

Review Paper on Educational Data Mining

Mr. Pradeep Nayak¹, Mohammed Sufiyan², Mohan Raju. V³, Monisha. N. S.⁴, Moollya Gautami Bhaskar⁵

Assistant Professor, Department of Information Science and Engineering¹

Students, Department of Information Science and Engineering^{2,3,4,5}

Alva's Institute of Engineering and Technology, Mijar, Mangalore, Karnataka, India

Abstract: Education and computer science are both involved in the burgeoning inter-disciplinary research field known as Educational Data Mining (EDM). EDM uses data mining software and ways to extract meaningful and practical data from big educational databases. EDM introduces better and more efficient learning techniques in an effort to enhance educational processes. The term "EDM methods" refers to a group of techniques for creating models and applications. This page provides a thorough literature review on EDM techniques. The essay also covers EDM research problems and trends. This EDM insight aims to provide researchers interested in furthering the field of EDM with useful and valuable information.

Keywords: Educational Data Mining (EDM); Prediction; Relationship Mining; Structure Discovery are some keywords

I. INTRODUCTION

Educational Data Mining (EDM) is the application of data mining tools and techniques to data collected during an educational activity. Prediction, association rule mining, clustering, and outlier analysis are some of the most important tasks in data mining[1]. Data mining methods are applied to

data from various domains. The datasets for EDM are drawn from the educational domain. EDM extracts useful information from large datasets related to education. EDM aids in dealing with educational environment and process issues and problems.

EDM is defined by the International Educational Data Mining Society (IEDMS) as "an emerging discipline concerned with developing methods for exploring the unique and increasingly large-scale data that come from educational settings, and using those methods to better understand students and the settings in which they learn." [2]

Education is one of the most important factors in a society's development. Learning methods that provide slower or ineffective education have an impact on society and the country's development. EDM strives to provide better and new methods of learning that are more effective than current methods of education delivery.

EDM research has examined important issues such as student dropout or poor academic performance and recommended solutions to help minimise such issues.

The following EDM research objectives are established by the literature:

- 1) To investigate 'Big Data' in the field of education.
- 2) To address issues concerning education.
- 3) To increase the learner's knowledge and learning.
- 4) Assisting students in achieving success.
- 5) To assist institutions in becoming successful.

EDM has its own set of data mining methods that are only used on educational data. The term 'EDM Methods' refers to a collection of all data mining techniques used to mine educational data. The collected educational data is pre-processed before EDM methods are used to extract useful information. EDM methods can be classified, as shown in figure 1, based on classifications proposed by various authors in [3]-[9]. The following are the six major divisions:

- 1) Data distillation for human judgement
- 2) Methods of prediction
- 3) Relationship mining techniques
- 4) Structure discovery techniques
- 5) Modeling for discovery
- 6) A variety of other methods

By reviewing significant contributions from the last ten years, this article comprehensively surveys all of the methods described in the EDM literature. The remainder of the article is structured as follows. Section II discusses the techniques used to distil data for human judgement. Section III goes over the various prediction methods used in EDM. Section IV goes over the various relationship mining techniques that are used in EDM. Section V examines the techniques used in EDM to discover structure in educational datasets. Section VI provides an overview of the models used in educational data mining. Section VII discusses the remaining methods used in EDM that could not be classified. Section VIII discusses the research trend and future research directions in EDM. Finally, section IX brings the article to a close.

II. DATA FILTRATION FOR HUMAN ASSESSMENT

Data distillation for human judgement entails highlighting useful information through data representation. It is possible to accomplish this through the use of summarization, visualisation, and interactive interfaces[10]. The most common tools for effective visualisation are basic statistical techniques and graph plotting methods. [11]Humans can identify patterns in visualised data, allowing for faster decision making[6]. Shukla[12] used statistical techniques to demonstrate the relationship between the symbols used in annotated academic resources and the associated sentiments for their use. Merceron and Yacef[13] demonstrated statistical analysis of the formed clusters using visualisation techniques in conjunction with clustering

