

Received October 26, 2018, accepted November 13, 2018, date of publication November 20, 2018,
date of current version December 19, 2018.

Digital Object Identifier 10.1109/ACCESS.2018.2882240

Integration of Data Mining Clustering Approach in the Personalized E-Learning System

SAMINA KAUSAR^{1,2}, XU HUAHU¹, IFTIKHAR HUSSAIN^{3,4},
ZHU WENHAO¹, AND MISHA ZAHID⁴

¹School of Computer Engineering and Science, Shanghai University, Shanghai 200444, China

²Department of CS & IT, University of Kotli Azad Jammu and Kashmir, Kotli Azad Kashmir 11100, Pakistan

³Instituut Voor Mobiliteit (IMOB), Hasselt University, 3500 Hasselt, Belgium

⁴School of Computer and IT, Beaconhouse National University, Terogil Campus, Lahore 53700, Pakistan

Corresponding author: Iftikhar Hussain (iftikhar.hussain@bnu.edu.pk)

This work was supported in part by the National Natural Science Foundation of China under Grant 61572434 and Grant 91630206, and in part by the Shanghai Science and Technology Committee under Grant 16DZ2293600.

ABSTRACT Educational data-mining is an evolving discipline that focuses on the improvement of self-learning and adaptive methods. It is used for finding hidden patterns or intrinsic structures of educational data. In the arena of education, the heterogeneous data is involving and continuously growing in the paradigm of big-data. To extract meaningful information adaptively from big educational data, some specific data mining techniques are needed. This paper presents a clustering approach to partition students into different groups or clusters based on their learning behavior. Furthermore, the personalized e-learning system architecture is presented, which detects and responds to teaching contents according to the students' learning capabilities. The primary objective includes the discovery of optimal settings, in which the learners can improve their learning capabilities. Moreover, the administration can find essential hidden patterns to bring the effective reforms in the existing system. The clustering methods K-Means, K-Medoids, Density-based Spatial Clustering of Applications with Noise, Agglomerative Hierarchical Cluster Tree and Clustering by Fast Search and Finding of Density Peaks via Heat Diffusion (CFSFDP-HD) are analyzed using educational data mining. It has been observed that more robust results can be achieved by the replacement of existing methods with CFSFDP-HD. The data mining techniques are equally effective in analyzing the big data to make education systems vigorous.

INDEX TERMS Big data, clustering, data-mining, educational data-mining, e-learning, profile learning.

I. INTRODUCTION

Educational data-mining (EDM) is a new perspective in modern educational systems. It is concerned with the study and development of new adaptive methods, instruments to artificially analyze and visualize the hidden patterns or intrinsic structures in educational datasets. Mostly, education related datasets contain structured, semi-structured and un-structured data with different geographical distribution [1]. EDM has emerged as a promising area of research aimed to analyze the intrinsic data structures, extracting hidden predictive information and finding insights into educational datasets [2]. In the field of education, EDM can be demarcated as an application of data-mining methods to exploit novel patterns and artificially analyze big data efficiently and effectively.

Recently, frontier technologies i.e. the Internet-of-Things (IoT), sensors, artificial intelligence, and social networks are being integrated with educational system for

effective learning [3], [4]. Web based systems are computer-aided virtual forms of instructions that are independent of geographical location. Sensors and IoT produce massive amount of data that lead towards the big data dilemma [5]. Moreover, big data has significant impact in scientific studies, public health, industrial applications, and in the field of education [6]–[10]. In educational field, the huge amount of data provides a new insight in improving the learning capabilities and decision making skills of teachers and students. The educational data mining may play an important role in improving the education system by (1) refining the individual based quality education, (2) discovering new areas of knowledge and finding associations among different fields and (3) finding new perspectives in experimental data.

With the advancement in communication technologies, nowadays many smart devices and sensors [11] are incorporated into educational systems to observe the overall

behavior of the education system. It contains rich information of people's thoughts about different events in semi-structured or unstructured form. Most of the web based learning methods are static and fail to take into account the diversity of students. These virtual educational systems can be improved by utilizing data mining techniques, in order to effectively meet the needs of diverse learners. In general, there is a wide variety of data mining methods that can be applied in the field of education. These methods can be categorized into classification, clustering, neural networks, and relationship mining. Clustering is a primary unsupervised method to partition datasets into separate groups (clusters) based on the estimated intrinsic characteristics or similarities [12] and has been applied in several fields [13]–[18]. Clustering methods can be considered and categorized as: partition-based, density-based, model-based and hierarchy-based [19]–[23] and cannot be directly used to handle the complexities of big data.

A. RESEARCH OBJECTIVES

This paper presents a data-mining clustering approach to partition students into different groups based on their learning behavior. The offered approach has been established on the basis of big relational databases. A personalized e-learning system architecture integrating data mining technique is presented; which creates and responds teaching content according to students' learning capability. The primary objective includes the discovery of optimal settings, in which learners can improve their learning capabilities. Additionally, the administration can find essential hidden patterns to bring the effective reforms in the existing system. This paper also analyzes the K-Means, K-Medoids, Density-based Spatial Clustering of Applications with Noise (DBSCAN), and Agglomerative Hierarchical Cluster Tree (AHCT) approaches for clustering and compares them with the recommended approach "Clustering by Fast Search and Finding of Density Peaks via Heat Diffusion (CFSFDP-HD)". A contrast between existing approaches and CFSFDP-HD in regard to academic performance of students is also examined. The undertaken research contributed by presenting and integrating an adaptive data-mining approach for clustering big-data in the field of education. The recommended approach gives accurate and efficient results as compared with some of the existing (described) clustering methods.

B. PAPER ORGANIZATION

This paper is organized as follows: Section 2 presents the literature review of data mining clustering techniques with some specific tools to deal with the educational data. Some of the literature on big data and e-learning systems is also presented. Section 3 describes the recommended data mining clustering approach. Some of the existing data mining clustering approaches are also discussed. A conceptual personalized e-learning system architecture integrating the recommended data-mining clustering approach is also presented.

Different steps of the presented approach are also discussed at the end of the section 3. Section 4 presents the experiments and some of the results with discussion by considering a specific case study. Finally, the conclusion and the future research are discussed in Section 5.

II. LITERATURE REVIEW

This section presents some of the literature on: (1) educational data-mining approaches with some specific tools to deal with educational data and (2) big data and e-learning systems.

A. CLUSTERING TECHNIQUES IN EDM

"Educational data-mining is emerging as a research area with a suite of computational and psychological methods, and research approaches for understanding how students learn" [24]. Various clustering methods have been applied in recent studies to predict students' performance. de Morais *et al.* [25] researched techniques to support the teacher decision-making process by grouping students and planning challenges accordingly and reached a positive conclusion. Prakash *et al.* [26] have researched learning analytic techniques for big data in educational data-mining to find out the Adaptive learning systems (ALS). The ALS empowers teachers to rapidly observe the adequacy of their adjustments and mediations, giving input to persistent change. The outcomes of the study are coherent with the conclusions of the study presented by Algarni [27].

In [27] author explored various studies and datasets revolving around the field of EDM. Author derived that EDM can be utilized as part of wide range zones including (a) recognizing the students who are at risk, (b) distinguishing needs for the adapting needs of various students' groups, (c) expanding graduation rates, (d) adequately surveying institutional execution, (e) boosting ground assets, and (f) upgrading subject educational modules reestablishment. The outcomes of the study are consistent with the conclusions of the study presented by Bovo *et al.* [28] after studying student records on different trainings and successfully predicting students who are falling behind. Another research study [29] consistent with [27] is conducted by Amjad Abu Saa examines and predicts student performance in different scenarios using clustering methods. In a similar study [30], Tommaso and Alex Bowers have analyzed various analytical techniques: Educational data-mining, Learning and Academic Analytics, and have reached the conclusion that the applications of data-mining methods (with responsibility and professionalism) yield positive results.

