



Predicting and Analysis of Student Performance Using Decision Tree Technique

Ankita A Nichat, Dr.Anjali B Raut

M.E Student, Department of Computer Science & Engineering, HVPM's College of Engineering & Technology, Amravati, India

HOD, Department of Computer Science & Engineering, HVPM's College of Engineering & Technology, Amravati, India

ABSTRACT: Analyzes data mining methods and techniques students' data to construct a predictive model for students' academic performance. Technically, data mining is the process of finding correlations or patterns among dozens of fields in large relational databases. Data mining is also use for sorting the educational problem by using analysis techniques for measuring the student performance, instructor performance. In this paper, measuring student performance using classification technique such as decision tree. The task can be processed based on the several attributes to predict the performance of the student activity respectively. In this research, the paper have been focused the improvement of Prediction/ classification techniques which are used to analyze the skill expertise based on their academic performance by the scope of knowledge. Giving the details about the results, and the specific needs of studies to improvement, such as the accompaniment of students along their learning process, and the taking of timely decisions in order to prevent academic risk and desertion. Lastly, some recommendations and thoughts are laid out for the future development of performance. Helps to analyze the slow learner that are likely study in poor which are used to improve their skill as early to achieve the goal.

KEYWORDS: Data mining, Classification algorithms, decision trees.

I. INTRODUCTION

Data mining is the analysis step of the "knowledge discovery in databases" process. Data mining is the process of analysing data from different perspectives and summarizing it into useful information. Data mining software is one of a number of analytical tools for analysing data. It allows users to analyse data from many different dimensions or angles, categorize it, and summarize the relationships identified. In recent years, there has been increasing interest in the use of data mining to investigate scientific questions within educational research, an area of inquiry termed educational data mining [8]. An ability of student performance is essential in education environment, which is influenced by many qualitative attributes like Student Identity, gender, age, Specialty, Lower class Grade, higher Class Grade, Extra knowledge or skill, Resource, Attendance, Time spend to study, Class Test Grade (Internal), Seminar Performance, Lab Work, Quiz, E-Exercise, E-Homework, Over all Semester exam Percentage are included for forming the data set. Educational data mining applied many techniques are K- nearestneighbor, decision tree, Naïve Bayes, Neural network, Fuzzy, Genetic and other techniques are applied in the environment [13].

One of the important facts in institution is the rapid growth of educational data. The main goal of any educational institution is to improve education quality. Prediction of student's performance in education institution is one way to reach the good quality in education system. Educational institution staff should identify students who are likely to fail in exams. But, it is hard to identify them early because of large number of enrolled students. Student's academic performance can be affected by many factors such as personal, social, demographic data etc. Finding useful information from large database is a difficult task [13]. Data mining in education is called educational data mining, EDM. Educational data mining (EDM) is describes a research field with the application of data mining, machine



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirccce.com

Vol. 5, Issue 4, April 2017

learning and statistics to information generated from educational settings e.g., universities and intelligent tutoring systems. In addition, these methods lack the ability to reveal useful hidden information. Online Exam is being launched because a need for a destination.

Classification techniques are using for education data modelling. There is an increase in their application within the last six years [14]. Researchers prefer to apply a single technique in their studies on student evaluations like those mentioned above. In addition, these methods lack the ability to reveal useful hidden information. Exam model is being constructed which is beneficial for both college and students. With this model, institutes can register and host exams. Students can give exams and view their results with guidance and instructor can evaluate the student performance. This model is an attempt to remove the existing flaws in the manual system of conducting exams. Examination System fulfils the requirements of the institutes to conduct the exams. Thus the purpose of the model is to provide a system that saves the efforts and time of both the institutes and the students. Institutes enter the questions they want in the exam. These questions are displayed as a test to the eligible students. The answers enter by the students are then evaluated and their score is calculated and saved. This score then can be accessed by the institutes to evaluate their performance. In this study measures the student performance by using data mining technique like classification, decision tree algorithm using to build the classifier model on base on dataset composed of responses of students to courses evaluation questions.