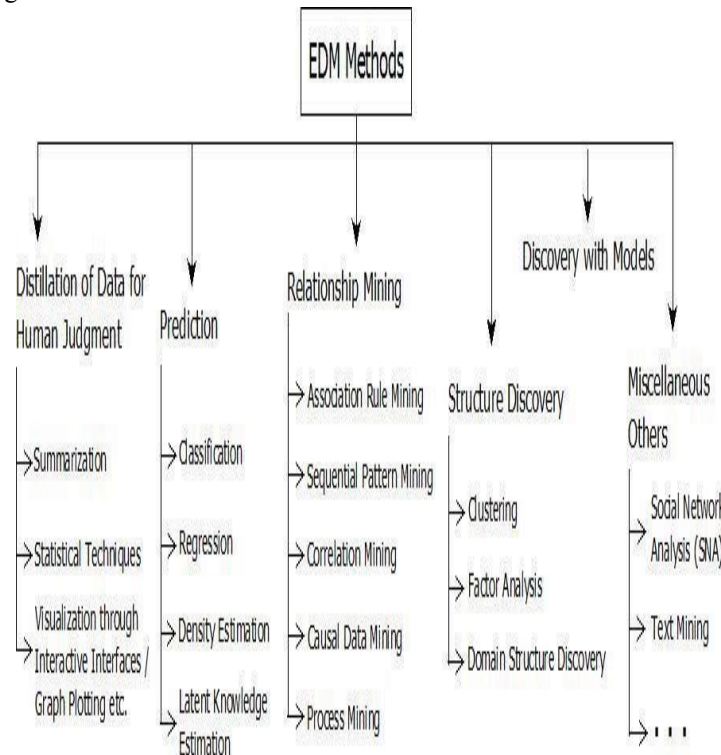


Figure 1: Categorization of EDM methods

Romero. [14] made use of a GISMO tool that can be installed in the Moodle system. The GISMO tool generated pictorial and graphical representations of Moodle logs, which instructors used to analyse various aspects of students enrolled in a distance education course Paiva et al. [15] used the RAG colours technique to create three visualisations: a segmented bar graph, ordered weights, and combined interactions. The visualisations effectively displayed the results of EDM operations performed on assessment data from a high school online mathematics course.

III. PROBABILITY TECHNIQUES

Prediction methods are the most commonly used techniques for dealing with issues such as student dropout or poor academic performance. Initially, prediction methods construct the model for the attribute whose value is to be predicted

based on the values of the other attributes in the dataset. This step is known as training, and it is performed on a dataset with known values of the attribute to be predicted. The values of the attribute are predicted for another dataset that may not have known values for the attribute to be predicted based on the model created in the initial phase.

The prediction methods used in EDM are discussed next.

3.1 Category

The variable whose value is to be predicted by classification is either categorical or binary. The predicted attribute values are class labels. Using k-fold validations, classification methods categorise the records in the dataset based on the class-labels. The classification methods used in the EDM literature are crude, decision rules, random forest, decision trees, step regression, logistic regression, statistical, rule induction, neural networks, and fuzzy rule learning classifiers [4], [13], [16], [17], [18]-[29]. These classifiers have algorithms (along with variants) that are executed on the dataset; for example, the most commonly used Decision Tree classifier algorithms are ID3, C4.5, and J48. Any classifier can be either a white-box or a black-box classifier.

Black-box classifiers do not reveal the classification mechanism used, but they are more accurate than white-box classifiers such as Neural Networks (NN), k-Nearest Neighbor (kNN), Support Vector Machine (SVM), and so on. White-box classifiers, such as decision tree classifiers and rule induction classifiers, produce the classification rules used in the classification process.

3.2 Regression

A continuous variable is the attribute whose value is to be predicted using regression. EDM employs regression techniques such as linear regression, regression trees, and locally weighted linear regression [4],[6]. Romero and Ventura[6] used regression models to predict student grades.

However, regression techniques are rarely used to predict continuous variables in EDM applications. In the last decade, no significant research work has used regression alone for educational data mining. Zimmermann et al. [30] developed a model for predicting students' performance in graduation (computer science) courses taught at ETH Zurich, Switzerland.

Strechetal.[31] demonstrated that regression algorithms are ineffective for modelling student performance when compared to classification algorithms. Al Hammadi and Aksoy[27] and Beal et al.[32] used logistic regression to predict student performance. It should be noted that logistic regression is a regression model that predicts the value of a binary attribute. As a result, it only serves as a classifier.

3.3 Density Estimation

Density estimation is a statistical technique for forecasting. Density estimators forecast an attribute's value based on the probability density function of that attribute hidden in the given dataset. Parzen windows, vector quantization, and density-based clustering methods are examples of prominent density estimators. Density estimation methods are based on kernel functions (including the Gaussian function)[8]. Ocumpaugh et al.[33] used statistical probability estimators - Cohen's Kappa and A - to perform population validity for classification models on students from urban, suburban, and rural areas. Minaei-Bidgoliet al.[16] used Parzen windows in conjunction with other classifiers to predict student performance.

3.4 Latent Knowledge Estimation

Latent knowledge estimation is another statistical technique used for forecasting. Latent knowledge estimators evaluate students' abilities based on their responses to a problem-solving exercise. A student's responses are mapped with the necessary skills for a problem.