The K-means is a state-of-the-art partition based clustering algorithm and have been applied in EDM. Such as, special selection of student's seat in lab or classroom and its impact on student's assessment has been evaluated by Ivančević *et al.* [31]. Another study presented by Ying *et al.* [32] has utilized K-means to understand the behavior of students based on the annotation dataset of 40 students. In a study conducted by Eranki *et al.* [33],

K-means was applied to examine the influence of human characteristics on student's performance while listening to music, yielding evident classifications.

B. BIG DATA AND E-LEARNING SYSTEMS

Big data has the capability to benefit students distinctly by providing them with a modern and dynamic education system. In the study [34], Athanasios and Panagiotis analysed the goals, purposes, and benefits of big-data and open-data in Education. Authors concluded that the education system can be enhanced by embracing new learning approaches to make it more effective and focused on. Moreover, Manohar *et al.* [35], support the same idea and anticipated that the big data can be effectively used in predicting student results, and improving both the teaching and the learning experience. The research conducted by Tulasi [36] and Daniel [37], targeted the higher education and explored the solutions proposed by big data systems to the challenges faced by higher education. Dede [38] further advanced the topic by studying "next steps" that can be taken using big data in education and concluded that the field has a lot of potential in the betterment of the individual learning experiences.

Numerous researchers have expressed that personalization, in an academic setting, permits executing more proficient and viable learning forms. Various works are attempting to enhance the quality and viability of e-learning by utilizing standards of other research zones. This pattern of personalization advancement additionally shows up in e-learning. Gaeta *et al.* [39] have introduced a new tool: Intelligent Web Teacher (IWT) to support Personalized E-Learning in their study on personalized e-learning. The comparison of traditional methods with IWT deduce that personalization permits executing more proficient and powerful e-Learning forms, featuring an expanding level of fulfilment by both educator and students. The system developed by Prawira *et al.* [40] using Moodle proved to be capable of improving the learning process and collaboration between the teacher and student in higher education.

Yarandi *et al.* [41] take individual learning capabilities of students to present an ontology-based technique to improve an adaptive e-learning scheme. The proposed e-learning system creates content permitting to the learner's knowledge. Maselena *et al.* [42] present a personalized e-learning framework and argue that the system encourages strong learning condition due to the realization of individual needs. The results are coherent with a study by Wu *et al.* [43] which presents a theoretical framework of adaptive e-learning, self-assessment and dynamic scaffolding theory. The system provides tailored learning material to students based on the student ability. A grid agent model was proposed by Liu and Liu [44] in their study for effective adaptation of e-learning systems using artificial psychology to individual students who would benefit from this personalization. Furthermore, Li and Chang [45] have proposed another

personalized e-learning system which is a feedback extractor with fusion capability to adjust the user preferences.

The above mentioned literature shows that personalized e-learning schemes are effective tools in individual learning. In the e-learning, a significant huge amount of data is continuously generated and ubiquitously available on the web. Therefore, more sophisticated and frontier clustering methods are required to benchmark on EDM data to get intrinsic insights. To cope with aforementioned issues, a data mining clustering approach is recommended and integrated with the personalized e-learning system. The integration of data-mining approaches makes the learning system more interesting.

III. DATA MINING CLUSTERING APPROACH AND THE PERSONALIZED E-LEARNING SYSTEM

In this paper, a data mining based clustering approach "CFSFDP-HD" is presented to partition students into groups and is established on the basis of big relational databases. The approach finds possible groups of students by comparing their similar learning behavior. It is sensitive to detect the understanding levels of students. Moreover, a conceptual Personalized E-Learning System Architecture (PESA) integrating data mining clustering approach has been described; which creates and responds to teaching content according to students' learning capability. For each group, the system generates different quizzes, assignments, study related games, and books' contents to improve the learning capabilities of students. To make groups and select appropriate teaching methods, system uses artificial intelligence and adaptive clustering methods. In the proposed architecture, the recommended approach can be used for clustering, profiling and content filtering to the group of students into appropriate classes. The traditional e-learning systems are mostly query-based and the queries are responded without any intelligence or heuristics.

A. THE EDUCATIONAL ENVIRONMENT

The primary agenda of higher education is to harness cross-disciplinary intelligence to improve syllabus, content delivery, enhancing learners' experiences and creating an atmosphere that integrates them with the skills and knowledge required to cope with the changes and challenges posed by the big data. In such a complex educational environment, it is tough for the human mind to identify patterns manually, but database projects have the ability to incorporate and link traditional and new data sources. Such compactness can create deeper insights into students learning capabilities and enhance classroom activities.

Grade Point Average (GPA) and percentage score are important indicators for the measurement of students' academic performance and capabilities. GPA is an important factor for academic planner to setup and analyze the learning environment in academic organizations [46]. The GPA or percentage score of students can be affected by different factors such as teaching methodology and attention of teachers

towards some particular students. It is a general phenomenon that teachers mostly focus on those students that take part in class activities and show satisfactory outputs. Moreover, there are some intrinsic hidden patterns that exist among the students. Students can be separated into different clusters on the basis of their progress. The same teaching method may not be effective for different clusters of students.

Institutional databases, having the teaching material and users' queries, are entertained according to the stored data. However, most of the updated knowledge is available at various places on the Internet. To reinforce the students learning capabilities, it might be credible to integrate the rest of data sources with e-learning system [47]. The data-mining clustering approaches can play a significant role in finding the relationships among different subjects available over the Internet, specifically in the e-learning systems.

Similarity measures and clustering are important tasks to find similar groups in educational big data. The similar patterns of data in different fields may be useful for researchers and learners to gain knowledge easily from various fields. For example, we can use partition based clustering, density based clustering, and hierarchical based clustering for text mining, to find the similarity between data points, outliers, and similar or related fields by clustering big data.

B. EXISTING CLUSTERING METHODS

Various clustering methods have been used in EDM such as Mean-shift clustering, K-means, K-medoids, Density Based Spatial Clustering of Applications with Noise (DBSCAN) in [21] and [48], and Hierarchical clustering in [49]–[52], however, these approaches are not robust in identifying significant clusters in ambiguous and noisy datasets [22]. A short description of these methods is as follows:

1) MEAN-SHIFT CLUSTERING

Mean-shift Clustering is a sliding-window based approach; tries to discover condensed areas of *data – points*. The goal of this approach is to detect the *center – points* of each cluster based on a centroid method. In the Mean-shift clustering method the candidates are updated for the center points (the mean of the points) within the sliding window. In the post processing stage, the candidate windows are filtered to remove near-duplicates and forms the final set of center-points and their matching clusters.

The Mean-shift clustering having radius 'r' (as the Kernel) begins with a circular sliding window centered at a randomly selected point C. Mean shift method shifts 'r' iteratively to a higher density region (in each step) until convergence. The density of each sliding window is proportional to its size (the points inside it). By shifting, the density of the points gradually moves towards areas of the high density point. The shifting of the sliding window continues until a shift cannot accommodate more points inside the kernel (no longer increasing the density). In case, when multiple sliding windows are overlapped; the data points are clustered according to the

sliding window in which they reside and the sliding window containing the most well-maintained points.

Mean-shift Clustering method has two disadvantages: (1) it is pretty calculation exhaustive, and (2) it trusts on satisfactory high data-density (with perfect gradient to find the cluster-center).

