

LEARNING PATH PLANNING ALGORITHM BASED ON KL DIVERGENCE AND D-VALUE MATRIX SIMILARITY

ZHAOYU SHOU, XIANYING LU, ZHENGZHENG WU, JUNLI LAI AND PAN CHEN

School of Information and Communication
Guilin University of Electronic Technology
No. 1, Jinji Road, Guilin 541004, P. R. China
guilinshou@guet.edu.cn

Received May 2020; accepted August 2020

ABSTRACT. *Aiming at the problem that the existing learning path planning algorithm fails to consider to which degree the online learner has mastered the knowledge points, a learning path planning algorithm based on KL divergence and D-value matrix similarity is proposed. The algorithm, based on the learner's online learning behavior data set, first establishes the conceptual interaction achievement model of knowledge points and the directed learning path network, and proposes a local structure similarity measurement method between the knowledge nodes of the directed learning path network. Second, based on the learner's KL divergence matrix, a learning behavior similarity calculation method on the basis of D-value matrix similarity is proposed, which is used to perform cluster analysis on learners with similar learning behaviors and to analyze the personalized optimal learning path of each kind of learners. Finally, comparison experiments on the real dataset demonstrate that our proposed algorithm is reliable.*

Keywords: Learning path planning, Directed learning path network, KL divergence, D-value matrix similarity

1. Introduction. Studying learning path planning can help find useful implicit learning behavior patterns from learners' online learning behavior data, which is conducive to helping beginners or learners with low participation to reasonably arrange the learning sequence of online knowledge points, so as to complete their learning goals efficiently and systematically [1-4].

Complex networks are widely used in learning path planning [5-7]. Shi et al. [8] proposed a learning path recommendation algorithm based on multi-dimensional knowledge graph to generate personalized learning paths that meet different learning objectives. Chungo [9] generated four learning styles based on Kolb's learning style scale. Yet, these learning path planning algorithms fail to take into account that the local structure of the knowledge nodes in the learner's learning path network and the learner's mastery of knowledge points will affect the reliability of the learning path planning algorithm. On the study of the similarity of nodes in complex networks, the local similarity model of undirected complex network nodes proposed by Zhang et al. [10] fails to take account of the similarity of local structure of directed network nodes. And the traditional European-based norm matrix similarity model will become less effective when processing high-dimensional time series data of learning behavior [11].

To solve the above problems, this paper studies the learning behavior data of a certain course as the research object, and proposes a learning path planning algorithm based on KL divergence and D-value matrix similarity. By comparisons, experimental results

indicate that our proposed methods are able to make sound recommendations on appropriate learning paths with significantly improved learning results in terms of accuracy and efficiency. The main contributions of this paper are as follows.

First, based on the KLD, a local structural similarity calculation method for directed network knowledge nodes is proposed to better characterize the similarity measure of the in-degree and out-degree of directed network nodes.

Second, based on the learner KLD matrix, a method for characterizing learning behavior similarity by using D-value matrix similarity is proposed.

The main content of this paper is as follows. Section 2 introduces the relevant definitions and related formulas of the algorithm. Section 3 proposes a directed learning-path-planning algorithm based on the interaction degree of knowledge points. Section 4 analyzes the proposed algorithm. Section 5 concludes the work and looks ahead.

2. Related Definitions. This section describes the relevant definitions and calculation methods of the proposed algorithm, and analyzes and explains some of the definitions.

Definition 2.1. Interaction degree of knowledge point concepts. *When learners study online knowledge points, whether they are proficient in using each knowledge point is portrayed by virtue of the concept interaction degree of knowledge points (CKP), the learner's mastery degree of knowledge points (MKP) and the relative difficulty coefficient of knowledge points (FKP), as demonstrated below:*

$$ckp_i = \frac{mkp_i}{fkp_i} \quad (1)$$

In the formula, ckp_i , mkp_i and fkp_i respectively represent the elements in the set of concept interaction degree of knowledge points $CKP = \{ckp_i | i = 1, \dots, m\}$, the set of learner's mastery degree of knowledge points $MKP = \{mkp_i | i = 1, \dots, m\}$ and the set of the relative difficulty coefficient of knowledge points $FKP = \{fkp_i | i = 1, \dots, m\}$ in the learner's online learning. And m represents the number of knowledge points of online videos that learners have learned.

