. Hybrid Recommendation Algorithm

The hybrid recommendation algorithm is a collaborative filtering algorithm that combines the user's historical rating information, feature information, and item score information by weight method. The basic idea of the algorithm is to combine the user score similarity with the user characteristic information similarity by weight method to obtain the nearest neighbor group of the user, so as to obtain the recommendation result based on the user-user relationship; on the other hand, the weight method is also used to combine the item score similarity with the item feature information so as to obtain the nearest neighbor, so as to obtain the recommendation result based on the item-item relationship. Finally, the combination between user recommendation result and item recommendation result is carried out to generate a final result.

4.2. Parameter Selection of Hybrid Recommendation Algorithm

The hybrid collaborative filtering recommendation algorithm described above

deals with the cold start and data sparse problems, but brings about the selection of threshold and weight. On the one hand, it is necessary to select the appropriate weight w_1 , w_2 , w_3 , and on the other hand, in the calculation of the user's nearest neighbor group and the nearest neighbor group of the item, there is also a similarity threshold value w_3 , w_4 selection problem. The value range of these parameters is $[0, 1]$ and if we use the exhaustive method to select the best combination of parameters, even if we choose the combination of the amplitude of 0.1 parameters, there are 11 to the 5th choices. Assuming that it is executed once and takes 10 seconds, and it will take more than half a month to find the optimal value selection through an exhaustive method, which will greatly reduce the efficiency of the algorithm [15]. This paper solves this problem through genetic algorithms.

4.3. Parameter Selection Optimization Process Based on Genetic Algorithm

The genetic algorithm is used to select the optimal combination of five parameters w_1 , w_2 , w_3 , w_4 , w_5 , so that the recommendation algorithm and the average absolute deviation MAE of the test data set tend to be the smallest [16]. These five parameters are unified chromosome coding, as shown in **Table 2**.

In genetic algorithm optimization, the fitness is defined as an improved recommendation algorithm based on the combination of parameters, and the reciprocal of the average absolute deviation MAE of the recommendation result and the test data set, *i.e.* $1/MAE$. It can be seen that the smaller the MAE value of the recommended result is, the greater the fitness is, indicating that the combination of the parameters is better.

The flow of parameter selection optimization of the hybrid collaborative filtering algorithm based on genetic algorithm is shown in **Figure 2**. It can be seen that the genetic algorithm translates the weight and threshold information stored

Table 2. Parameter value chromosome coding.

Numerical Value	Coding
0	0
0.1	1
0.2	2
0.3	3
0.4	4
0.5	5
0.6	6
0.7	7
0.8	8
0.9	9
1	*