The discovered correctness or incorrectness pattern for various skills estimates a student's knowledge of that skill. Nave Bayes, Bayes net, Bayesian Knowledge Tracing (BKT), and performance factor assessment are popular methods for estimating latent knowledge [4],[6]. BKT was used by Corbett and Anderson [34] to model the student's knowledge states. In [19],[22],[23],[27],[31], Nave Bayes was used as a classifier alongside other classifiers for a comparative study of classification algorithms. Lopez and co.

[24] classified using both BayesNet and naive Bayes. Sethi and Singh[29] enhanced naive Bayes with additional parameters such as 'support' and 'confidence' to predict with an improved naive Bayes algorithm.

IV. RELATIONSHIP MINING METHODS

Relationship mining methods detect associations between attributes in a dataset. These methods are used to discover relationships between educational factors in a dataset. Some enigmatic relationships are also discovered, which may be useful for improving a learning methodology or an educational system. The following section discusses the relationship mining methods used in the field of EDM.

4.1 Association Rule Mining (ARM)

ARM is a widely used technique in data mining that produces association rules that relate the attributes of a dataset. The most commonly used algorithms for association rule mining are Apriori, ECLAT, and FP-Growth [35]. ARM has been used to detect associations between subjects (strong or weak)[36-38]. The FP-Growth algorithm [39] was used to detect associations within various questions. Shi et al.[40] managed the curriculum using Apriori[40] based on the detected relationships. DeCarlo and Rizk[41] proposed an expert system prototype based on rules obtained from the application of ARM on a dataset that recorded the effectiveness of the course materials and students' exam performance. Luna and colleagues [42] examined student Moodle usage data to discover rare association rules corresponding to unusual student grades. The G3P-Rare algorithm was used, which combined ARM and genetic programming.

G3PARM is a new technique developed by Romero et al. [43] that combines ARM and Grammar Guided Genetic Programming (G3P). The G3PARM algorithm was used to detect association rules in the evaluation data of a multiple-choice quiz, which served as student feedback.

Dimi et al.[44] improved the e-testing process using the Predictive Apriori algorithm and the data structure used in [43]. Aleem and Gore[35] extracted confidence-based feedback from MCQ evaluation data by running FP-Growth, ECLAT, and Apriori algorithms. Aleem et al. [45] used ARM to detect the relationships between common grammatical errors made by researchers while writing manuscripts. Deng and Que[46] used ARM on student feedback for teachers to discover relationships useful for teaching assessment. Aleem and Gore[47] used ARM to detect examinees' answering patterns for optional questions in an engineering semester examination.

4.2 Sequential Pattern Mining (SPM)

SPM is used to discover relationships between the attributes of a sequential dataset[48]. The sequential dataset is a database that contains sequence. One such sequence includes the values of an attribute at various points in time. Many fields, including medical science and web technology, store data as a series of symbols. A sequential rule is any mined relationship from a sequential database that is a subsequence of the dataset.

The well-known sequential pattern mining algorithms are Generalized Sequential Pattern (GSP), AprioriAll, Clospan, PrefixSpan, FreSpan, and Sequential PAttern Discovery using Equivalence classes (SPADE). SPM methods are used in EDM to detect periodic changes in student behaviour, build student models, and explore concept maps to discover valuable learning patterns. [4],[6]. Romero et al. [14] used the PrefixSpan algorithm on Moodle logs from various forums to find associations between the answers, which were saved as sequences.

4.3 Correlation Mining

Correlation mining is a statistical method for calculating the correlation between different attributes in a dataset. The correlation value ranges from -1 to +1 depending on the direction and closeness of the relationship between the two sets of attributes being examined. A correlation value of +1 (or -1) indicates that both sets of attributes are perfectly related (or inversely associated). A correlation value close to zero indicates a weaker association. The established statistical formulas are used to calculate correlation.

Correlation analysis is frequently used in conjunction with other methods to achieve better results. Tissera et al. [36] calculated the strength of the discovered association between two university-level subjects using the Pearson correlation coefficient.

The discovered relationships were discovered by applying ARM to the grades of ICT students in Sri Lanka. The Pearson correlation coefficient has also been used to investigate correlations in datasets used for ARM and clustering techniques [39].

4.4 Causal Data Mining

Causal data mining combines statistical techniques with association rule mining to determine the cause of a discovered relationship in large datasets[49]. Causal Data Mining is a broad ARM technique. Association rule mining discovers the associational dependency between a dataset's attributes but does not investigate the cause of the dependency. If the cause of the relationship is not investigated, the association rules obtained from one dataset may differ in another. A dataset-discovered relationship may be causally related to another relationship. Causal data mining aids in the exploration of all such dependencies among the relationships mined in a dataset and provides more reliable and generalizable association rules.