Definition 2.2. A directed learning path network based on the interaction degree of knowledge points. *The directed learning path network (DLPN) based on the interaction degree of the knowledge points is a topological network generated based on the learning behavior time series data of the learner's online learning. The nodes in the network represent the knowledge points of the learners' online learning, the knowledge node values are characterized by the interaction degree of the knowledge points, and the knowledge edges are characterized by the sequential order of the learners' learning the knowledge points, and the weight of the edges is jointly portrayed through the interaction degree of various knowledge nodes.*

$$DLPN = G(M, CKP, E, W) \quad (2)$$

In the formula, DLPN characterizes a directed learning path network based on the concept interaction degree of knowledge points; M represents the set of knowledge nodes in DLPN; CKP characterizes the concept interaction degree of knowledge points; E represents the set of edges between the knowledge nodes in DLPN; and $W = [w_{ij}]_{m \times m}$ represents the weight matrix between the knowledge nodes in DLPN, of which the calculation method of the edge weight from knowledge node i to knowledge node j is shown in Formula (3):

$$w_{ij} = \frac{ckp_i}{ckp_j} \quad (i, j = 1, \dots, m) \quad (3)$$

In the formula, ckp_i and ckp_j respectively represent the concept interaction degree of knowledge points of knowledge node i and knowledge node j .

Definition 2.3. Knowledge point local structure. In the DLPN of learner e , the directly connected in-degree (DCID) knowledge node set with the knowledge node i is defined as $DKN_{ID}(i)$, and the indirectly connected in-degree (ICID) knowledge node set with knowledge node i is defined as $IKN_{ID}(i)$. The directly connected out-degree (DCOD) knowledge node set with the knowledge node i is defined as $DKN_{OD}(i)$, and the indirectly connected out-degree (ICOD) knowledge node set with the knowledge node i is defined as $IKN_{OD}(i)$. $dkn_{ID}(k)$, $ikn_{ID}(k)$, $dkn_{OD}(k)$ and $ikn_{OD}(k)$ represent the DCID knowledge node, ICID knowledge node, DCOD knowledge node, and ICOD knowledge node, respectively. And l_{ID}^{DKN} , l_{ID}^{IKN} , l_{OD}^{DKN} and l_{OD}^{IKN} respectively represent the number of elements of the set $DKN_{ID}(i)$, $IKN_{ID}(i)$, $DKN_{OD}(i)$ and $IKN_{OD}(i)$. $dci(k)$, $ici(k)$, $dco(k)$ and $ico(k)$ correspond to DCID, ICID, DCOD and ICOD respectively. The in-degree set and out-degree set of knowledge node i are defined as $ID(i)$ and $OD(i)$ respectively, and the numbers of elements of $ID(i)$ and $OD(i)$ for knowledge node i are respectively defined as $L_{ID}(i)$ and $L_{OD}(i)$.

$$\begin{cases} ID(i) = DCI(i) + ICI(i) = \{id(i, 1), \dots, id(i, k), \dots, id(i, L_{ID}(i))\} \\ OD(i) = DCO(i) + ICO(i) = \{od(i, 1), \dots, od(i, k), \dots, od(i, L_{OD}(i))\} \end{cases} \quad (4)$$

In the formula, $id(i, k)$ represents the elements of $ID(i)$, and $od(i, k)$ represents the elements of $OD(i)$. $L_{ID}(i)$ and $L_{OD}(i)$ are expressed by Formula (5), which is as follows:

$$\begin{cases} L_{ID}(i) = \left(\sum_{\alpha=1}^{l_{ID}^{DKN}} dci(\alpha) + \sum_{\beta=1}^{l_{ID}^{IKN}} ici(\beta) \right) \\ L_{OD}(i) = \left(\sum_{\alpha=1}^{l_{OD}^{DKN}} dco(\alpha) + \sum_{\beta=1}^{l_{OD}^{IKN}} ico(\beta) \right) \end{cases} \quad (5)$$