2) K-MEANS

The *K-means* [19] is a state-of-the-art partition based clustering algorithm. It is considered as a better approach than the Mean-shift clustering because it does not have the above mentioned problems. In K-means, input data is divided into *k* distinct groups, where *k* is an input parameter used to specify the number of output clusters. K-means iteratively improves the initial partitions until the optimized clusters are not found. Mathematically we can express K-means using the following expression:

$$\underset{S}{\operatorname{argmin}} \sum_{i=1}^n \sum_{x \in S_i} \|x - \mu_i\|^2 \quad (1)$$

where, μ_i is mean of data-points in S . S_i is initial partition of dataset $\{x_1, x_2, x_3, \dots, x_n\}$.

K-means is considered as the best choice to discover the signal of interest from educational datasets if significant number of clusters is already defined. However, it might be a hectic job to discover appropriate groups using K-means without prior knowledge of existing number of clusters or in presence of noisy or complex data. In EDM, it is rigid to setup the selection of clusters, and initial centroids setting of K-means. It is also obscure to find significant signal of interest.

3) K-MEDOIDS

The *K-medoids* [53] is used to partition the dataset into clusters similar to the *k*-means. The aim of both (K-means and K-medoids) methods is to minimize the sum of distances between data-points of a cluster and a central data-point of the same cluster. In distinction to the *k*-means, *k*-medoids picks data-points through centers (the medoids) and workings by a generality of the Manhattan Norm to express distance among data-points. It clusters the dataset of 'n' data-points into 'k' groups or clusters. To determine 'k', the silhouette is considered as a useful tool. K-medoids reduces an amount of pairwise variations instead of summation of squared Euclidean distances. That's why it is considered as more robust (to noise and outliers) than the *k*-means. The *k*-medoids method is effective to return the actual data-points (the medoids) of the dataset. It is also suitable for clustering the definite data where a mean does not exist.

4) DBSCAN

DBSCAN [21] is a density – based clustering method, begins through a random starting data point that has not stayed visited. The neighbors of the data-point are mined by a Distance – Epsilon ϵ . If the sufficient amount of *data – points*

is available within neighborhood, the clustering procedure starts and the current *data – point* is considered as the first *data – point*. The other data-points will be labeled as *noisy*. Later on, those noisy data-points may become the portion of the cluster. In both situations the data-points are noted as *visited*. For current cluster, the procedure of making all the data-points in the ε neighborhood is repeated to add all the new data-points. This process is recurring until all data-points in the current cluster are recognized and labeled as *visited*. Same process is repeated for all clusters.

DBSCAN doesn't execute noisy data when the clusters of variable density are established. When the density varies, (1) the distance threshold epsilon ε and minimum data-points for identifying the neighborhood data-points will differ from cluster to cluster (2) for very high-dimensional data, the distance threshold epsilon ε becomes challenging to estimate.

5) HIERARCHICAL CLUSTERING

In this section we will discuss only the Bottom-up hierarchical clustering. It treats each data-point as a single cluster at the beginning and then continuously combines pairs of clusters till entire clusters have been combined into the single cluster. It is also known as Agglomerative Hierarchical Cluster Tree (AHCT) [54]. The AHCT is represented as a tree. The roots of the tree are considered as a unique cluster.

In the beginning, each data-point is treated as a single cluster; where '*k*' data-points are treated as '*k*' clusters. In the next step a *distance metric* is selected to measure the distance between two clusters. Furthermore the two clusters are combined into one iteratively for each pair. The combined clusters are selected with the smallest average linkage; both the clusters having (i) the smallest distance and (ii) the most similar data-points. This step will continue until the root of the tree which is explicitly given in the start. The number of clusters can be selected by recognizing the given root number, which helps to stop combining the clusters.

Bottom-up hierarchical clustering method does not need to specify the number of clusters and has the ability to select the best cluster because of using a tree.

C. RECOMMENDED DATA MINING CLUSTERING METHOD

Clustering by Fast Search and Finding of Density Peaks (CFSFDP) has been proposed by Rodriguez and Laio [22]. It has characteristics to discover significant clusters in a more spontaneous manner as compared with the K-means. It empowers clustering procedure, in which high-density regions are identified as potential clusters, outliers are automatically identified and arbitrary shape of clusters is organized. In K-means to obtain meaningful clusters users are required to repeat clustering process multiple times with different parametric setting. While the unique approach utilized in CFSFDP to discover clusters and noise adaptively would be a significant clustering tool to analyze the educational data. The CFSFDP uses the following given methodology to discover significant clusters.

CFSFDP calculates local density (ρ_i) and a minimum distance (δ_i) for each given data-point *i*, with its nearest high density point. Table 1 shows different symbols with their description, used in different equations.

TABLE 1. Symbols used and their description.

Symbols	Meanings
n	The data-points ' <i>n</i> ' of the dataset.
k	<i>k</i> an input parameter used to specify the number of output clusters.
S_i	S_i is initial partition of dataset: $\{x_1, x_2, x_3, \dots, x_n\}$
x	x represents the data-point.
μ_i	μ_i is the mean of data-points in <i>S</i> .
ε	Epsilon ε is the distance threshold.
ρ_i	ρ_i is the local density of <i>i</i> .
δ_i	δ_i is the minimum distance
d_c	d_c is the cut-off distance: is an important parameter used to estimate the ρ_i of each <i>i</i> .
d_{ij}	d_{ij} is the distance from the <i>data – point i</i> to <i>data – point j</i> .

The ρ_i is equivalent to the number of *data – points* that are closer than the cut-off distance d_c to *i*. The d_c is a vital parameter used to estimate the ρ_i of each *i*. The usefulness of CFSFDP depends upon the proper choice of d_c . The local density can be estimated by utilizing the Equation (2) where the d_{ij} is the distance from the *data – point i* to *–point j*.

$$\rho_i = \sum_j X(d_{ij} - d_c) \quad (2)$$

where,

$$X(x) = \begin{cases} 1 & x < 0 \\ 0 & \text{otherwise.} \end{cases} \quad (3)$$

The local distance δ_i , can be computed by Equation (4), of a *data – point* to the nearest extremely condensed *data – point* \max_{ρ_i} . It is achieved for the purpose of assigning *i* to the nearby cluster. The value of δ is considered as maximum, when the *data – points* with high local or global density are found. On the other hand, the *data – points* having high ρ_i and large δ_i (compared to other *data – points*) are considered as cluster centers.

$$\delta_i = \begin{cases} \min_{j: \rho_j > \rho_i} (d_{ij}) & \text{if } \exists j \text{ s.t. } : \rho_j > \rho_i \\ \max_{j: \rho_j > \rho_i} (d_{ij}) & \text{otherwise} \end{cases} \quad (4)$$

Cluster centers are attained by plotting estimated values of ρ_i and δ_i , which is called the decision graph. Moreover, the CFSFDP allocates the remaining *data – points* to the closest cluster center and based on their δ values. In cluster analysis, the key challenge is to discover correct cluster centers in the datasets [1]. However, CFSFDP uses a decision graph to identify the correct *cluster – centers* with the least human interaction, which makes it more worthy to analyze big data / streaming data.