Students' poor performance can be effectively analysed using causal data mining techniques[4]. Li et al. [49] used a retrospective cohort study in conjunction with ARM to determine causal association rules. Any association rule that has an odds ratio greater than one represents a causal association. To detect the causal relationship, local causal discovery (LCD) algorithms could be used. Mani and Cooper[50] used LCD techniques in the Bayesian framework and checked the causal faithfulness condition to determine the association's dependency. The causal Markov condition, on the other hand, specified whether or not the association was independent.

4.5 Process Mining

Educational process mining is a new technology that incorporates process mining techniques into EDM[51, 52]. Educational process mining is an extension of ARM techniques that improves on the model created by other EDM methods. The model is designed to represent the process of completing an educational task, such as students taking an online quiz. Educational process mining methods mine event logs for knowledge about related processes. Educational process mining provides a clear visualisation of the entire process and aids in the improvement of educational processes[6]. Cairns et al.[51] investigated the interaction process of trainers with the used training mate using the Phidias platform, a two-step clustering method, and social mining techniques.

[53] used a process mining technique to discover the behavioural patterns of students participating in game-based learning. The activity logs generated while remotely playing educational games for making decisions and interacting with players were used for testing purposes.

V. STRUCTURE DISCOVERY METHODS

Structure discovery methods find structure in a dataset using an unsupervised classification mechanism. The dataset's data items are classified based on attribute values rather than class labels, as in prediction methods. After prediction methods, structure discovery methods are the most commonly used EDM methods. The following section discusses the structure discovery methods used in EDM.

5.1 Clustering

Clustering methods group together data items with similar attributes, dividing the entire dataset into clusters. Distance measures such as Manhattan, Euclidean, and others are used to assess the similarity/dissimilarity of items[54]. Clustering is used in EDM to group students who have similar intelligence and behaviour. Clustering is sometimes used to group similar learning materials of a subject[55]. Educational Data Clustering is the method of using clustering in educational data mining (EDC).

EDC is also used to analyse learning styles and improve collaborative learning among students[56]. K-Means[13],[14],[24],[39],[57], X-means[24], Hierarchical-Clusterer[13], [24],[32], Expectation-Maximization (EM)[24],[57], neural clustering [28], demographic clustering[28], and FarthestFirst[24] are clustering methods used in the literature.

Chalariset al.[39] used clustering, ARM, and correlation to identify knowledge that could be used to improve the quality of educational processes. Lopez et al. [24] ran clustering algorithms EM, Farthest-First, Hierarchical Clusterer, Simple-KMeans, and X-means on Moodle forum usage data and compared them to various classifiers. The EM

algorithm outperformed and outperformed the other classifiers in terms of accuracy. Beal et al.[32] classified learning styles using hierarchical cluster analysis on learner motivational data and problem-solving logs gathered from ITS. Delavari and PhonAmnuaisuk[28] used neural and demographic clustering methods to group students based on academic and demographic information provided by lecturers. K means and hierarchical clustering were used by Merceron and Yacef[13] to identify the specific behaviour of students who lacked problem-solving ability. Ramanathan et al. [58] used Bayesian fuzzy clustering and a Lion-Wolf-based Deep Belief Network (LW-DBN) in a distributed architecture to predict student performance in higher education.

5.2 Factor Analysis

Factor analysis is a dimensionality reduction utility that employs the statistical technique. There are numerous attributes from various fields in educationally related data, such as family background, personal details, demographic details, academic details, social details, cultural details, psychometric details, and so on. This is known as the One Thousand Factors Problem[21] by the researchers. To address this issue, relationships between directly observable attributes are examined in order to identify latent factors. Variations in the observable attributes can be reflected in fewer latent attributes, reducing the number of attributes. Following the discovery of the latent attributes, the dataset is clustered based on the derived latent factors.

The most widely used technique for factor analysis is Principle Component Analysis (PCA). Following factor analysis, other EDM techniques are applied to the reduced dataset. Goyal and Vohra[59] suggested performing factor analysis on the dataset before employing EDM techniques such as classification, regression, clustering, and association rule mining.

5.3 Domain Structure Discovery

Domain structure discovery methods map dataset attributes to domain knowledge states. Domain structure discovery is used in EDM to assess skill expertise by mapping key educational concepts (knowledge states) to the structure of tasks performed by students. It is best expressed as follows: Assume a basic arithmetic questionnaire contains four questions. The first two questions involve addition, the third question involves subtraction, and the fourth question involves number multiplication. In this case, the task structure would be the student's marks for the answers to the four questions.