Definition 2.4. Knowledge node local structure similarity measure. In order to calculate the local structural similarity between the knowledge nodes in the learner's DLPN, the in-degree probability set $P_{ID}\{i\}_{i=1, \dots, N}$ and the out-degree probability set $P_{OD}\{i\}_{i=1, \dots, N}$ of each knowledge node in the learner's DLPN should have the same length, and the maximum length of the in-degree and out-degree sets of all the knowledge nodes in the learner's DLPN should be selected and defined as L . When the length of the in-degree node set or the out-degree node set is less than L , the remaining elements are set to 0. That is:

$$L = \max((l_{ID}^{DKN} + l_{ID}^{IKN}), (l_{OD}^{DKN} + l_{OD}^{IKN})) \quad (6)$$

$$\begin{cases} P_{ID}(i) = \{p_{ID}(i, k) | k = 1, \dots, L\} \\ P'_{ID}(i) = \{p'_{ID}(i, k) | k = 1, \dots, L\} \end{cases} \quad (7)$$

where $p_{ID}(i, k)$ are represented by Formula (8). To better describe the local structural similarity between knowledge nodes, according to the literature review [10], $P_{ID}(i)$ of each knowledge node need to be ordered, and the sorted in-degree probability degree set is $P'_{ID}(i)$.

$$p_{ID}(i, k) = \begin{cases} w_{ik} \cdot \left(\frac{id(i, k)}{\left(\sum_{\alpha}^{l_{ID}^{DKN}} dci_{ID}(\alpha) + \sum_{\beta}^{l_{ID}^{IKN}} ici_{ID}(\beta) \right)} \right) & k \leq L \\ 0 & k > L \end{cases} \quad (8)$$

The in-degree KLD of the knowledge node i and knowledge node j calculated according to the sorted probability set should follow:

$$H_{KL}(P'_{ID}(i) | P'_{ID}(j)) = \begin{cases} \sum_{k=1}^{l'_{ID}} \left(p'_{ID}(i, k) \ln \frac{p'_{ID}(i, k)}{p'_{ID}(j, k)} \right) & (p'_{ID}(j, k) \neq 0) \\ 0 & \text{else} \end{cases} \quad (9)$$

To avoid the calculation results to infinitive, l'_{ID} is represented as follows:

$$l'_{ID} = \min((l_{ID}^{DKN}(i) + l_{ID}^{IKN}(i)), (l_{ID}^{DKN}(j) + l_{ID}^{IKN}(j))) \quad (10)$$

In the same way, the out-degree probability set $P'_{OD}(i)$ and the out-degree KL divergence of the knowledge nodes i and j , namely, $H_{KL}(P'_{OD}(i)|P'_{OD}(j))$ can also be obtained.

The KLD matrix in the DLPN of learner e is obtained:

$$KL(e) = [kl_{ij}]_{m \times m} \quad (11)$$

where kl_{ij} are represented by Formula (12).

$$kl_{ij} = H_{KL}(P(i)|P(j)) = H_{KL}(P'_{ID}(i)|P'_{ID}(j)) + H_{KL}(P'_{OD}(i)|P'_{OD}(j)) \quad (12)$$

Definition 2.5. Learning behavior similarity. Based on the learner's KL divergence matrix, a learning behavior similarity model of D-value matrix similarity is proposed, and the corresponding elements of the learner e 's KL divergence matrix $KL(e)$ and the learner f 's KL divergence matrix $KL(f)$ are subtracted to obtain the difference matrix DVM, namely:

$$DVM = KL(e) - KL_B(f) \quad (13)$$

Use s_{ef} to describe the similarity between matrix $KL(e)$ and matrix $KL(f)$, namely:

$$s_{ef} = 1 - \frac{\sum_{i,j=1}^m |DVM(i,j)|}{z_e + z_f - z_{ef}} \quad (14)$$

wherein z_e and z_f respectively represent the number of elements in the learner e 's KL divergence matrix $kl \neq 0$, the learner f 's KL divergence matrix $kl \neq 0$. z_{ef} represents the learner e 's KL divergence matrix $kl \neq 0$ and the number of elements in learner f 's KL divergence matrix $kl \neq 0$. $|DVM|$ represents the absolute value of each element of the difference matrix DVM. According to Formulas (13) and (14), the learning similarity matrix of N learners is defined as $S = [s_{ij}]_{N \times N}$.