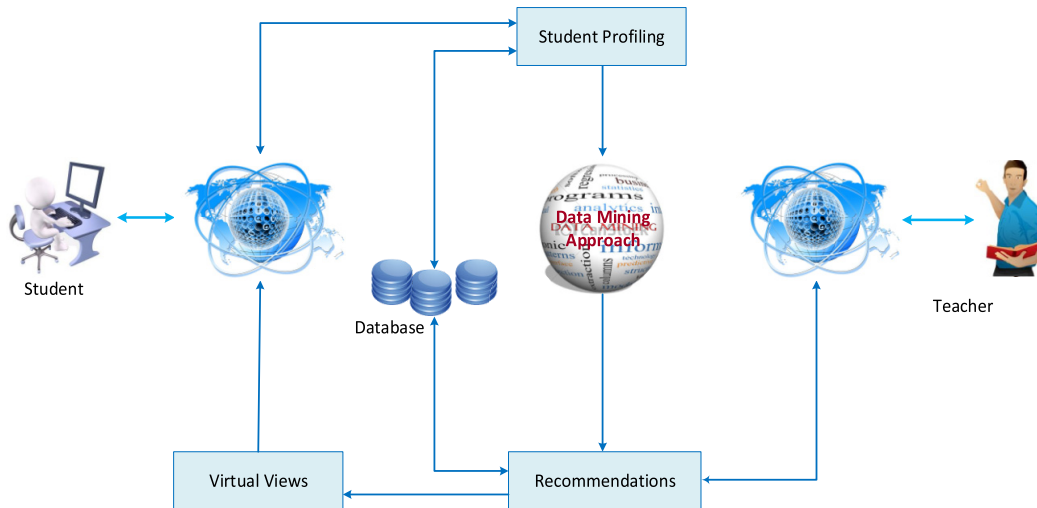


FIGURE 1. Personalized E-learning Architecture. A profile is created for each learner and is automatically updated based upon the activities of the learner.

The *CFSFDP* has characteristics to discover intrinsic hidden signal of interest from ambiguous data; it can be applied in existing education data mining systems and e-learning systems to produce more significant clusters. It can be used to cluster the similar documents, find plagiarism in documents, and analyze the students’ profiles to find the similar insights in different research areas. The *CFSFDP* via heat diffusion (*CFSFDP – HD*) [20] was proposed as a variant of *CFSFDP*, where limitations of *CFSFDP* are improved and users can analyze data without any prior domain knowledge. In *CFSFDP – HD*, an adaptive method was used to estimate density of underlying dataset, which is given in Equation (5):

$$\hat{f}(x;t) = \frac{1}{n} \sum_{j=1}^n \sum_{k=-\infty}^{\infty} e^{-k^2\pi^2t/2} \cos(k\pi x) \cos(k\pi x_j) \quad (5)$$

Where x represents the data-points and preparatory probability is distributed through the data-points $\{x_1, x_2, x_3, \dots, x_n\}$. The method evolves for a time t . The function \hat{f} in Equation 5 can be expressed as

$$\hat{f}(x;t) \approx \sum_{j=0}^{n-1} a_k e^{-k^2\pi^2t/2} \cos(k\pi x), \quad (6)$$

where n is a positive large interger and a_k is

$$a_k = \begin{cases} 1 & k = 0 \\ \frac{1}{n} \sum_{i=1}^n \cos(k\pi x_i) & k = 1, 2, \dots, n - 1, \end{cases} \quad (7)$$

Equation (6) is fully adaptive and may consider; firstly the *optimal bandwidth selection* and secondly the *boundary corrections*. It delivers *enhanced performance* and is also consistent with the *actual density*.

D. PERSONALIZED E-LEARNING SYSTEM ARCHITECTURE

The e-learning architecture responds to the individual demands of users, and is able to predict user preferences

or interests. E-learning not only allows the instructors and learners to meet virtually, but also makes sharing of resources possible electronically.

The overall Personalized E-learning System Architecture is shown in Figure 1. The major steps of the PESA are described as follows:

1) STUDENT PROFILING

The student interacts and manages his/her profile through the interface deployed on a desktop laptop or a smartphone. The user profile and other information change through the Internet. According to [55], [56], student profile or sometimes a student model called a typical group of students. Its function is to determine the user-learner needs and preferences automatically.

Student related data works like a seed for personalization of student queries and intelligent response of queries. Student profiling is an ongoing process which contains both static and dynamic data. Data collected in a static way [56] includes personal, personality, cognitive, pedagogical and preference data. Individual data define the bio-graphical information about the students. Personality data enlighten the students’ attention, cooperation and coordination skills. Student profile reflects the overall interest and behavior of the student. Cognitive data inform about the students’ cognition while pedagogical data describe different learning styles and methods. If the profile maintaining system detects any unusual behavior in student activities, it updates the profile accordingly.

2) DATA MINING CLUSTERING APPROACH

Data mining clustering is responsible in finding association, recommendation, and intelligence to provide customized and powerful learning mechanism for students. For example, appropriate content selection on the basis of the students’ interest and understanding is a big problem. This can be resolved by grouping whole contents by simply applying clustering approach to filter contents according to individual

student profile. Moreover, the key inference components in such e-learning systems are based on data mining approaches, which analyze the user's profile and suggest some sort of actions with the application of artificial intelligence. Moreover, especially, when we talk about clustering methods in existing systems are mostly based on the naïve clustering approaches such as Mean-shift clustering, K-means, and DBSCAN as mentioned earlier. Unlike existing e-learning systems, we recommended to use CFSFDP-HD methods to achieve feasible results. The recommended data mining clustering approach is already explained in the section 3.3.

3) RECOMMENDATIONS

This process is responsible in collecting data of interest from relational databases filtered according to student profile with the help of data mining approach. It also has the ability to prevent duplication of the information created before. This process recommends or proposes the solution to the instructor implicitly based on the output of recommended data mining clustering approach.

4) DATABASE

Database contains large datasets of courses and other education related activities. This component contains all the information that the student received from the instructor and also recommends or proposes instructions to the instructor.

5) VIRTUAL VIEWS

By considering the large datasets from the academic databases and by using the recommended data mining clustering approach, the intelligent analysis of records and selection of appropriate contents, virtual views are created and delivered to the students in the form of electronic documents. It can be achieve to identify different patterns which will help students to study, predict and improve their academic performance. The recommended clustering method is able to find groups of items so that the items that are in a cluster are similar to each other than to the items in another cluster. This may help to arrange different items which are under consideration. The clustering data mining approach would help in analyzing different profiles and may implicitly proposes the suitable educational items/materials to each student.

E. IMPLEMENTATION: STEPS INVOLVED IN THE RECOMMENDED EDM CLUSTERING APPROACH

The recommended data mining approach (CFSFDP-HD) is implemented using MATLAB to analyze the behavior and to simulate the educational data. The simulated educational data consists of students' (1) obtained marks and (2) class-attendance. The obtained marks consists of (1) three quizzes, (2) two assignments, (3) one midterm, and (4) one final-term exams while the class attendance is calculated on the basis of students' record of presence and absence. Only two percent weightage of class attendance is considered. The presented approach takes *distance matrix D* of dataset as input: *D* is the pairwise distance matrix of educational data.

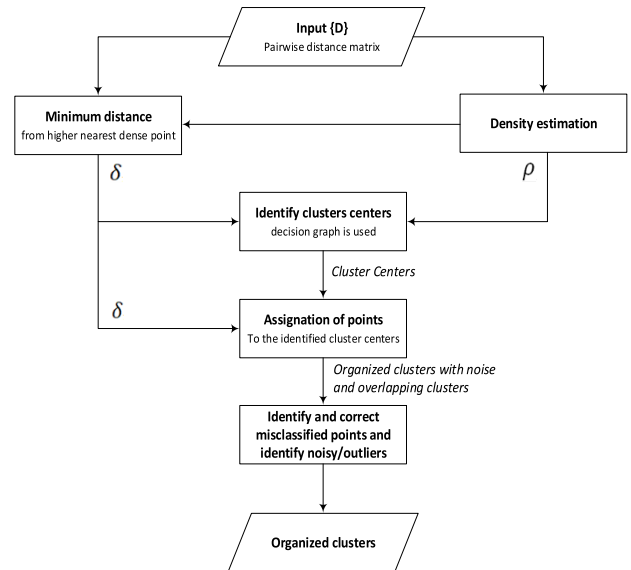


FIGURE 2. Key steps involved in the presented data mining approach (CFSFDP-HD).