II. RELATED WORK

Mustafa Agaoglu [1] research in educational mining focuses on modeling student's performance instead of instructors' performance. One of the common tools to evaluate instructors' performance is the course evaluation questionnaire to evaluate based on students' perception. In this study, four different classification techniques, –decision tree algorithms, support vector machines, artificial neural networks, and discriminant analysis– are used to build classifier models. Their performances are compared over a dataset composed of responses of students to a real course evaluation questionnaire using accuracy, precision, recall, and specificity performance metrics. Although all the classifier models show comparably high classification performances, C5.0 classifier is the best with respect to accuracy, precision, and specificity. In addition, an analysis of the variable importance for each classifier model is done. Accordingly, it is shown that many of the questions in the course evaluation questionnaire appear to be irrelevant. Furthermore, the analysis shows that the instructors' success based on the students' perception.

Tripti Mishra, Dr. Dharminder Kumar, Dr. Sangeeta Gupta [2] use different classification techniques to build performance prediction model based on students' social integration, academic integration, and various emotional skills which have not been considered so far. Two algorithms J48 (Implementation of C4.5) and Random Tree have been applied to the records of MCA students of colleges affiliated to Guru Gobind Singh Indraprastha University to predict third semester performance. Random Tree is found to be more accurate in predicting performance than J48 algorithm.

Keno C. Piad, Menchita Dumlao, Melvin A. Ballera, Shaneth C. Ambat [3] predicts the employability of IT graduates using nine variables. First, different classification algorithms in data mining were tested making logistic regression with accuracy of 78.4 is implemented. Based on logistic regression analysis, three academic variables directly affect; IT_Core, IT_Professional and Gender identified as significant predictors for employability. The data were collected based on the five year profiles of 515 students randomly selected at the placement office tracer study.

Bipin Bihari Jayasingh [4] initiates a sample study that is taken for a particular institution, in the particular environment, for the particular batch and particular set of students. The sample data are collected from a classroom by distributing the questionnaire attempted by two different batches of student having questions pertaining to Inquiry based and deductive learning. The system is developed and tested twice after teaching the content using inductive method and implemented using attribute relevance, discriminant rules of class discrimination mining. The results are visualized through bar charts and shows that the two batches of learners of different years have different learning characteristics.

S. M. Merchán [5] presents and analyzes the experience of applying certain data mining methods and techniques on 932 Systems Engineering students' data, from El Bosque University in Bogotá, Colombia; effort which has been



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirccce.com

Vol. 5, Issue 4, April 2017

pursued in order to construct a predictive model for students' academic performance. As an iterative discovery and learning process, the experience is analyzed according to the results obtained in each of the process' iterations. Each obtained result is evaluated regarding the results that are expected, the data's input and output characterization, what theory dictates and the pertinence of the model obtained in terms of prediction accuracy. Said pertinence is evaluated taking into account particular details about the population studied, and the specific needs manifested by the institution.

Konstantina Chrysafiadi and Maria Virvou [6] considered a novel approach of web-based education that performs individualized instruction on the domain of programming languages is presented. This approach is fully implemented and evaluated in an educational application module, which is called Fuzzy Knowledge State Definer (FuzKSD). In particular, FuzKSD performs user modeling by dynamically identifying and updating the student's knowledge level for all the concepts of the domain knowledge. The operation of FuzKSD is based on Fuzzy Cognitive Maps (FCMs) that are used to represent the dependencies among the domain concepts. FuzKSD uses fuzzy sets to represent the students' knowledge level as a subset of the domain knowledge.

M. Mayilvaganan, D. Kalpanadevi [7] research, the paper have been focused the improvement of Prediction/classification techniques which are used to analyze the skill expertise based on their academic performance by the scope of knowledge. Also the paper shows the comparative performance of C4.5 algorithm, AODE, Naïve Bayesian classifier algorithm, Multi Label K-Nearest Neighbor algorithm to find the well suited accuracy of classification algorithm and decision tree algorithm to analysis the performance of the students which can be experimented in Weka tool.