3. The Pseudo Code of Learning Path Planning Algorithm Based on KL Divergence and D-Value Matrix Similarity. It is the main idea of the algorithm that the learner's online learning behavior data is modeled to obtain the learner's concept interaction degree of online video knowledge points which is combined with the directed weighted complex network theory. Then, the clustering algorithm and the optimal path algorithm are used to generate a personalized learning path.

Learning path planning algorithm based on KL divergence and D-value matrix similarity

Input: Learner's online learning behavior data set $D = \{d_1, \dots, d_e, \dots, d_N\}$, MKP , FKP

Output: Three types of learner's optimal learning path

1: According to Definition 2.1, set CKP is generated

2: for each $d_e \in D$ do

3: Construct learner e 's DLPN according to Definition 2.2

4: for $i = 1 : m$ do

5: According to Definition 2.3, get $P_{ID}\{i\}_{i=1, \dots, N}$ and $P_{OD}\{i\}_{i=1, \dots, N}$ of all the knowledge nodes in the learner e 's DLPN

6: end for

7: Get the KLD matrix and the learning behavior similarity matrix S of all learners

8: end for

9: According to the learner's KLD matrix and the learning behavior similarity matrix S , the DNSCAN algorithm is used to derive the three types of learners of the knowledge point concept interaction level, primary, intermediate and advanced, and draw the optimal learning path for each type of learners.

10: Return the optimal learning path of three types learners.

4. Experimental Evaluation. To verify the reliability of the algorithm proposed in this paper, the authors first used the ACC and ARI clustering indicators to judge the quality of the clustering algorithm. The average weighting and average path length were used to determine the quality of the DLPN. Last, the results of the empirical data of the learning behavior of the students enrolled in 2020 were taken to verify the reliability of the proposed algorithm.

4.1. Dataset of the experiment. This experiment analyzed the online learning behavior data of learners, based on the learning behavior data of software engineering students on the online learning platform. The selected course is *Data Structures and Algorithms*. The course contains 207 video knowledge points, covering 1,198 learners enrolled in 2017, 2018, 2019 and 2020. There are 293,751 learning behavior data. The deadline for obtaining data is April 1, 2020. All experiments are conducted in Matlab R2018b.

4.2. Experimental results and analysis.

4.2.1. Clustering experiment results and analysis. Clustering learners through a clustering algorithm is to plan a learning path for learners of different levels that is more in line with their initial cognitive level. In this section, we use the spectral clustering algorithm based on distance matrix similarity (SC-D), DBSCAN algorithm based on distance matrix similarity (DBSCAN-D), the spectral clustering algorithm based D-value matrix similarity (SC-DVM) and DBSCAN algorithm based on D-value matrix similarity (DBSCAN-DVM) which were applied to conducting cluster analysis on the learning behavior data of the students enrolled in 2017, 2018 and 2019 (hereinafter referred to as students of 2017, 2018 and 2019) respectively, so as to classify the students of different grades into three types of learners in terms of their concept interaction degree of knowledge points: primary, intermediate and advanced to be specific. According to the learner's initial cognitive level, the learners were divided into three different groups, with primary, intermediate and advanced concept interaction degree of knowledge points. ACC and ARI analysis were performed on the initial classification of learners' concept interaction degree of knowledge