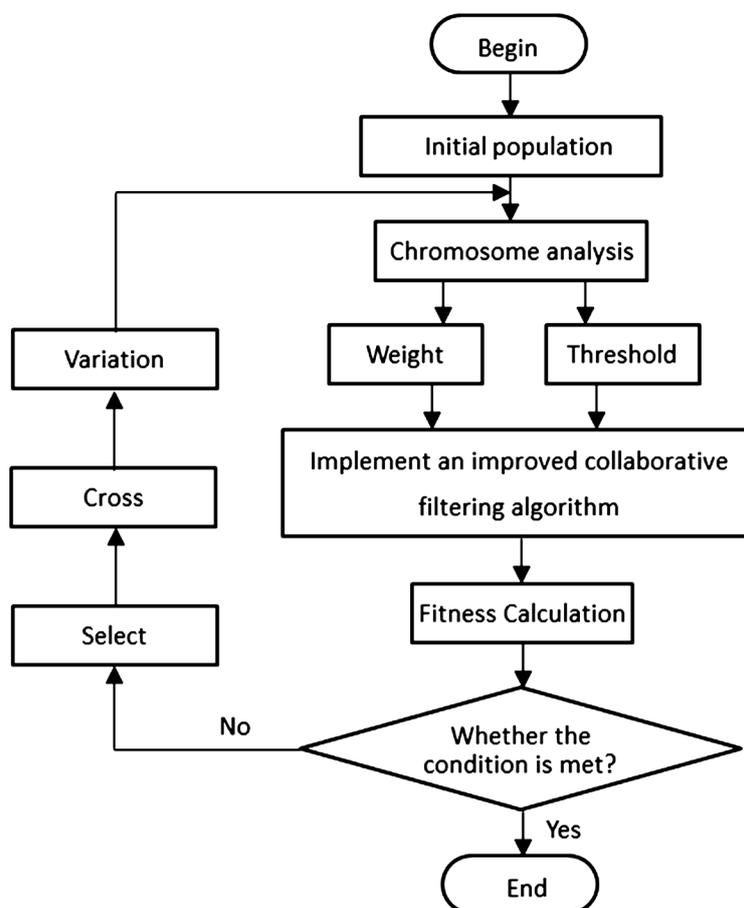


Figure 2. Flow chart of genetic algorithm optimizing algorithm.

in the chromosome, and then substitutes them into the hybrid collaborative filtering algorithm described above until the end condition of the algorithm is satisfied, and the next generation inheritance will be implemented if it is not satisfied operating. The improved algorithm uses weights and thresholds in five places, and the values of these five parameters are all $[0, 1]$. By using genetic algorithm, the parameter combination of the recommendation algorithm will be improved and the pseudo code of the main function optimization is shown in **Table 3**.

5. Conclusions

In this paper, the traditional collaborative filtering algorithm is improved, and the user and item scores are combined with the feature that attributes to generate recommendations. The cold start problem in the traditional collaborative filtering algorithm is solved, and the sparseness of user rating data is alleviated to some extent. The paper also demonstrates the parameter combination optimization of improved collaborative filtering algorithm and introduces the algorithm process combining improved algorithm as well as genetic algorithm. What's more, the flow and pseudo code of the combined algorithm is also given to solve the parameter combination optimization issue of collaborative filtering algorithm

Table 3. Genetic algorithm optimizing pseudo code.

Function name: main GA
Function: Optimized parameter combination
<pre> BEGIN % the ipop array variable represents a population and consists of multiple chromosome individuals and is used to store the chromosomes of the i-th generation population. % the evals array stores the fitness values of the chromosomes at the corresponding locations in the ipop array. % uses the inialPops function to create the initial population. ipop = inialPops(); % implements genetic manipulation cyclically with a minimum of 50 generations of genetic ma-nipulation. for i = 1 to 50 for k = 1 to m % m is the number of groups evals[k] = calculatefitnessvalue(ipop[k]); % calculates the fitness of the corresponding chromosome. if(evals[k] > best.fitness) { best.fitness = evals[k]; % records the most adaptive value with the corresponding chromosome. best.pop = ipop[k]; } end if end if(best.fitness >= n) { % n is the preset MAE threshold. select(); % select operation. cross(); % cross operation. matation(); % variation operation. } else stop; % reaches the threshold requirement and stops executing the genetic algorithm. end if end END </pre>

while genetic algorithm is introduced.

The recommendation algorithm plays a pivotal role in the development of ecommerce. As the number of user is soaring, the user model becomes more and more complex, and the recommendation effect of a single recommendation algorithm becomes worse and worse. The hybrid collaborative filtering algorithm based on genetic algorithm proposed in this paper combines multiple recommendation algorithms and can process a large amount of user data with good scalability, and can achieve better recommendation results.

Acknowledgements

This work was supported in part by a grant from the characteristics innovation project of colleges and universities of Guangdong Province (Natural Science, No. 2016KTSCX182, 2016), a grant from the Youth Innovation Talent Project of colleges and universities of Guangdong Province (No. 2016KQNCX230, 2016).

Conflicts of Interest

The authors declare no conflicts of interest regarding the publication of this paper.

References

- [1] Zhang, J.W. and Yang, Z. (2014) Collaborative Filtering Recommendation Algorithm Based on Improved User Clustering. *Computer Science*, **41**, 176-178.