The key educational concept for the first two questions is addition, subtraction for the third, and multiplication for the last. As a result, the first two questions' scores would correspond to the first key educational concept - addition, the third question's scores would correspond to the second key educational concept - subtraction, and the fourth question's scores would correspond to the third key educational concept - multiplication. A student's knowledge of addition, subtraction, and multiplication can be demonstrated by analysing the mapped concept-skill matrix. The Q-Matrix is a well-known method for mapping key educational concepts to students' skills[4]. QMatrix stores the degree to which a question/item is associated with various concepts. Non-negative By decomposing a matrix of positive numbers into two smaller matrices, the Matrix Factorization (NMF) technique allows for a straightforward interpretation in terms of Q-matrix[6]. Desmarais[60] used NMF, followed by clustering on quiz evaluation data for four subjects, to visualise the students' skill sets. Romero et al. [43] proposed a method for providing feedback to instructors based on quiz data evaluation using the G3PARM algorithm. The knowledge matrix is obtained by multiplying the score matrix of students' quiz marks by the relationship matrix during the data preprocessing stage. The relationship matrix is a Q-Matrix that contains the degree to which quiz questions are associated with the related concepts. To reach further conclusions, the G3PARM algorithm was applied to the knowledge matrix.

VI. DISCOVERY WITH MODELS

The model's technique for discovery entails two stages of data analysis. In the first phase, either educational data mining methods or human reasoning are used to build a model[4].

To construct the model, EDM methods such as prediction, relationship mining, and clustering can be used. In the second phase of complex attribute analysis, the model obtained in the first phase is used as a component/model. Delavari and Phon-Amnuaisuk[28] proposed DM-HEDU, an acronym for Data Mining in Higher Education Systems, a model that performs classification, clustering, and association rule algorithms on academic and demographic data of

students and lecturers for various education-related activities in order to improve the higher education system. Al-Twijri and Noaman[61] proposed a new model, DMAM (Data Mining Admission Model), to aid in decision making for the admission process in Saudi universities.

The classification rules developed using statistical tools and human knowledge are implemented by DMAM. Zimmermann et al.[30] used a model built with regression methods and various variable selection methods to predict Master's course performance using Bachelor's performance parameters. Jeong and Biswas[62] created a learning-by-teaching model in which students teach a computer agent.

VII. MISCELLANEOUS METHODS

Miscellaneous methods are those that have been used less frequently in the EDM literature and do not fit into either of the previously mentioned categories.

One such technique is text mining, which processes unstructured text data rather than structured relational data.

The majority of unstructured web data is mined using text mining.

Text mining has given researchers a tool for analysing the content of web pages and forums[10]. Romero et al. [14] used the text-mining tool KEA to search for specific key phrases in Moodle forum posts. Another method is social network analysis (SNA), which analyses the structure of a social network and interprets the relationships between nodes. SNA is used to examine educational processes involving collaborative learning [10].

Many other methods and procedures, such as optimization techniques[16], genetic algorithms[21],[22],[26], semi-supervised learning methods[63], active learning methods[64], HMM[62], NMF[60], and so on, have been used in EDM.

VIII. RESEARCH TREND AND DIRECTIONS FOR FUTURE RESEARCH

The most commonly used technique is classification, followed by clustering methods. The most commonly used EDM methods for discovering structure in educational data are clustering techniques. In the literature, association rule mining methods have been used less frequently than clustering methods. Data distillation for human judgement is primarily used for data analysis and visualisation. Secondary collaboration with other EDM methods may be used to develop more purposeful applications of educational data mining. Statistics and visualisation are used less frequently than ARM but more frequently than other machine learning techniques. Rest EDM methods are either in development or are used infrequently.

Before 2007, relationship mining methods were the most commonly used techniques in the early stages of EDM[65]. Prediction methods (classification and regression) and structure discovery methods (clustering) gained prominence and were widely used between 2005 and 2009[4]. Latent knowledge estimation and domain structure discovery appeared in the majority of the articles published between 2009 and 2012. Most published articles in recent years, i.e., 2013-2020, use process mining methods, domain structure discovery methods, and discovery with models methods. Another recent popular trend is the use of two or more EDM methods in tandem. A significant future development could be the creation of a free EDM preprocessing tool capable of performing all preprocessing functions in a user-friendly environment[66]. EDM applications such as user modelling, domain modelling, and educational recommender systems must be targeted for effective EDM method utilisation. Advanced machine learning methods and optimization techniques could also be used to solve educational problems. A lot of research could be done in this area in the future.

IX. CONCLUSION

Educational Data Mining is a rapidly expanding interdisciplinary research field with numerous opportunities, researchers. EDM researchers are working toward the betterment of educational processes through the introduction of new technologies and effective learning practises. The literature is discussed in this article. A survey of various EDM methods has been conducted, as well as a brief discussion of research trends and future directions research. This survey provides readers with useful information. information to proceed with the research work in this area.