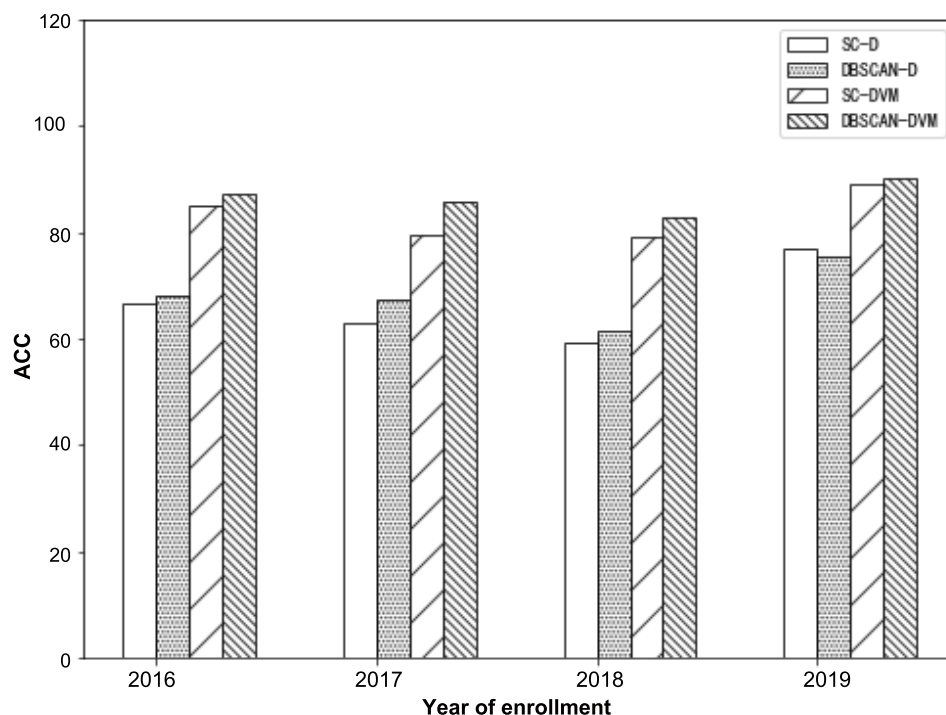


FIGURE 1. ACC comparison chart

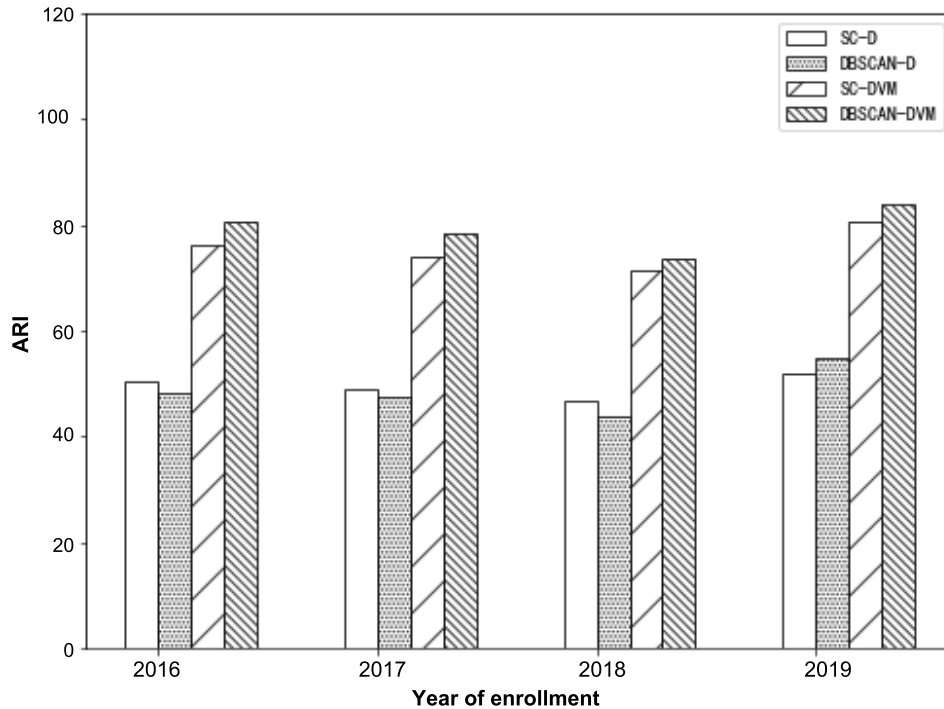


FIGURE 2. ARI comparison chart

points based on their cognitive levels and the classification results of four algorithms, as shown in Figure 1 and Figure 2, respectively.

As shown in Figure 1 and Figure 2, DBSCAN-DVM algorithm is superior to SC-D, DBSCAN-D and SC-DVM algorithm in ACC and ARI clustering indicators. Therefore, DBSCAN-DVM algorithm has better clustering effect.

