



Student Academic Performance Monitoring and Evaluation using Data Mining Techniques

Harpreet Kaur

Assistant Professor

Guru Kashi University, Talwandi Sabo, Punjab, India

Abstract: *Predicting students' performance is an essential activity towards the success of the world's education sector. However, the action continues to present itself as a challenging task due to the existing large data amounts in educational databases. On the other hand, some Institutions lack systems that are capable of analyzing and monitoring students' performance. This problem could be partially due to a lack of awareness about the importance of predicting students' performance. In addition to that, the existing studies on performance prediction methods are still inadequate in identifying and convincing educators with the most suitable method for predicting students' performance. This review explores the commonly used data mining techniques to predict students' performance in previous studies to find out the most suitable technology that can be trusted with predicting students' performance. The result of the study showed that the decision trees algorithm is the best classification technique that gives trusted and accurate results when it comes to student performance prediction. Predicting students' performance helps in monitoring the students' progress, both pass and fail, and therefore provides a gap for early interventions and decision making by educators. This opportunity dramatically helps in promoting the education sector by raising the academic standards of educational Institutions.*

Keywords: EDM, Data, mining, feature, accuracy etc

I. INTRODUCTION

Throughout the study life of students, performance is one of the most critical aspects of their success. This condition is because usually, performance is an essential part of learning from junior school to higher institutions of learning. Providing quality education to students requires continuous monitoring of learning and teaching activities in an education environment. However, this tends to be difficult because of large amounts of data in educational databases. Nevertheless, the new technology of big data and machine learning comes in to solve all these problems. Big data is a set of methods and techniques that necessitate a new form of combinations to find large hidden values from various complex datasets on a considerable scale. This technology is very advantageous in a way that it can determine the patterns obtained from data analysis to understand hidden information and ease decision making [1], [2] Currently, many methods have been proposed by researchers to evaluate the performance of students in the education field. However, one of the most recently popular proposed and applied procedures to predict the performance of students is data mining. This method has several advantages; for example, data mining techniques can be used in e-learning since it is easy to access students' log data, and automatic data analysis can be done [3]. In big data and machine learning technology, the application of data mining in the education domain is known as 'education data mining' (EDM). This aims at converting education data into very useful information which brings a positive impact to the education sector and promotes research practices such as predicting student performance [4]. Performance prediction can help strengthen the education sector and assist responsible leaders in providing practical teaching approaches, thus encouraging higher achievements since student performance is monitored right from the early stages of education. Therefore, students are also equipped with the understanding and improving their learning activities to avoid risks such as school dropouts, poor performance, and failure to graduate. To ease the process of students' performance prediction, we propose a systematic review to help us achieve this study's objectives that are;



1. To identify the commonly used methods to predict student performance.
2. To find out the most suitable method for student performance prediction among the commonly used methods.

The improvement in the quality of education is one of the most significant aspects of forming a successful member of society. The data stored in educational institutions repository plays an important role in order to extract hidden and interesting patterns to assist every stakeholder of an educational process [1]. There are many techniques being anticipated to assess the student academic performance in way of making fruit full future of a student. Predicting performance of student has been continued to a hot topic in the Educational data mining domain. Data mining is considered to be one of the best choices for the researchers to analyse student's performance. The techniques of data mining are extensively used on educational data now- a day's [2, 3]. It is called educational data mining. Educational Data Mining (EDM) explores the educational data to better understand the issues of student's performance using the fundamental nature of data mining techniques [4]. EDM manipulates educational data to help educational institutions to plan educational strategies, in order to improve the educational quality. Prediction is one of the main areas in EDM. Prediction and analysis of student academic performance are essential for student academic growth. Identifying the factors affecting the student academic performance is complicated research task [5]. The original of academic data contains many irrelevant and redundant data .This redundant data effects the results of prediction. Feature selection methods minimize the redundancy and maximize relevancy of features without any loss of crucial data [6]. Feature Selection is very dynamic and productive field and research area of machine learning and data mining. The main goal of feature selection is to choose a subset by eliminating non-predictive data. Furthermore, it increases the predictive accuracy and reduces the complexity of learned results .The effectiveness of student performance prediction models can be increased in connection with feature selection techniques. Feature Selection techniques can be classified in to three groups: filter, wrapper, and embedded models. Filter method depends upon general characteristics of training data, this method is done on pre-processing stage and not dependent on a learning algorithm. Wrapper method uses learning algorithms to evaluate the features. Embedded methods are specific to some given learning algorithms, and these methods are performed on training process of classifiers. Previously, alot of work is done to predict the performance of student using different feature selection techniques [5]. In recent studies, researchers use different feature selection techniques and the combination of classifiers to produce efficient prediction models. A research is required to identify the performance analysis in terms of prediction accuracy in combination of different feature selection algorithms with differently classifiers. This paper is a step towards identifying the prediction accuracy of different available feature selection algorithm in the context of classifiers being used on educational data.