The key steps of CFSFDP-HD along with the flow control are as follows (see Figure 2):

- 1) Step 1: In the first step, the proposed approach estimates the density ρ_i via heat diffusion using Eq. (5).
- 2) Step 2: the proposed approach calculates the minimum distance δ_i from the highest nearest dense points by using Eq. (3).
- 3) Step 3: the identification of cluster centers is achieved by the use of decision graph. In the decision graph, the ρ_i and δ_i are plotted. The output of this step is the Cluster Centers.
- 4) Step 4: The assignment of the remaining points to the identified cluster centers. The output of this step is the organized clusters with noise and overlapping clusters.
- 5) Step 5: In this step, the presented approach identifies and fixes the misclassified points and also identifies the noisy or outliers of the organized clusters (noisy and overlapping clusters).

The output of the proposed approach is presented by the *organized clusters*.

IV. RESULTS

The synthetic dataset of 600 students (enrolled in different sessions) is simulated by CFSFDP-HD approach to partition students into appropriate groups and is based on the students' class attendance and also on the obtained marks in: quizzes, assignments, mid and final exams. Progress-based segmentation of students is necessary to design appropriate teaching methods to address the weakness of a particular group in the class. In the Figure 3(a), the decision graph based heuristic approach is visualized to select the exact number of clusters intuitively. The full black points in Figure 3(a) are treated as non-cluster central-points while the colored points (the outliers) are considered as the central-points of the expected cluster. The decision graph is established after

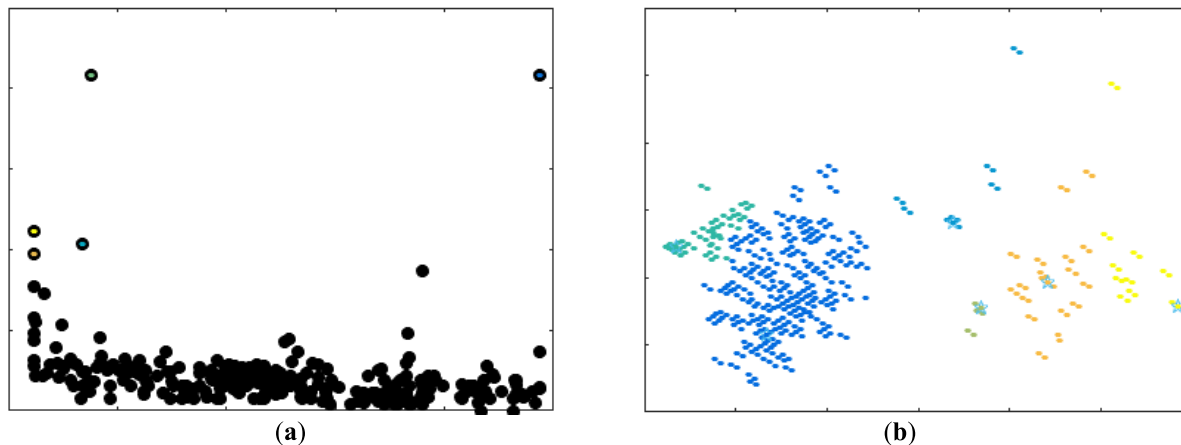


FIGURE 3. (a) In the decision graph, the parameters ρ (x-axis) and δ (y-axis) are plotted. The identification of cluster centers is achieved by the use of decision graph. (b) Assigning the remaining data-points to the identified cluster centers are shown with different color schemes; where different colors represent different groups.

(1) the estimation of density ρ_i via heat-diffusion and (2) the calculation of the minimum distance δ_i from the nearest point with higher density. The identification of cluster centers is achieved by the use of decision graph as shown in the Figure 3(a) where the ρ and δ are plotted.

With the minimum interpretation of heuristic approach to select the exact number of clusters, four distinct clusters are identified effectively, as shown in Figure 3(a), where outliers are treated as potential cluster centers and are represented with different colors. After identification of potential cluster centers, the remaining data-points are assigned to the identified cluster centers. The discovered clusters are shown with different colors scheme in Figure 3(b), where 2D Non-classical multidimensional scaling is used to visualize the dataset. The recommended approach “CFSFDP-HD” is adaptive in nature; so there is no need to set any parameter explicitly.

The aforementioned partition of students into four significant groups can play an important role to enhance the learning skills by paying special attention to a particular group of students. The self-motivator and talented students are separated from students with low and below to the average grading students. Based on the obtained different categories of the students, the instructors can adapt different teaching approaches to deal with appropriate group of students. Hence performance of students can be enhanced by applying different methods for each group of students. From the aforementioned case study of GPA, clustering has the potential to partition education data into appropriate groups and those groups can be used for further analysis to improve the overall education system. This application is simple to understand and exercise in a class at small level effectively.

The existing clustering methods i.e. K-means, K-medoids, DBSCAN, and AHCT are simulated using the same dataset of 600 students, and by considering the constraints described in the table 2.

TABLE 2. Clustering approaches with their parameters’ values.

Approach	Parameter settings
K-Means	No. of clusters ($k = 4$), number of iterations ($n = 50$)
K-Medoids	No. of clusters ($k = 4$), Predefined number of iterations ($n = 50$).
DBSCAN	Epsilon ($\epsilon = 0.5$): it defines the radius of neighborhood around the data-point “x”. Minimum points ($\text{minPts} = 10$): minimum number of neighbors within “eps”. Does not need to specify the number of clusters and iterations.
AHCT	No. of clusters ($k = 4$)
CFSFDP-HD	Does not need to specify explicitly, the number. of clusters and the number of iterations.

By comparing the recommended approach “CFSFDP-HD” with K-means, K-medoids, DBSCAN, and AHCT; the decision graph based approach “CFSFDP-HD” provides a deep insight to select potential clusters intuitively. In general practice, users run K-means and K-medoids more than 1000 times with various input settings (i.e. number of clusters and iterations) to get the meaningful clusters, however, the decision graph based approach used in CFSFDP-HD provides heuristics to get exact solutions within few iterations. While the DBSCAN approach shows some noisy data and also needs to define some explicit parameters’ values i.e. *Epsilon* “eps” and the *minimum points* “minPts”. The *Epsilon* defines the radius of the neighborhood around the data-point and the *minimum data-points* represent the minimum number of neighbors within the radius of the *Epsilon* value. In AHCT approach, there is also a need to mention the number of clusters explicitly. Furthermore, four distinct groups shown in the Figure 3(b) can easily be examined,