4.2.2. Directed weighted complex networks experimental results and analysis. In this section, the DBSCAN-DVM algorithm was used to cluster online learning behavior data of students of 2017, 2018 and 2019, so as to get the optimal learning path for the primary, intermediate and advanced learners of the concept interaction degree of knowledge points in each grade, for the reference of students of 2020. Figure 3 shows the optimal learning path for the advanced level of concept interaction degree of knowledge points for the students of 2019.

The average weighted degree and average path length of the primary, intermediate and advanced DLPN of the concept interaction degree of knowledge points in each grade are shown in Table 1, Table 2 and Table 3. The average weighted degree of grade 2019 is the highest and the length of level road is the smallest, so the learning path of grade 2019 can be the best.

In order to test whether the primary, intermediate and advanced optimal learning paths of students of 2019 in terms of the concept interaction degree of knowledge points have improved the learning effect of students of 2020. Before the learners of 2020 learned the course, the 108 students who took the course had been divided into three categories: primary, intermediate and advanced, according to their initial learning ability. After the reference path for learning the course was given, the DBSCAN-DVM algorithm was used to perform cluster analysis on 108 learners. Table 4 shows the distribution of the number of primary, intermediate and advanced learners before and after referring to the reference path of course learning. Group A is the number distribution before the reference optimal learning path, and group B is the number distribution after the reference optimal learning path.

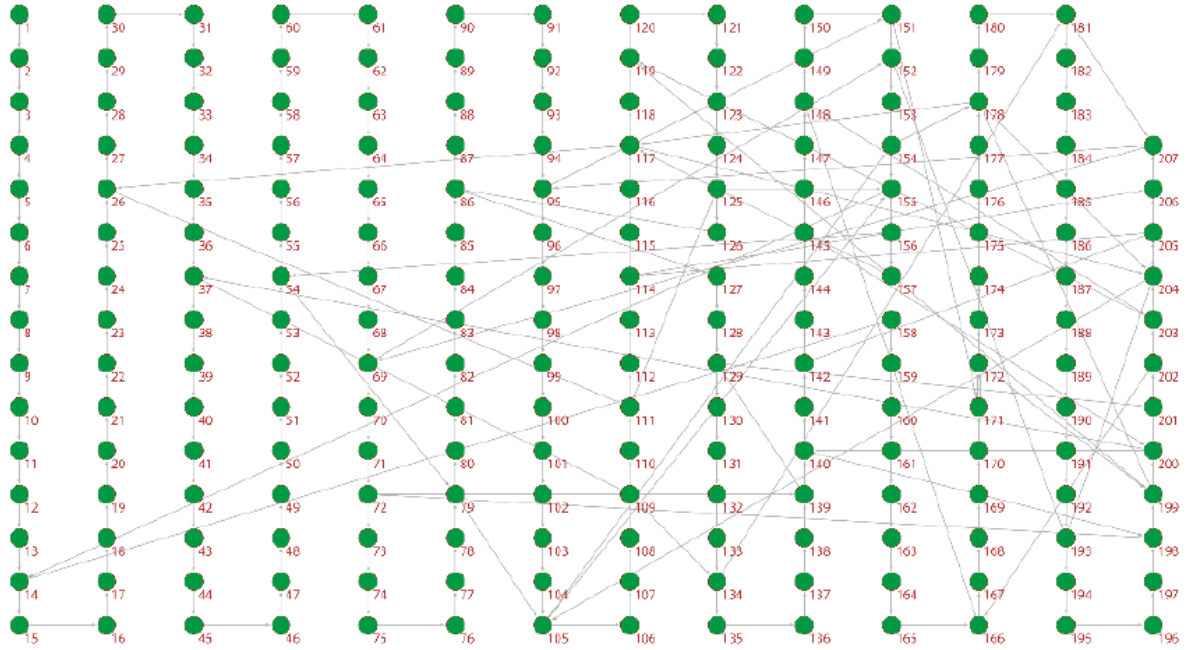


FIGURE 3. The optimal learning path map of the advanced learners of students of 2019 according to their concept interaction degree of knowledge points

TABLE 1. Learning path network structure evaluation index table based on knowledge point interaction level (Primary)

Year of enrollment	Average weighting	Average path length
2017	123.47	38.27
2018	114.15	46.51
2019	199.52	13.25

TABLE 2. Learning path network structure evaluation index table based on knowledge point interaction level (Intermediate)