II. FEATURE SELECTION

Feature selection has been an active and fruitful field of research area in pattern recognition, machine learning, statistics and data mining communities [2, 3]. The main objective of feature selection is to choose a subset of input variables by eliminating features, which are irrelevant or of no predictive information. Feature selection has proven in both theory and practice to be effective in enhancing learning efficiency, increasing predictive accuracy and reducing complexity of learned results [4, 5]. Feature selection in supervised learning has a main goal of finding a feature subset that produces higher classification accuracy. As the dimensionality of a domain expands, the number of features N increases. Finding an optimal feature subset is intractable and problems related feature selections have been proved to be NP-hard [6]. At this juncture, it is essential to describe traditional feature selection process, which consists of four basic steps, namely, subset generation, subset evaluation, stopping criterion, and validation [7]. Subset generation is a search process that produces candidate feature subsets for evaluation based on a certain search strategy. Each candidate subset is evaluated and compared with the previous best one according to a certain evaluation. If the new subset turns to be better, it replaces best one. This process is repeated until a given stopping condition is satisfied. Ranking of features determines the importance of any individual feature, neglecting their possible interactions. Ranking methods are based on statistics, information theory, or on some functions of classifier's outputs [8]. Algorithms for feature selection fall into two broad categories namely wrappers that use the learning algorithm itself to evaluate the usefulness of features and filters that evaluate features according to heuristics based on general characteristics of the data [7].



Several justifications for the use of filters for subset selection have been discussed [6] and it has been reported that filters are comparatively faster than wrappers. Many student performance prediction models have been proposed and comparative analyses of different classifier models using Decision Tree, Bayesian Network, and other classification algorithms have also been discussed. But, they reveal only classifier accuracy without performing the feature selection procedures.

III. EXISTING ISSUE

The aim of HSC is to provide quality education to their students. The quality of education in HSC can be increased by discovering new knowledge or patterns for prediction of students' performances especially in the academic aspect. The prediction on SAP can be used as a guideline for the faculty management or educators to prevent students from dropout. The objective of this study is to get the patterns of SAP focusing on the first semester of the first year Bachelor of Computer Science with specialization in Software Development at the FIT. At the beginning of the semester for new students, educators face difficulty to know and analyse the student's performance because there are lack of information about their students' previous background. All the information about students is stored in a database at Academic Department, Student Entry Management Department (SEMD), Ministry of Higher Education based in different location (Kuala Lumpur, Malaysia). The selected parameters from Academic Department, database are gender, hometown, race, and GPA. The parameters on students' university entry mode and family income are taken from different database which is located in SEMD database. The study is made to determine whether or not the selected parameters contribute to the SAP. Besides, this study is also conducted to find out the relationship between the independent parameters and the dependent parameter (GPA). The discovered pattern can be used by educators the first year bachelor students at FIT. The development of effective systems to predict SAP is very important to provide more information about the students to the educators. Therefore, the educators would know how the first year students will perform in their academic from the earliest moment. A prototype system will be developed using the discovered patterns that are extracted from the DM process. In addition, the system can work as a helping tool for educators to plan the teaching materials in order to improve students' performance, and to decrease the failure rate in computer science course. Artificial Neural Network (ANN) is based on the human brain architecture that consists of multiple processing layers connected with nodes. ANN method is used in the educational field for predicting and classifying SAP.