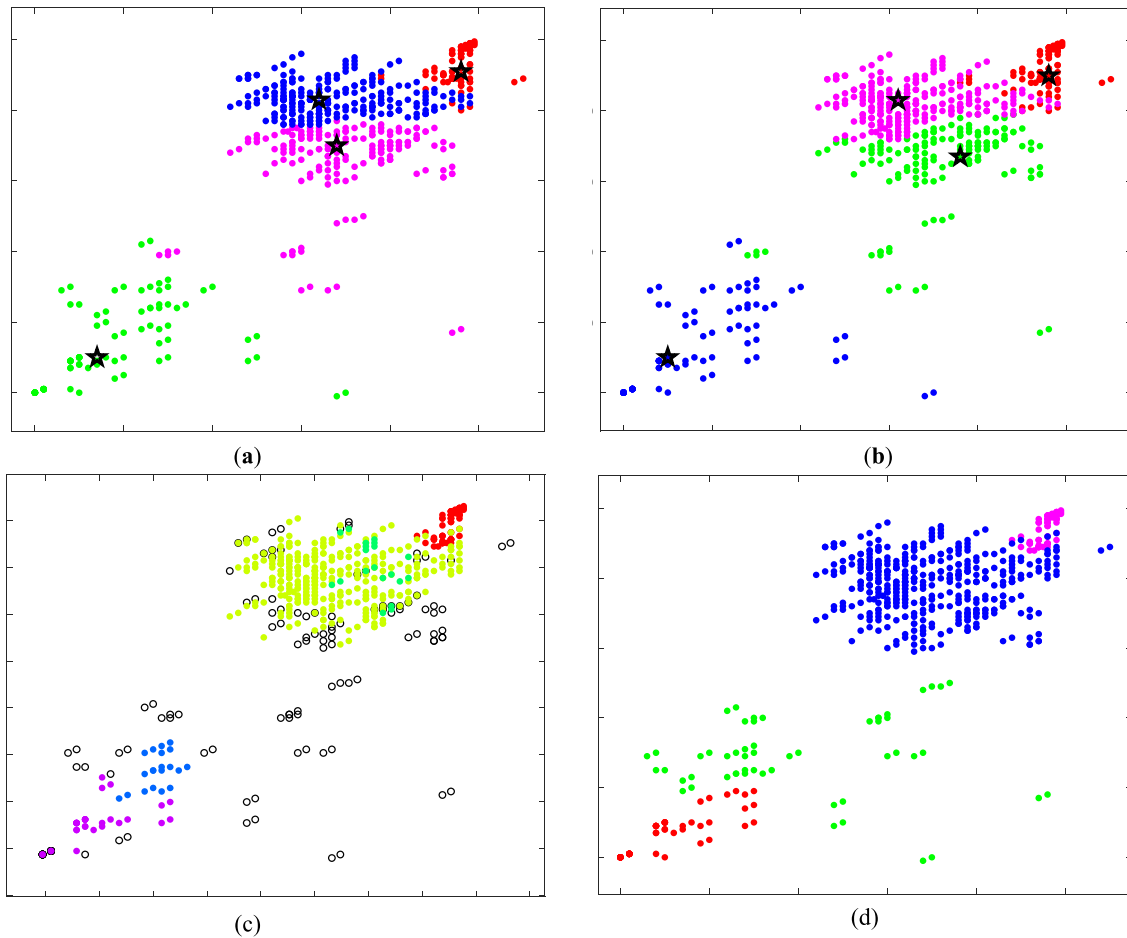


FIGURE 4. (a) K-means clustering results are visualized; the black-stars are the centroids of the clusters while the colored datapoints represent clusters (b) K-medoids clusters are shown; the black-stars are the medoids while different colors show different clusters. Both K-means and K-medoids are of similar nature and almost have similar results. (c) The DBSCAN clusters are represented with different colors while the noisy data-points are represented by the black-circles, (d) different colors are used to represent different clusters identified by the Agglomerative Hierarchical Cluster Tree.

visualized and compared with the K-means, K-medoids, DBSCAN, and AHCT in the Figure 4(a), 4(b), 4(c), 4(d) respectively. The representation of data-points in Figures 3(b) and Figure 4 differs because of their displayed-layout. The students having good grades are shown at left side in Figure 3(b) on the basis of ρ and δ values while the students having good grades are shown at top right of the graphs in Figure 4. In order to get appropriate clusters using the discussed approaches, users must have prior knowledge of existing clusters (number of clusters) and also unable to detect very low and below average students. This limitation makes it inappropriate to discover all intrinsic hidden patterns in data. To tackle technical issues, the CFSFDP-HD method is recommended to discover the existing patterns without knowing technical knowledge of the underlying data.

It is observed from the experiments that the CFSFDP-HD is more adaptive in nature and its results are more significant as compared with some of the existing approaches.

One of the goals of our research is to calculate the execution-time of the “CFSFDP-HD” approach and to determine whether optimization is required when accurately

predict different number of clusters from the large dataset. Figure 5 displays the average computation-time of K-means, K-medoids, DBSCAN, AHCT and the recommended “CFSFDP-HD” approaches, on an Intel®Core™CPU i5-7200U CPU @2.50GHz 2.71 GHz (2 processors), with 4GB RAM, 64 bits operating system and x64-based processor. The benchmark is completed by taking different numbers of students as: 600, 1200, 1800, 2400, 3000, and 3600. Each approach is executed multiple times, and took the average time for each run and shown in the Figure 5. The x-axis represents the number of students while the execution time (in seconds) is shown on the y-axis.

The result shown (in Figure 5) that the recommended approach is quite efficient as compared with the K-means and K-medoids. It also has taken the less execution time than AHCT on the smaller dataset. The DBSCAN is faster than the recommended approach on the smaller dataset. The recommended approach is less efficient than the DBSCAN and AHCT, when the dataset is large. To measure the accurate behavior, the presented approach “CFSFDP-HD” is to be executed on the large and real datasets.

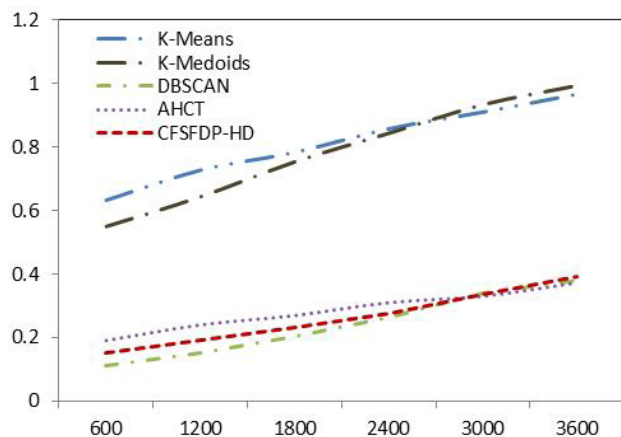


FIGURE 5. Execution time of K-means, K-medoids, DBSCAN, AHCT and CFSFDP-HD against different number of students.

V. CONCLUSIONS

The data mining approaches provide a sense of intelligence in existing e-learning systems, efficiently and effectively. This paper has presented an adaptive data mining clustering approach “CFSFDP-HD” that is integrated with the conceptual personalized e-learning system architecture. It has been observed from the literature that traditional e-learning systems are mostly query-based and the queries are responded without any intelligence or heuristics. The potential application of clustering in educational big data has also been examined. The existing discussed clustering approaches are suitable in clustering educational data where the number of cluster are known and face drawbacks when applied to unknown cluster sizes. It has been evaluated that the recommended data mining clustering approach is efficacious in analyzing the big data to make education systems robust. Results show that the recommended approach has formed accurate clusters and also took less execution time as compared with some of the existing approaches. The recommended approach also has the potential to solve the challenges of interdisciplinary research, emotional learning, and e-learning in the field of education.

The data mining approaches can further be improved to generate knowledge and provide intelligent assistance to the students. Real and larger datasets can be simulated to analyze the behavior of the recommended data mining approach. The learning capabilities of the students can be further improved by introducing the intelligent games and integrating the recommended approaches. Student collaboration is an important aspect of learning by group discussion and by sharing personal thoughts: Intelligent techniques can be introduced in different students’ groups with significant attributes for problem solving.