Cristóbal Romero [8] Educational data mining (EDM) is an emerging interdisciplinary research area that deals with the development of methods to explore data originating in an educational context. EDM uses computational approaches to analyze educational data in order to study educational questions. This paper surveys the most relevant studies carried out in this field to date. First, it introduces EDM and describes the different groups of user, types of educational environments, and the data they provide. It then goes on to list the most typical/common tasks in the educational environment that have been resolved through data-mining techniques, and finally, some of the most promising future lines of research are discussed.

R. V. Mane, Priyanka Patil [9] proposes Generalized Sequential Pattern mining algorithm for finding frequent patterns from student's database and Frequent Pattern tree algorithm to build the tree based on frequent patterns. This tree can be used for predicting the student's performance as pass or fail. Once the student is found at the risk of failure he/she can be provided guidance for performance improvement.

III. PROPOSED SYSTEM

We are going to propose the system by using which the user can give a test on specific educational or subject categories. When student complete the test, system will calculate the performance of the user by using the algorithm decision tree. The system will suggest to the teacher that on which topics the user is weak or need to study again. To solve the problems faced with manual examination writing, there is need for a computerized system to handle all the works. We propose an application that will provide a working environment that will be flexible and will provide ease of work and will reduce the time for report generation and other paper works. Today many organizations are conducting online examinations worldwide successfully and issue results online but they are not measuring the performance of the student and teacher not know about the weak points of the students and we are focusing on this issue. The main advantage is that it can be evaluation of answers can be fully automated for all questions and other essay type questions can be evaluated manually or through automated system, depending on the nature of the question's and the requirements. To bring, efficiency, transparency and reliability, universities should also adopt this new technology for managing the examination system.

Student modeling can be defined as the process of gathering relevant information in order to infer the current cognitive state of the student, and to represent it so as to be accessible and useful to evaluating their performance. The proposed student model consists of two categories of user. The first user represents the students who can giving the test, managing profile, student can also view their performance, scores, recommendation given by expert, etc. The second user plays a role of teacher/instructor who can view all data about students, their result, and performance of each student. Model help to analysis the student performance, their weakness, need to improving score. Each time the teacher interacts with the system, s/he takes a test, the results of which determine the student's knowledge.

International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirccce.com

Vol. 5, Issue 4, April 2017

In proposed model users have to complete registration then they have authenticated username and password respectively. When student login then he/she can give test according to which sub or category is selected by them. After completion of test result is automatically computed by system with help of mining and giving scored with some recommendation for improving the scored and giving some useful guidance.

This record is also seen by teacher for analyzing the performance of student and they can take a suitable action to improve the student performance before the student in the critical area. This technique will monitor and evaluate the student academic performance at different year levels before the final test in order to forecast the weaknesses of the students. Teacher can play the role of admin who have authority to adding subjects, topics and questions for test purpose. All of this working is simply design in form of flow chart as shown in Fig. 1.

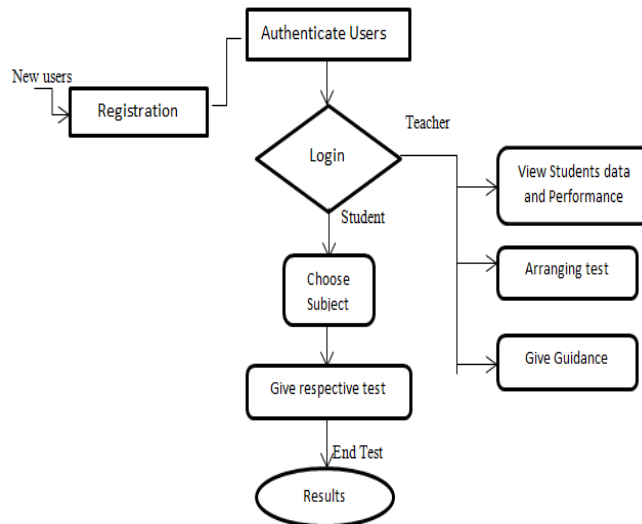


Fig. 1. Data Flow Diagram

IV. PROPOSED METHODOLOGY

The major objective of the proposed methodology is to build the classification model that classifies a students' performance. The classifiers, has been built by combining the Standard for Data Mining that includes student data and finally application of data mining techniques which is classification in present study. In other words, using this Decision tree algorithm, we wanted to be able to guide student towards achievement of good score that we felt they would enjoy doing. Tree-based methods classify instances by sorting the instances down the tree from the root to some leaf node, which provides the classification of a particular instance. Each node in the tree specifies a test of some attribute of the instance and each branch descending from that node corresponds to one of the possible values for this attribute [15]. The benefits of having a decision tree are as follows –

- It does not require any domain knowledge.
- It is easy to comprehend.
- The learning and classification steps of a decision tree are simple and fast.