IV. LITERATURE SURVEY

4.1 General Applications of Data Mining

Analyzing large amounts of data is a necessity. Presently, data mining continues to achieve significant success in many areas of human activities such as medicine, business, and robotics. It is, therefore, beneficial to all Institutions and enterprises of all kinds to collect and analyze their data. With proper analysis, administrators in health sectors can learn trends and identify issues that come up in the patients' records, find ways of managing the hospital's human resources better and enhance the performance of hospital staff [2]. In addition to that, nowadays, the technology of big data gives a chance to professional health workers to screen and predict mental illness, derive essential health trends, and support timely preventive care. With the significant increase of data in the agricultural domain, data mining is the only capable technology employed to transform agriculture into smart agriculture [3]. Data mining technology can also be employed to monitor energy consumption of household appliances thorough monitoring of electricity consumption patterns (ECPs) and their relationship with household features, thereby economizing the electricity consumption in residential areas [4]. The technology of machine learning and data mining can be applied in cyber-manufacturing systems to detect and effectively manage cyber security and computer network intrusions by using big data tools. In addition to the above, the technique of data mining can be useful for financial firms to correct their material in case of an inaccuracy. Applying algorithms such as SVM, naive bayes, and ANN on financial data gives an easy way of revising past published financial statements [5]. Big data applications can be also used to spot trends in blogs and understand the



similarities and differences in social media messages to make an easy way of operation for online businesses to operate [6].

Maryam Zaffar et.al.[2017] have studied Student's academic performance is the main focus of all educational institutions. Educational Data Mining (EDM) is an emerging research area help the educational institutions to improve the performance of their students. Feature Selection (FS) algorithms remove irrelevant data from the educational dataset and hence increases the performance of classifiers used in EDM techniques. This paper present an analysis of the performance of feature selection algorithms on student data set. The obtained results of the different FS algorithms and classifiers will also help the new researchers in finding the best combinations of FS algorithms and classifiers. Selecting relevant features for student prediction model is very sensitive issue for educational stakeholders, as they have to take decisions on the basis of results of prediction models. Furthermore our paper is an attempt of playing a positive role in the improvement of education quality, as well as guides new researchers in making academic intervention.[1]

R. Sasi Regha et.al.[2016] have studied Technology has revolutionized the field of education. As a result, the education related data is been increasing rapidly. This made data mining approaches to spot over educational data ended in Educational data mining (EDM). The regulation focuses on investigating educational data to build models for enhancing learning experiences and improving institutional effectiveness. In this paper, the data mining techniques is used for predicting the student performance in different educational levels. Irrelevant features, along with redundant features, rigorously influence the accuracy of the classification of student performance. Therefore, feature selection should be able to detect and eliminate both irrelevant and redundant features as hard as possible. After feature selecting process, two effective classification techniques i.e., Prism and J48 is used for predicting the student performance. Experimentation result is shown that the feature selection method is well effective.[2]

Dr.M.Chidambaram et.al.[2016] have studied Feature Selection is a fundamental problem in machine learning and data mining . Feature Selection is an effective way for reducing dimensionality, removing irrelevant data increasing learning accuracy. Feature Selection is the process of identifying a subset of the most useful features that produce compatible results as the original entire set of features .A Feature Selection techniques may be evaluated from both efficiency and effectiveness point of view. While the efficiency concerns the time required to find a subset of features, the effectiveness is related to the quality of the subset of the subset of features. Feature Selection is different from dimensionality reduction. Both methods search for to reduce the number of attributes in the dataset. But dimensionality reduction method creating new combination of attributes. Feature Selection methods include and exclude attributes present in the data without change. The central assumption when using a Feature Selection technique is that the data contains many redundant or irrelevant features. This paper actually a survey on various technique of feature selection and its advantages disadvantages.[3]

Mital Doshi et.al.[2014] have studied Feature selection is a preprocessing step to machine learning which is effective in reducing dimensionality, helps in removing irrelevant data, increasing learning accuracy, and improving result. In this paper we have shown different feature selection approaches, applications and the relation between them and the various machine learning algorithms. [4]