- <https://doi.org/10.11896/j.issn.1002-137X.2014.12.038>
- [2] Xiao, Q., Zhu, Q.H., Zheng, H. and Wu, K.W. (2013) Design and Implementation of Distributed Collaborative Filtering Algorithm on Hadoop. *Data Analysis and Knowledge Discovery*, No. 1, 83-89.
- [3] Zhang, L., Teng, P.Q. and Qin, T. (2014) Using Key Users of Social Network to Enhance Collaborative Filtering Performance. *Journal of Intelligence*, **33**, 196-200.
- [4] Wen, S.Q., Wang, C., Su, F.F., Liu, J.F., Chen, Y.W. and Zheng, G.Q. (2017) Using Users' Unfavorable Item Attributes to Improve the Efficiency and Accuracy of Item-Based Collaborative Filtering Algorithm. *Journal of Chinese Computer Systems*, **38**, 1735-1740.
- [5] Wang, J.H. and Han, J.T. (2017) Collaborative Filtering Algorithm Based on Item Attribute Preference. *Computer Engineering and Applications*, **53**, 106-110.
- [6] Zhu, B. (2017) The Improvement and Empirical Analysis of the Collaborative Filtering Algorithm about the Recommendation System of Digital Library. *Library and Information Service*, **61**, 130-134.
<https://doi.org/10.13266/j.issn.0252-3116.2017.09.017>
- [7] Zhong, D.H., Du, R.X., Cui, B., Wu, B.P. and Guan, T. (2018) Real-Time Spreading Thickness Monitoring of High-Core Rockfill Dam Based on K -Nearest Neighbor Algorithm. *Transactions of Tianjin University*, **24**, 282-289.
<https://doi.org/10.1007/s12209-017-0115-5>
- [8] Chen, A.P. and Wang, S. (2014) A Hybrid Collaborative Filtering Algorithm Based on User-Item. *Computer Technology and Development*, **24**, 88-91.
- [9] Stacewicz, P. (2015) Evolutionary Schema of Modeling Based on Genetic Algorithms. *Studies in Logic. Grammar and Rhetoric*, **40**, 219-239.
<https://content.sciendo.com/view/journals/slgr/40/1/article-p219.xml>
<https://doi.org/10.1515/slgr-2015-0011>
- [10] Kuzelewska, U. (2014) Clustering Algorithms in Hybrid Recommender System on Movie Lens Data. *Studies in Logic, Grammar and Rhetoric*, **37**, 125-139.
<https://doi.org/10.2478/slgr-2014-0021>
- [11] Wang, W.J. and Lu, Y.M. (2018) Analysis of the Mean Absolute Error (MAE) and the Root Mean Square Error (RMSE) in Assessing Rounding Model. *IOP Conference Series. Materials Science and Engineering*, No. 1, 1-10.
<http://iopscience.iop.org/article/10.1088/1757-899X/324/1/012049/meta>
<https://doi.org/10.1088/1757-899X/324/1/012049>
- [12] Jafari, M., Ghavami, B. and Sattari, V. (2017) A Hybrid Framework for Reverse Engineering of Robust Gene Regulatory Networks. *Artificial Intelligence in Medicine*, No. 6, 15-27. <https://linkinghub.elsevier.com/retrieve/pii/S0933365716304882>
<https://doi.org/10.1016/j.artmed.2017.05.004>
- [13] Liu, W.F. and Li, B.L. (2016) Projectively Flat Finsler Metrics Defined by the Euclidean Metric and Related 1-Forms. *Differential Geometry and Its Applications*, No. 6, 14-24. <https://linkinghub.elsevier.com/retrieve/pii/S0926224516300109>
<https://doi.org/10.1016/j.difgeo.2016.01.007>
- [14] Chiara, M. and Silvia, S. (2017) Reliability of TMS Phosphene Threshold Estimation: Toward a Standardized Protocol. *Brain stimulation*, **10**, 609-617.
[https://www.brainstimjrn.com/article/S1935-861X\(17\)30605-8/fulltext](https://www.brainstimjrn.com/article/S1935-861X(17)30605-8/fulltext)
- [15] Zhang, X. and Zhang, X.U. (2017) Circular Antenna Design by Adaptive Position Inheritance Artificial Bee Colony Algorithm. *Physical Communication*, **25**, 369-375.
<https://www.sciencedirect.com/science/article/pii/S1874490716302646>

<https://doi.org/10.1016/j.phycom.2017.06.004>

- [16] Li, D.J., Li, Y.Y., Li, J.X. and Fu, Y. (2018) Gesture Recognition Based on BP Neural Network Improved by Chaotic Genetic Algorithm. *International Journal of Automation and Computing*, **15**, 267-276.

<https://link.springer.com/article/10.1007%2Fs11633-017-1107-6>