Decision Tree Induction Algorithm

A machine researcher named J. Ross Quinlan in 1980 developed a decision tree algorithm known as ID3 (Iterative Dichotomiser). Later, he presented C4.5, which was the successor of ID3. ID3 and C4.5 adopt a greedy approach. In this algorithm, there is no backtracking; the trees are constructed in a top-down recursive divide-and-conquer manner.



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: www.ijircce.com

Vol. 5, Issue 4, April 2017

Decision Trees

Decision trees are commonly used in operations research, specifically in decision analysis, to help identify a strategy most likely to reach a goal. A decision tree is a flow-chart-like tree structure, where each internal node is denoted by rectangles, and leaf nodes are denoted by ovals. All internal nodes have two or more child nodes. All internal nodes contain splits, which test the value of an expression of the attributes. Arcs from an internal node to its children are labelled with distinct outcomes of the test. Each leaf node has a class label associated with it.

V. DATA MINING TECHNIQUE AND CLASSIFICATION

Data mining is very promising as a new effective technique for decision making processes. Through Educational data mining is an analysis of discipline to developing the methods for exploring the unique types of data from educational settings and it is used for improvement of students in better way [8]. Data mining techniques are applied in higher education more and more to give insights to educational and administrative problems in order to increase the managerial effectiveness. However, most of the educational mining research focuses on modelling student's performance. Data mining technique can give the input for the teachers and students about the student academic results. This technique can analysis the database patterns to forecast student performance, so this allows the teachers to prepare like a remedial program (needing extra help for learning) or more additional assignments for the students.

Several DM algorithms (naive Bayes, Bayes net, support vector machines, logistic regression, and decision trees) have been compared to detect student mental models in ITSs [16]. Unsupervised (clustering) and supervised (classification) machine learning have been proposed to reduce development costs in building user models and to facilitate transferability in intelligent learning environments. Clustering and classification of learning variables have been used to measure the online learner's motivation.

Classification is one of the most studied problems in machine learning and data mining. It consists in predicting the value of a categorical attribute (the class) based on the values of other attributes (predicting attributes). A search algorithm is used to induce a classifier from a set of correctly classified data instances called the training set. Another set of correctly classified data instances, known as the testing set, is used to measure the quality of the classifier obtained. Different kinds of models, such as decision trees or rules, can be used to represent classifiers. In Classification process, the derive model is to predict the class of objects whose class label is unknown. Generally, the classification of data has two step process are learning and a classification step which is used to predict class labels for training data. In classification step, test data are used to estimate the accuracy of classification rules. There are many techniques that can be used for classification techniques such as decision tree, Bayesian methods, Bayesian network, rule based algorithms, neural network, support vector machine, association rule mining, k-nearest- neighbor, case-based reasoning, genetic algorithms, rough sets and fuzzy logic. In this study, we focus on classification techniques such as decision tree [7].

A. Data Preparations

The data set used in this study was obtained from a student's database which we are created for our application. In this step data stored in different tables was joined in a single table after joining process errors were removed.

B. Data selection and transformation

In this step only those fields were selected which were required for data mining. A few derived variables were selected. While some of the information for the variables was extracted from the database.



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirccce.com

Vol. 5, Issue 4, April 2017

C. Decision Tree Algorithm

C4.5 is an algorithm used to generate a decision tree developed by Ross Quinlan. C4.5 is an extension of Quinlan's earlier ID3 algorithm. The decision trees generated by C4.5 can be used for classification. At each node of the tree, C4.5 chooses the attribute of the data that most effectively splits its set of samples into subsets enriched in one class or the other. The splitting criterion is the normalized information gain (difference in entropy). The attribute with the highest normalized information gain is chosen to make the decision. The C4.5 algorithm then recurs on the smaller sub lists.