Year of enrollment	Average weighting	Average path length
2017	264.29	19.98
2018	222.99	30.82
2019	506.32	11.76

TABLE 3. Learning path network structure evaluation index table based on knowledge point interaction level (Advanced)

Year of enrollment	Average weighting	Average path length
2017	361.28	39.46
2018	305.94	39.63
2019	660.58	6.63

TABLE 4. Comparison of online testing scores between group A and group B

Group category	Primary	Intermediate	Advanced
Group A	21	82	5
Group B	15	84	9

It can be seen from Table 4 that the number of learners with primary level of concept interaction degree of knowledge points has decreased after learning the course with reference to the optimal learning path. On the other hand, the number of intermediate and advanced students has increased. In summary, the optimal path maps for the students of 2019 in terms of primary, intermediate and advanced concept interaction degree of knowledge points can be recommended to both learners and teachers, so as to improve students' academic performance and teachers' teaching effect.

5. Conclusions. Aiming at the limitation of traditional collaborative filtering on similarity calculation, we propose a knowledge point recommendation algorithm based on similarity optimization. It fully considers the differences in learners' relative difficulty coefficient of knowledge points in different dimensions, and at the same time incorporates information about knowledge points learned by non-associated learners, improving the performance and quality of the recommendation algorithm. The future work will focus on artificial intelligence and optimize the performance indicators of the recommended algorithm with the combination of relevant knowledge of deep learning.

Acknowledgments. This work was supported by the National Natural Science Foundation of China (61967005, 61662013, U1501252), Innovation Project of GUET Graduate Education (2020YCX022), the Key Laboratory of Cognitive Radio and Information Processing Ministry of Education (CRKL190107).

REFERENCES

- [1] I. Kamsa, R. Elouahbi and F. El khoukhi, The combination between the individual factors and the collective experience for ultimate optimization learning path using ant colony algorithm, *International Journal on Advanced Science, Engineering and Information Technology*, vol.8, no.4, pp.1198-1208, 2018.
- [2] V. Vanitha, P. Krishnan and R. Elakkiya, Collaborative optimization algorithm for learning path construction in E-learning, *Computers and Electrical Engineering*, vol.77, pp.325-338, 2019.
- [3] P. Dwivedi, V. Kant and K. K. Bharadwaj, Learning path recommendation based on modified variable length genetic algorithm, *Education & Information Technologies*, pp.106-120, 2018.
- [4] Y. Zhu, P. Wang, Y. Fan and Y. Chen, Research of learning path recommendation algorithm based on knowledge graph, *Proc. of the 6th International Conference on Information Engineering*, 2017.
- [5] H. Liu and X. Li, Learning path combination recommendation based on the learning networks, *Soft Computing*, vol.24, no.6, pp.4427-4439, 2020.
- [6] H. Zhu, Y. Liu, F. Tian et al., A cross-curriculum video recommendation algorithm based on a video-associated knowledge map, *IEEE Access*, vol.6, pp.57562-57571, DOI: 10.1109/ACCESS.2018.2873106, 2018.
- [7] H. Zhu, F. Tian, K. Wu et al., A multi-constraint learning path recommendation algorithm based on knowledge map, *Knowledge-Based Systems*, vol.143, pp.102-114, 2018.
- [8] D. Shi, T. Wang, H. Xing et al., A learning path recommendation model based on a multidimensional knowledge graph framework for e-learning, *Knowledge-Based Systems*, vol.195, pp.201-215, 2020.
- [9] S. Chungo, Designing and developing a novel hybrid adaptive learning path recommendation system (ALPRS) for gamification mathematics geometry course, *Eurasia Journal of Mathematics Science & Technology Education*, vol.13, no.6, pp.2275-2298, 2017.
- [10] Y. Zhang, G. Lin and J. Li, Research of knowledge mapping construction method based on scientific research results, *The 7th International Conference on Computer Engineering and Networks*, 2017.
- [11] Z. Shou, H. Tian, S. Li et al., Outlier detection with enhanced angle-based outlier factor in high-dimensional data stream, *International Journal of Innovative Computing, Information and Control*, vol.14, no.5, pp.1633-1651, 2018.