Anal Acharya et.al.[2014] have studied web based learning has emerged as a new field of research due to growth of network and communication technology. These learning systems generate a large volume of student data. Data mining algorithms may be applied on this data set to study interesting patterns. As an example, student enrollment data and his past examination records could be used to predict his grades in the term end examination. However this prediction could mean examining a lot of features of the student data resulting in creation of a model with high computational complexity. In this context this work first defines a student data set with 309 records and 14 features collected by a survey from various graduation level students majoring in Computer Science under University of Calcutta. Different feature selection algorithms are applied on this data set. The best results are obtained by Correlation Based Feature Selection algorithm with 8 features. Subsequently classification algorithms may be applied on this feature subset for predicting student grades.[5]



M. Ramaswami et.al.[2009] have studied Educational data mining (EDM) is a new growing research area and the essence of data mining concepts are used in the educational field for the purpose of extracting useful information on the behaviors of students in the learning process. In this EDM, feature selection is to be made for the generation of subset of candidate variables. As the feature selection influences the predictive accuracy of any performance model, it is essential to study elaborately the effectiveness of student performance model in connection with feature selection techniques. In this connection, the present study is devoted not only to investigate the most relevant subset features with minimum cardinality for achieving high predictive performance by adopting various filtered feature selection techniques in data mining but also to evaluate the goodness of subsets with different cardinalities and the quality of six filtered feature selection algorithms in terms of F-measure value and Receiver Operating Characteristics (ROC) value, generated by the NaïveBayes algorithm as base-line classifier method. The comparative study carried out by us on six filter feature section algorithms reveals the best method, as well as optimal dimensionality of the feature subset. Benchmarking of filter feature selection method is subsequently carried out by deploying different classifier models. The result of the present study effectively supports the well known fact of increase in the predictive accuracy with the existence of minimum number of features. The expected outcomes show a reduction in computational time and constructional cost in both training and classification phases of the student performance model.[6]

V. STUDENT PERFORMANCE PREDICTION

Student performance prediction consists of identifying an unknown mark of a student [6]. However, some factors can affect students' performance, making this task challenging to accomplish. The elements may include the economic status of a student, demographic characteristics, students' psychological profiles, past scholar experiences, different cultural backgrounds, and interactions between fellow students [7]. Despite the above, for any institution to carry out a right performance prediction, majorly, they must first identify the risk factors that may affect prediction results. All educational institutions should always answer this critical question before taking prediction steps; "What are the important risk factors or variables for predicting student performance?" These risk factors may include employed student performance prediction methods, and the data sets considered to obtain prediction results. The data sets may include attributes of the previous semester grade, gender, attendance, GPA, parents' education, parents' income, scholarship, first child, etc. Different Institutions consider different data sets to predict student performance. However, diverse data sets work well with varying methods of prediction. Therefore, the predicted results may vary depending on the techniques used. Although these factors that affect prediction results are apparent, some variables may be delicate and difficult to understand and identify without applying a more sophisticated analysis. Therefore, using modern data mining techniques such as decision trees, naive bayes, and neural networks may accurately predict student performance (pass/fail) compared to other current models. The results from accurate predictions can help Institutions to attain quality education. Below, the conventional methods used for predicting student performance.

VI. CONVENTIONAL METHODS USED FOR PREDICTING STUDENT PERFORMANCE

The process of predictive modeling or predictive analytics are employed in educational data mining to predict outcomes, i.e., it is used in predicting the performance of students. In this domain, predictive analytics acts as an effective process that decision-makers and institution heads may use to optimize the limited resources and plan their Institution strategy and policy effectively [38], [30]. However, building predictive modeling requires the use of different tasks, such as categorization, regression, and classification. Among all these tasks involved in predictive modeling, classification is the most commonly used task for student performance prediction. Classification uses several algorithms to perform this task of student performance prediction. The algorithms include naive bayes, decision trees, support vector machine, and artificial neural networks. The most frequently used data mining methods for student performance grouped by their algorithms are explained in detail below.