REFERENCES

- [1]. J. Han, J. Pei, and M. Kamber, Data mining: concepts and techniques. Elsevier, New York, USA, 2011.
- [2]. Website of International Educational Data Mining Society (last accessed on 19 February 2020). Available at www.educationaldatamining.org.
- [3]. P. Nithya, B. Umamaheswari, and A. Umadevi, "A survey on educational data mining in field of education," International Journal of Advanced Research in Computer Engineering & Technology (IJARCET), vol. 5, no. 1, 2016.
- [4]. R. S. Baker and P. S. Inventado, "Educational data mining and learning analytics," in Learning Analytics, pp. 61–75, Springer, New York, USA, 2014.
- [5]. C. Romero and S. Ventura, "Educational data mining: A survey from 1995 to 2005," Expert Systems with Applications, vol. 33, no. 1, pp. 135–146, 2007.
- [6]. C. Romero and S. Ventura, "Educational data mining: a review of the state of the art," IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews), vol. 40, no. 6, pp. 601–618, 2010.
- [7]. R. S. Baker and K. Yacef, "The state of educational data mining in 2009: A review and future visions," JEDM-Journal of Educational Data Mining, vol. 1, no. 1, pp. 3–17, 2009.
- [8]. R. S. Baker, "Data mining for education," International encyclopedia of education, vol. 7, no. 3, pp. 112–118, 2010.
- [9]. G. Kashyap and E. Chauhan, "Review on educational data mining techniques," International Journal of Advance Technology in Engineering and Science, vol. 3, no. 11, 2015.
- [10]. C. Romero and S. Ventura, "Data mining in education," Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery, vol. 3, no. 1, pp. 12–27, 2013.
- [11]. C. Vieira, P. Parsons, and V. Byrd, "Visual learning analytics of educational data: A systematic literature review and research agenda," Computers & Education, vol. 122, pp. 119–135, 2018.
- [12]. A. Shukla, "Ph.d. thesis report - a study of relationship between symbols and sentiments for management of annotated academic resources," tech. rep., MNNIT Allahabad, India, 2014.
- [13]. A. Merceron and K. Yacef, "Educational data mining: a case study," in AIED, pp. 467–474, 2005.
- [14]. C. Romero, S. Ventura, and E. García, "Data mining in course management systems: Moodle case study and tutorial," Computers & Education, vol. 51, no. 1, pp. 368–384, 2008.
- [15]. R. Paiva, I. I. Bittencourt, W. Lemos, A. Vinicius, and D. Dermeval, "Visualizing learning analytics and educational data mining outputs," in International Conference on Artificial Intelligence in Education, pp. 251–256, Springer, 2018.
- [16]. B. Minaei-Bidgoli, D. A. Kashy, G. Kortemeyer, and W. F. Punch, "Predicting student performance: an application of data mining methods with an educational web-based system," in Frontiers in Education, 2003. FIE 2003 33rd Annual, vol. 1, pp. T2A–13, IEEE, 2003.
- [17]. B. K. Baradwaj and S. Pal, "Mining educational data to analyze students' performance," Computing Research Repository - CoRR, vol. abs/1201.3417, 2012.
- [18]. C. Romero, S. Ventura, P. G. Espejo, and C. Hervás, "Data mining algorithms to classify students," in Educational Data Mining, 2008.
- [19]. P. Kaur, M. Singh, and G. S. Josan, "Classification and prediction based data mining algorithms to predict slow learners in education sector," Procedia Computer Science, vol. 57, pp. 500–508, 2015.
- [20]. C. Romero, P. G. Espejo, A. Zafra, J. R. Romero, and S. Ventura, "Web usage mining for predicting final marks of students that use moodle courses," Computer Applications in Engineering Education, vol. 21, no. 1, pp. 135–146, 2013.
- [21]. C. Márquez-Vera, A. Cano, C. Romero, and S. Ventura, "Predicting student failure at school using genetic programming and different data mining approaches with high dimensional and imbalanced data," Applied Intelligence, vol. 38, no. 3, pp. 315–330, 2013.
- [22]. C. Márquez-Vera, A. Cano, C. Romero, AY.M. Noaman, H. Mousa Fardoun, and S. Ventura, "Early dropout prediction using data mining: a case study with high school students," Expert Systems, vol. 33, no. 1, pp. 107–124, 2016.