REFERENCES

- [1] W.-C. Wong and A. W.-C. Fu, “Incremental document clustering for Web page classification,” in *Enabling Society with Information Technology*. Tokyo, Japan: Springer, 2002, pp. 101–110.
- [2] R. S. J. D. Baker and K. Yacef, “The state of educational data mining in 2009: A review and future visions,” *J. Edu. Data Mining*, vol. 1, no. 1, pp. 3–17, 2009.
- [3] L.-Y. Li and Y.-L. Zheng, “The application of the Internet of Things in education,” *Mod. Educ. Technol.*, vol. 2, no. 5, 2010. [Online]. Available: http://en.cnki.com.cn/Article_en/CJFDTotat-XJJS201002005.htm
- [4] R. Baker, “Data mining for education,” *Int. Encyclopedia Edu.*, vol. 7, no. 3, pp. 112–118, 2010.
- [5] T. Blanco, R. Casas, E. Manchado-Pérez, Á. Asensio, and J. M. López-Pérez, “From the islands of knowledge to a shared understanding: Interdisciplinarity and technology literacy for innovation in smart electronic product design,” *Int. J. Technol. Des. Educ.*, vol. 27, no. 2, pp. 329–362, 2017.
- [6] G. Siemens and P. Long, “Penetrating the fog: Analytics in learning and education,” *EDUCAUSE Rev.*, vol. 46, no. 5, p. 30, 2011.
- [7] D. Howe et al., “Big data: The future of biocuration,” *Nature*, vol. 455, no. 7209, pp. 47–50, 2008.
- [8] G.-H. Kim, S. Trimi, and J.-H. Chung, “Big-data applications in the government sector,” *Commun. ACM*, vol. 57, no. 3, pp. 78–85, 2014.
- [9] M. Chen, S. Mao, Y. Zhang, and V. C. M. Leung, “Big data applications,” in *Big Data*. Cham, Switzerland: Springer, 2014, pp. 59–79.
- [10] C. L. P. Chen and C.-Y. Zhang, “Data-intensive applications, challenges, techniques and technologies: A survey on big data,” *Inf. Sci.*, vol. 275, pp. 314–347, Aug. 2014.
- [11] N. Noury et al., “Monitoring behavior in home using a smart fall sensor and position sensors,” in *Proc. 1st Annu. Int., Conf. Microtechnol. Med. Biol.*, 2000, pp. 607–610.
- [12] R. Bie, R. Mehmood, S. Ruan, Y. Sun, and H. Dawood, “Adaptive fuzzy clustering by fast search and find of density peaks,” *Pers. Ubiquitous Comput.*, vol. 20, no. 5, pp. 785–793, 2016.
- [13] G. Qian, Y. Wu, D. Ferrari, P. Qiao, and F. Hollande, “Semisupervised clustering by iterative partition and regression with neuroscience applications,” *Comput. Intell. Neurosci.*, vol. 2016, Mar. 2016, Art. no. 4037380.
- [14] U. Markowska-Kaczmar, H. Kwasnicka, and M. Paradowski, “Intelligent techniques in personalization of learning in e-learning systems,” in *Computational Intelligence for Technology Enhanced Learning*. Berlin, Germany: Springer, 2010, pp. 1–23.
- [15] M. Cordeiro and J. Gama, “Online social networks event detection: A survey,” in *Solving Large Scale Learning Tasks. Challenges and Algorithms*. Cham, Switzerland: Springer, 2016, pp. 1–41.
- [16] G. H. Shah, C. K. Bhensadadia, and A. P. Ganatra, “An empirical evaluation of density-based clustering techniques,” *Int. J. Soft Comput. Eng.*, vol. 2, no. 1, pp. 216–223, 2012.
- [17] S. Engström, “Differences and similarities between female students and male students that succeed within higher technical education: Profiles emerge through the use of cluster analysis,” *Int. J. Technol. Des. Edu.*, vol. 28, no. 1, pp. 239–261, 2018.
- [18] P. Nerurkar, A. Shirke, M. Chandane, and S. Bhirud, “Empirical analysis of data clustering algorithms,” *Procedia Comput. Sci.*, vol. 125, pp. 770–779, Jan. 2018.
- [19] J. MacQueen, “Some methods for classification and analysis of multivariate observations,” in *Proc. 5th Berkeley Symp. Math. Statist. Probab.*, vol. 1, 1967, pp. 281–297.
- [20] R. Mehmood, G. Zhang, R. Bie, H. Dawood, and H. Ahmad, “Clustering by fast search and find of density peaks via heat diffusion,” *Neurocomputing*, vol. 208, pp. 210–217, Oct. 2016.
- [21] M. Ester, H.-P. Kriegel, J. Sander, and X. Xu, “A density-based algorithm for discovering clusters in large spatial databases with noise,” in *Proc. KDD*, 1996, pp. 226–231.
- [22] A. Rodriguez and A. Laio, “Clustering by fast search and find of density peaks,” *Science*, vol. 344, no. 6191, pp. 1492–1496, Jun. 2014.
- [23] L. Zappia and A. Oshlack, “Clustering trees: A visualisation for evaluating clusterings at multiple resolutions,” *BioRxiv*, pp. 1–8, May 2018.
- [24] A. R. Anaya and J. G. Boticario, “A data mining approach to reveal representative collaboration indicators in open collaboration frameworks,” in *Proc. Int. Work. Group Educ. Data Mining*, 2009, pp. 210–219.
- [25] A. M. de Moraes, J. M. F. R. Araújo, and E. B. Costa, “Monitoring student performance using data clustering and predictive modelling,” in *Proc. IEEE Frontiers Edu. Conf. (FIE)*, Oct. 2014, pp. 1–8.
- [26] B. R. Prakash, M. Hanumanthappa, and V. Kavitha, “Big data in educational data mining and learning analytics,” *Int. J. Innov. Res. Comput. Commun. Eng.*, vol. 2, no. 12, pp. 7515–7520, 2014.
- [27] A. Algarni, “Data mining in education,” *Int. J. Adv. Comput. Sci. Appl.*, vol. 7, no. 6, pp. 456–461, 2016.
- [28] A. Bovo, S. Sanchez, O. Héguay, and Y. Duthen, “Analysis of students clustering results based on Moodle log data,” in *Proc. 6th Int. Conf. Educ. Data Mining (EDM)*, Memphis, TN, USA, 2013, pp. 306–307.