6.1 Neural Networks

Neural network is a biologically inspired analytical method capable of modeling extremely complex nonlinear functions. Neural networks are part of the popularly employed algorithms in the education domain for student performance prediction. ANN is mostly desired because it can classify patterns without requiring any training. Being inherently parallel and thus able to speed up the computational process makes ANN suitable for prediction activities in the educational data mining domain.

6.2 Decision Trees

Decision trees are a classification approach that works by constructing decision trees following a top-down recursive divide-and-conquer way. These algorithms are increasingly becoming a well-known prediction approach in data science due to their working characteristics and usefulness in discovering useful models [3]. Some researchers have opted for this method due to its ease and ability to use all size data structures to estimate values. Decision trees are easy to interpret because their tree-like structure brings the classification rules to real-life human reasoning. Some researchers used a decision tree algorithm to show the impact of data mining technology in education by predicting the dropout students, segmenting students according to their performances, student retention management, and predicting drop out of students. In a particular prediction study, bagged trees, adaptive boosting trees, and random forests attained an accuracy of 88.7%, 95.7%, and 96.1%, respectively [41].

6.3 Naive Bayes

The naive bayes method is a supervised classifier that bases on applying the bayes' theorem with strong naive independence assumptions between the explanatory variables and uses two simplifications where one uses the conditional independence assumption and the other ignores the denominator [4]. This algorithm is also one of the commonly preferred methods by researchers to carry out prediction activities since it learns fast, predicts equally, and does not need ample storage.

6.4 K- Nearest Neighbor

K- nearest-neighbor classifiers are an analogy-based algorithm that learns by comparing themselves with similarly given test training tuples. The algorithms can be employed in numeric predictions to return a real-valued forecast for an unidentified tuple. The algorithm, therefore, restores the reasonable value associated with the k-nearest neighbors of the unidentified tuple. When used for student performance prediction, KNN gave good results [10].

6.5 Support Vector Machine (SVM)

This algorithm is a classification technique for linear and nonlinear data. In the SVM supervised learning method, the original training data is transformed into a higher dimension using nonlinear mapping. SVM algorithm takes long training time but has a high accuracy and ability to perform excellently with small data sets as well as predict at-risk and marginal students [26]. Table 1 summarizes the prediction accuracy of artificial neural networks, decision trees, naive bayes, k- nearest neighbor, and support vector machine.

VII. PROBLEM FORMULATION

In the research work Feature Selection Algorithm for Educational Data Mining different problems are faced that are given below:

- As the dimensionality of a domain expands, the number of features N increases. Finding an optimal feature subset is intractable and problems related feature selections have been proved to be NP-hard.
- It is to be noted that Receiver Operating Characteristic (ROC) parameter is not used here since it applies to a 2-class problem.



- Many kinds of redundancy, irrelevant and noisy feature as well as Relief algorithm limitations in high datasets.
- Another problem is high rate feature selection and accuracy problem.

VIII.METHODOLOGY

The main aim of the research is to evaluate the performance of different FS algorithms on different classification algorithms using student dataset. The comparison between different FS algorithms give a deep insight to new educational data miners about the performance of different feature selection algorithms on student data .To achieve the objective of the research , a student dataset is taken from a valid sources, and then different FS algorithms are applied on it , which was not used earlier on this dataset. Different classification algorithms are applied by using selected FS algorithms, and furthermore evaluated to check the best performance among all the combinations applied on student data set.

8.1 Data set Description

The dataset used in this study is taken from the source www.kaggle.com, and is comprised of 500 students 16 features. This dataset has been used in the study [11], to check the learner's interactivity with e- learning management system, bagging and boosting methods are applied on the given dataset, however, only information gain based feature selection algorithm is used previously. In this paper, the main aim of using the dataset is to identify the best combinations of FS algorithms and classifiers, in order to identify the key performance factors on the academic achievements of students. WEKA (Waikato Environment for Knowledge Analysis) is used as a data mining tool. It has a rich source of Machine learning algorithms. WEKA is developed by the University of Waikato in New Zealand. It is an open source software developed in JAVA language, that provides facility for developing machine learning techniques for data mining tasks.