- [23]. E. B. Costa, B. Fonseca, M. A. Santana, F. F. de Araújo, and J. Rego, "Evaluating the effectiveness of educational data mining techniques for early prediction of students' academic failure in introductory program- ming courses," *Computers in Human Behavior*, vol. 73, pp. 247–256, 2017.
- [24]. M. I. Lopez, J. Luna, C. Romero, and S. Ventura, "Classification via clustering for predicting final marks based on student participation in forums," *International Educational Data Mining Society*, 2012
- [25]. J. L. Olmo, C. Romero, E. Gibaja, and S. Ventura, "Improving metalearning for algorithm selection by using multi-label classification: A case of study with educational data sets," *International Journal of Computational Intelligence Systems*, vol. 8, no. 6, pp. 1144–1164, 2015.
- [26]. A. Zafrá and S. Ventura, "Multi-instance genetic programming for predicting student performance in web based educational environments," *Applied Soft Computing*, vol. 12, no. 8, pp. 2693– 2706, 2012.
- [27]. D. A. Al Hammadi and M. S. Aksoy, "Data mining in education-an experimental study," *International Journal of Computer Applications*, vol. 62, Jan. 2013.
- [28]. N. Delavari and S. Phon-Amnuaisuk, "Data mining application in higher learning institutions," *Informatics in Education*, 2008.
- [29]. M. A. Sethi and M. C. Singh, "Data mining for prediction and classifica- tion of engineering students achievements using improved naïve bayes," *International Journal of Advanced Research in Computer Engineering & Technology (IJARCET)*, vol. 6, pp. 966– 971, July 2017.
- [30]. J. Zimmermann, K. H. Brodersen, H. R. Heinimann, and J. M. Buhmann, "A model-based approach to predicting graduate-leve performance using indicators of undergraduate-level performance," *Journal of Educational Data Mining*, vol. 7, no. 3, pp. 151–176, 2015.
- [31]. P. Strecht, L. Cruz, C. Soares, J. Mendes-Moreira, and R. Abreu, "A comparative study of classification and regression algorithms for modelling students' academic performance.," *International Educational Data Mining Society*, 2015.
- [32]. C. R. Beal, L. Qu, and H. Lee, "Classifying learner engagement through integration of multiple data sources," in *Proceedings of the National Conference on Artificial Intelligence*, vol. 21, pp. 151–156, Menlo Park, CA; Cambridge, MA; London; AAAI Press; MIT Press; 1999, 2006.
- [33]. J. Ocumpaugh, R. Baker, S. Gowda, N. Heffernan, and C. Heffernan, "Population validity for educational data mining models: A case study in affect detection," *British Journal of Educational Technology*, vol. 45, no. 3, pp. 487–501, 2014.
- [34]. A. T. Corbett and J. R. Anderson, "Knowledge tracing: Modeling the acquisition of procedural knowledge," *User Modeling and UserAdapted Interaction*, vol. 4, pp. 253–278, Dec 1994.
- [35]. A. Aleem and M. M. Gore, "C-BEAM: A confidence-based evaluation of MCQs for providing feedback to instructors," *Computer Applications in Engineering Education*, vol. 27, no. 1, pp. 112–127, 2019.
- [36]. W. Tissera, R. Athauda, and H. Fernando, "Discovery of strongly related subjects in the undergraduate syllabi using data mining," in *International Conference on Information and Automation, 2006. ICIA 2006.*, pp. 57– 62, IEEE, Shandong, China, 2006.
- [37]. A. Buldu and K. U" c,gu"n, "Data mining application on students' data," *Procedia-Social and Behavioral Sciences*, vol. 2, no. 2, pp. 5251–5259, 2010.
- [38]. E. Chandra and K. Nandhini, "Knowledge mining from student data," *European Journal of Scientific Research*, vol. 47, no. 1, pp. 156–163, 2010.
- [39]. M. Chalaris, S. Gritzalis, M. Maragoudakis, C. Sgouropoulou, and A. Tsolakidis, "Improving quality of educational processes providing new knowledge using data mining techniques," *Procedia-Social and Behavioral Sciences*, vol. 147, pp. 390–397, 2014.
- [40]. F. Shi, Q. Miao, and D. Mei, "The application of data association mining technology in university curriculum management," in *2012 IEEE Symposium on Robotics and Applications (ISRA)*, pp. 521– 524, IEEE, Kuala Lumpur, Malaysia, June 2012.
- [41]. P. J. DeCarlo and N. Rizk, "The design and development of an expert system prototype for enhancing exam quality," *International Journal of Advanced Corporate Learning (iJAC)*, vol. 3, pp. 10–13, 2010.