- [29] A. A. Saa, "Educational data mining & students' performance prediction," *Int. J. Adv. Comput. Sci. Appl.*, vol. 7, no. 5, pp. 212–220, 2016.
- [30] T. Agasisti and A. J. Bowers, "9. Data analytics and decision making in education: Towards the educational data Scientist as a key actor in schools and higher education institutions," in *Handbook of Contemporary Education Economics*, 2017, p. 184.
- [31] V. Ivančević, M. Čeliković, and I. Luković, "The individual stability of student spatial deployment and its implications," in *Proc. Int. Symp. Comput. Edu. (SIIE)*, Oct. 2012, pp. 1–4.
- [32] K. Ying, M. Chang, A. F. Chiarella, Kinshuk, and J.-S. Heh, "Clustering students based on their annotations of a digital text," in *Proc. IEEE 4th Int. Conf. Technol. Educ. (T4E)*, Jul. 2012, pp. 20–25.
- [33] K. L. N. Eranki and K. M. Moudgalya, "Evaluation of Web based behavioral interventions using spoken tutorials," in *Proc. IEEE 4th Int. Conf. Technol. Educ. (T4E)*, Jul. 2012, pp. 38–45.
- [34] A. S. Drigas and P. Leliopoulos, "The use of big data in education," *Int. J. Comput. Sci. Issues*, vol. 11, no. 5, p. 58, 2014.
- [35] A. Manohar, P. Gupta, V. Priyanka, and M. F. Uddin, "Utilizing big data analytics to Improve education," in *Proc. ASEE North-east Sect. Conf.*, Kingston, RI, USA, 2016. [Online]. Available: <https://scholarworks.bridgeport.edu/xmlui/handle/123456789/1616>
- [36] B. Tulasi, "Significance of big data and analytics in higher education," *Int. J. Comput. Appl.*, vol. 68, no. 14, pp. 21–23, 2013.
- [37] B. Daniel, "Big Data and analytics in higher education: Opportunities and challenges," *Brit. J. Educ. Technol.*, vol. 46, no. 5, pp. 904–920, 2015.
- [38] C. Dede, "Next steps for 'big data' in education: Utilizing data-intensive research," *Educ. Technol.*, vol. 56, no. 2, pp. 37–42, 2016.
- [39] M. Gaeta, S. Miranda, F. Orciuoli, S. Paolozzi, and A. Poce, "An approach to personalized e-learning," *J. Edu., Informat.*, vol. 11, no. 1, pp. 15–21, 2013.
- [40] A. K. Prawira, T. D. Sofianti, and Y. Indrayadi, "Developing E-learning system to support teaching and learning activities using DSDM approach," *PERFORMA: Media Ilmiah Teknik Industri*, vol. 14, no. 1, pp. 41–52, 2015.
- [41] M. Yarandi, H. Jahankhani, and A.-R. Tawil, "A personalized adaptive e-learning approach based on semantic Web technology," *Webology*, vol. 10, no. 2, 2013, Art. no. 110.
- [42] A. Maselena, N. Sabani, M. Huda, R. Ahmad, K. A. Jasmi, and B. Basiron, "Demy stifying learning analytics in personalised learning," *Int. J. Eng. Technol.*, vol. 7, no. 3, pp. 1124–1129, 2018.
- [43] C.-H. Wu, Y.-S. Chen, and T.-C. Chen, "An adaptive e-learning system for enhancing learning performance: Based on dynamic scaffolding theory," *EURASIA J. Math., Sci. Technol. Educ.*, vol. 14, no. 3, pp. 903–913, 2017.
- [44] Z. Liu and Y. Liu, "Research on personalization e-learning system based on agent technology," in *Proc. 3rd WSEAS Int. Conf. Circuits, Syst., Signal Telecommun.*, Ningbo, China, 2009, pp. 110–114.
- [45] X. Li and S.-K. Chang, "A personalized e-learning system based on user profile constructed using information fusion," in *Proc. DMS*, 2005, pp. 109–114.
- [46] M. Hedayetull, I. Shovon, and M. Haque, "An approach of improving student's academic performance by using k-means clustering algorithm and decision tree," *Int. J. Adv. Comput. Sci. Appl.*, vol. 3, no. 8, pp. 146–149, 2012.
- [47] L. Shen, M. Wang, and R. Shen, "Affective e-learning: Using 'emotional' data to improve learning in pervasive learning environment," *J. Educ. Technol. Soc.*, vol. 12, no. 2, pp. 176–189, 2009.
- [48] A. Dutt, M. A. Ismail, and T. Herawan, "A systematic review on educational data mining," *IEEE Access*, vol. 5, pp. 15991–16005, 2017.
- [49] C. Romero, M.-I. López, J.-M. Luna, and S. Ventura, "Predicting students' final performance from participation in on-line discussion forums," *Comput. Educ.*, vol. 68, pp. 458–472, Oct. 2013.
- [50] W.-C. Chang, S.-L. Chen, M.-F. Li, and J.-Y. Chiu, "Integrating IRT to clustering student's ability with K-means," in *Proc. 4th Int. Conf. Innov. Comput., Inf. Control (ICICIC)*, Dec. 2009, pp. 1045–1048.
- [51] P. Džáždilová, J. Martinovic, K. Slaninová, and V. Snašel, "Analysis of Relations in eLearning," in *Proc. IEEE/WIC/ACM Int. Conf. Web Intell., Intell. Agent Technol. (WI-IAT)*, vol. 3, Dec. 2008, pp. 373–376.
- [52] G. Cobo, D. García-Solórzano, E. Santamaría, J. A. Morán, J. Melenchón, and C. Monzo, "Modeling students' activity in online discussion forums: A strategy based on time series and agglomerative hierarchical clustering," in *Proc. EDM*, 2011, pp. 253–258.
- [53] L. Kaufman and P. J. Rousseeuw, "Clustering by means of medoids," in *The Statistical Data Analysis Based on the L1-Norm and Related Methods*. Amsterdam, The Netherlands: North Holland, 1987.
- [54] L. Rokach and O. Maimon, "Clustering methods," in *Data Mining and Knowledge Discovery Handbook*, O. Maimon and L. Rokach, Eds. Boston, MA, USA: Springer, 2005, pp. 321–352.
- [55] F. Esposito, O. Licchelli, and G. Semeraro, "Extraction of user profiles in e-learning systems," in *Proc. I-KNOW*, Graz, Austria, 2003, pp. 238–243.
- [56] P. Gomes, B. Antunes, L. Rodrigues, A. Santos, J. Barbeira, and R. Carvalho, "Using ontologies for elearning personalization," *Commun. Cognition*, vol. 41, no. 1, p. 127, 2008.



SAMINA KAUSAR received the M.S. degree in computer science from the International Islamic University Islamabad, Pakistan, in 2007. She is currently working as a Researcher and also a Ph.D. Scholar with the School of Computer Engineering and Science, Shanghai University, China. She has been working as an Assistant Professor with the University of Management Sciences & Information Technology Kotli, Azad Kashmir, Pakistan. Her research interests are in the fields of big data,

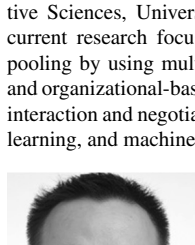
bioinformatics, computer networks, cloud computing, data mining, and machine learning algorithms.



XU HUAHU is currently a Doctoral Supervisor and a Professor with the School of Computer Engineering and Science, Shanghai University, where he is also the Director of the Information Office of Shanghai University. He is the Chairman of the Shanghai Security and Technology Association. His research interests include multimedia technology, CIMS, and computer network technology.



IFTIKHAR HUSSAIN received the Ph.D. degree from Universiteit Hasselt, Belgium, in 2017, and graduated in computer science, specialized in software engineering from Iqra University, Islamabad Campus, in 2009. He is currently working as an Assistant Professor with the School of Computer & IT, Beaconhouse National University, Lahore, Pakistan. After his graduation, he has worked for three years as a Lecturer with the Department of CS & IT, Faculty of Administrative Sciences, University of Kotli Azad Jammu and Kashmir, Kotli. His current research focuses on the coordination and negotiation in the car-pooling by using multi-agents. He has research experience in agent-based and organizational-based modeling, multi-agents in transportation, modeling interaction and negotiation, software engineering, data mining, IoT, machine learning, and machine algorithms.



ZHU WENHAO was born in 1979. He received the bachelor's, master's, and Ph.D. degrees from Zhejiang University in 2002, 2006, and 2009, respectively. He was a Research Scholar at the Computer Laboratory, University of Cambridge, from 2012 to 2013. He is currently working as an Associate Professor with the School of Computer Engineering and Science, Shanghai University, China. His research is in the areas of text representation, information extraction, and web data mining.



MISHA ZAHID is currently pursuing a bachelor's degree in software engineering. She is working as Teacher's Assistant and also an undergraduate student at Beaconhouse National University, Lahore, Pakistan. She is also the Founder and the CEO of a software house: Tech Designers. She has research interests in the field of e-learning, data-mining, Internet of Things, and machine learning.