8.2 Feature Selection Algorithm and Classifiers

In this research work six FS algorithm CfsSubsetEval, ChiSquaredAttributeEval, Filtered Attribute Eval, GainRatioAttributeEval, Principal Components, and ReliefAttributeEval are evaluated. The classification algorithm BayesNet(BN), Naïve Bayes(NB), NaiveBayesUpdateable(NBU), MLP, Simple Logistic(SL), SMO, Decision Table(DT), Jrip, OneR, OneR, DecsionStump(DS), J48, Random Forest(RF), RandomTree(RT), REPTree(RepT) are evaluated through the educational data set.

IX. CONCLUSION

Student Performance is a crucial factor that requires proper monitoring if the goal of education in higher educational institutions and at all levels of education is to be achieved. This condition is because student performance prediction helps institutional leaders in improving their educational systems. This study aimed at reviewing the commonly used classification techniques for predicting student performance. Among the widely used methods to predict student performance, the decision trees method proved to be the best method for predicting student performance compared to Naive bayes, k-nearest neighbor, support vector machines and artificial networks, due to its simplicity of use and ability to uncover small or large data structures and predict values which gives high accuracy. In conclusion, the results from this review will help educators to monitor students' performance systematically by using the easiest and most accurate method to predict student performance. The authors believe that using the best prediction method helps educators to infer students' performance, which allows early interventions that may bring an increase in excellent academic performance rates, thus promoting education with high-quality achievers. In the future work, authors can use the data in this review as a basis for other related studies in the educational data mining field. Focusing on how to improve the accuracy of other methods would also be Important.

**REFERENCES**

- [1]. Maryam Zaffar et.al. "Performance Analysis of Feature Selection Algorithm for Educational Data Mining" IEEE Conference on Big Data and Analytics (ICBDA)-2017.
- [2]. R. Sasi Regha et.al. "Optimization Feature Selection for classifying student in Educational Data Mining" International Journal of Innovations in Engineering and Technology (IJET) , Volume 7 Issue 4 December 2016.
- [3]. Dr.M.Chidambaram et.al. "A Survey on Feature Selection in Data Mining " International Journal of Innovative Research in Computer Science & Technology (IJRCST) ISSN: 2347-5552, Volume-4, Issue-1, January 2016.
- [4]. Mital Doshi et.al. "Survey of Feature Selection Algorithms in Higher Education" International Journal of Computer Applications in Engineering Sciences [VOL IV, ISSUE I, MARCH 2014].
- [5]. Anal Acharya et.al. "Application of Feature Selection Methods in Educational Data Mining" International Journal of Computer Applications © 2014 by IJCA Journal Volume 103 - Number 2 Year of Publication: 2014.
- [6]. M. Ramaswami et.al. "A Study on Feature Selection Techniques in Educational Data Mining" Journal Of Computing, Volume 1, Issue 1, December 2009.
- [7]. E. Osmanbegović, M. Suljić, and H. Agić, "DETERMINING DOMINANT FACTOR FOR STUDENTS PERFORMANCE PREDICTION BY USING DATA MINING CLASSIFICATION ALGORITHMS," Tranzicija, vol. 16, pp. 147-158, 2015.
- [8]. M. Shahiri and W. Husain, "A review on predicting student's performance using data mining techniques," Procedia Computer Science, vol. 72, pp. 414-422, 2015.
- [9]. 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, pp. 601-618, 2010.
- [10]. M. Ramaswami and R. Bhaskaran, "A study on feature selection techniques in educational data mining," arXiv preprint arXiv:0912.3924, 2009.
- [11]. A. Mueen, B. Zafar, and U. Manzoor, "Modeling and Predicting Students' Academic Performance Using Data Mining Techniques," International Journal of Modern Education and Computer Science, vol. 8, p. 36, 2016.
- [12]. W. Punlumjeak and N. Rachburee, "A comparative study of feature selection techniques for classify student performance," in Information Technology and Electrical Engineering (ICITEE), 2015 7th International Conference on, 2015, pp. 425- 429.