- [42]. J. M. Luna, C. Romero, J. R. Romero, and S. Ventura, "An evolutionary algorithm for the discovery of rare class association rules in learning management systems," *Applied Intelligence*, vol. 42, no. 3, pp. 501–513, 2015.
- [43]. C. Romero, A. Zafra, J. M. Luna, and S. Ventura, "Association rule mining using genetic programming to provide feedback to instructors from multiple-choice quiz data," *Expert Systems*, vol. 30, no. 2, pp. 162–172, 2013.
- [44]. G. Dimić, B. Predić, D. Rančić, V. Petrović, N. Maček, and P. Spalević, "Association analysis of moodle e-tests in blended learning educational environment," *Computer Applications in Engineering Education*, vol. 26, no. 3, pp. 417–430, 2018.
- [45]. A. Aleem, A. Kumar, and M. M. Gore, "A study of manuscripts evolution to perfection," in *2nd International Conference on Advanced Computing and Software Engineering - ICACSE 2019*, pp. 278–282, 2019.
- [46]. S.Y. Deng and X. Que, "Research on the teaching assessment of students of science and engineering teachers in a university," *Computer Applications in Engineering Education*, vol. 27, no. 1, pp. 5–12, 2019.
- [47]. A. Aleem and M. M. Gore, "The choice is yours: The effects of optional questions in engineering examinations," *Computer Applications in Engineering Education*, vol. 27, no. 5, pp. 1087–1102, 2019.
- [48]. P. Fournier-Viger, J. C.-W. Lin, R. U. Kiran, Y. S. Koh, and R. Thomas, "A survey of sequential pattern mining," *Data Science and Pattern Recognition*, vol. 1, no. 1, pp. 54–77, 2017.
- [49]. J. Li, T. D. Le, L. Liu, J. Liu, Z. Jin, and B. Sun, "Mining causal association rules," in *13th IEEE International Conference on DataMining Workshops (ICDMW)*, pp. 114–123, IEEE, Dallas, Texas, USA, 2013.
- [50]. S. Mani and G. F. Cooper, "A simulation study of three related causal data mining algorithms," in *Eighth International Workshop on Artificial Intelligence and Statistics (AISTATS)*, pp. 73–80, Society for Artificial Intelligence and Statistics, Key West, Florida, USA, Jan. 2001.
- [51]. A. H. Cairns, B. Gueni, M. Fhima, A. Cairns, and S. N. K. David, "Process mining in the education domain," *International Journal on Advances in Intelligent Systems*, vol. 8, no. 1 & 2, pp. 219–232, 2015.
- [52]. A. Bogarín, R. Cerezo, and C. Romero, "A survey on educational process mining," *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, vol. 8, no. 1, p. e1230, 2018.
- [53]. T. S. Uhlmann, E. A. P. Santos, and L. A. Mendes, "Process mining applied to player interaction and decision taking analysis in educational remote games," in *International Conference on Remote Engineering and Virtual Instrumentation*, pp. 425–434, Springer, 2018.
- [54]. A. Aleem, R. Srivastava, A. K. Singh, and M. M. Gore, "GCLOD: A clustering algorithm for improved intra-cluster similarity and efficient local outliers detection.," in *DMIN*, pp. 524–530, 2009.
- [55]. R. Ahuja, A. Jha, R. Maurya, and R. Srivastava, "Analysis of educational data mining," in *Harmony Search and Nature Inspired Optimization Algorithms*, pp. 897–907, Springer, 2019.
- [56]. A. Dutt, M. A. Ismail, and T. Herawan, "A systematic review on educational data mining," *IEEE Access*, 2017.
- [57]. S. H. Ganesh and A. J. Christy, "Applications of educational data mining: A survey," in *Innovations in Information, Embedded and Communication Systems (ICIIECS)*, 2015 International Conference on, pp. 1–6, IEEE, 2015.
- [58]. L. Ramanathan, G. Parthasarathy, K. Vijayakumar, L. Lakshmanan, and S. Ramani, "Cluster-based distributed architecture for prediction of student's performance in higher education," *Cluster Computing*, pp. 1–16, 2018.
- [59]. M. Goyal and R. Vohra, "Applications of data mining in higher education," *International Journal of Computer Science*, vol. 9, no. 2, pp. 113–120, 2012.
- [60]. M. C. Desmarais, "Mapping question items to skills with nonnegative matrix factorization," *ACM SIGKDD Explorations Newsletter*, vol. 13, no. 2, pp. 30–36, 2012.
- [61]. M. I. Al-Twijri and A. Y. Noaman, "A new data mining model adopted for higher institutions," *Procedia Computer Science*, vol. 65, pp. 836–844, 2015.

- [62]. H. Jeong and G. Biswas, "Mining student behavior models in learning- by-teaching environments," in Educational Data Mining, 2008.
- [63]. G. Kostopoulos, S. Kotsiantis, and P. Pintelas, "Predicting student performance in distance higher education using semi-supervised techniques," in Model and Data Engineering, pp. 259–270, Springer, 2015.
- [64]. G. Kostopoulos, A.-D. Lipitakis, S. Kotsiantis, and G. Gravvanis, "Pre- dicting student performance in distance higher education using active learning," in International Conference on Engineering Applications of Neural Networks, pp. 75–86, Springer, 2017.
- [65]. R. Jindal and M. D. Borah, "A survey on educational data mining and research trends," International Journal of Database Management Systems (IJDMS), vol. 5, no. 3, p. 53, 2013.
- [66]. C. Romero, J. R. Romero, and S. Ventura, "A survey on preprocessing education