All the samples in the list belong to the same class. When this happens, it simply creates a leaf node for the decision tree saying to choose that class. None of the features provide any information gain. In this case, C4.5 creates a decision node higher up the tree using the expected value of the class. Instance of previously-unseen class encountered. Again, C4.5 creates a decision node higher up the tree using the expected value. To find an optimal way to classify a learning set, what we need to do is to minimize the questions asked (i.e. minimizing the depth of the tree). Thus, we need some function which can measure which questions provide the most balanced splitting.

Algorithm: Generate_decision_tree

```
Step 1- Start
Step 2-Take input which is given by User
       $I_n = \{I_1, \dots, I_n\}$ 
Step 3-Dataset preparation
       $D_n = \{ \{I_1, \dots, I_n\} D \}$ 
Step 4-Dataset elaboration
       $D_1 = \{S_1, \dots, S_n, C_1, \dots, C_n, I_1, \dots, I_n, a_1, \dots, a_n\}$ 
Step 5- Processing
      While(  $D_n \neq 0$  )
      {
          If (  $a_n == I_n$  )
              Check  $C_n, S_n$ ;
      }
Step 6- Result Generation
```

$R = \{ S_c, S_n, C_n \}$;

Where,

I_n = Input given by users

D_n = Dataset

D = Database

D_1 = Dataset contents

S_c = Score

a_1, \dots, a_n = Answer

S_1, \dots, S_n = Subject

C_1, \dots, C_n = Category

We also using Generalized Sequential Pattern mining algorithm for predicting the student's performance as pass or fail. Once the student is found at the risk of failure he/she can be provided guidance for performance improvement.

International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: www.ijircce.com

Vol. 5, Issue 4, April 2017

Generation Sequential pattern Mining Algorithm

Step 1-Start

Step 2- Take input from Dataset

Step 3- Processing

While (Cn!=0)

{

 Qn={{ Q1,.....Qn } Cn }

 PData= {count(Qn), So(Qn), Cr(Qn)}

 Rc= PData{(count(Qn)-Cr(Qn)) Cn }

 }

Step 4- Result Generation

 R=Rc;

R will show the weak category of student.

Where,

Qn=Questions(Total)

{(Qn)Cn}= Questions regarding to category.

So= Solved Questions

Cr= Correct questions

Rc= Result Category

It provides weak categories or subjects of students from his/her performance.

VI. SIMULATION AND RESULTS

The simulation studies involve comparison of ID3 and C4.5 accuracy with different data set size, this comparison is presented graphically in Fig.2.

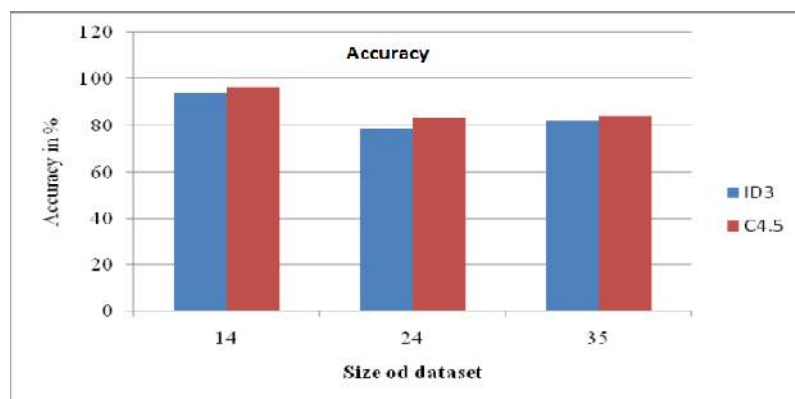


Fig.2. Comparison of Accuracy for ID3 & C4.5 Algorithm

The 2nd parameter compared between ID3 and C4.5 is the execution time which is show in Fig.3.

International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: www.ijircce.com

Vol. 5, Issue 4, April 2017

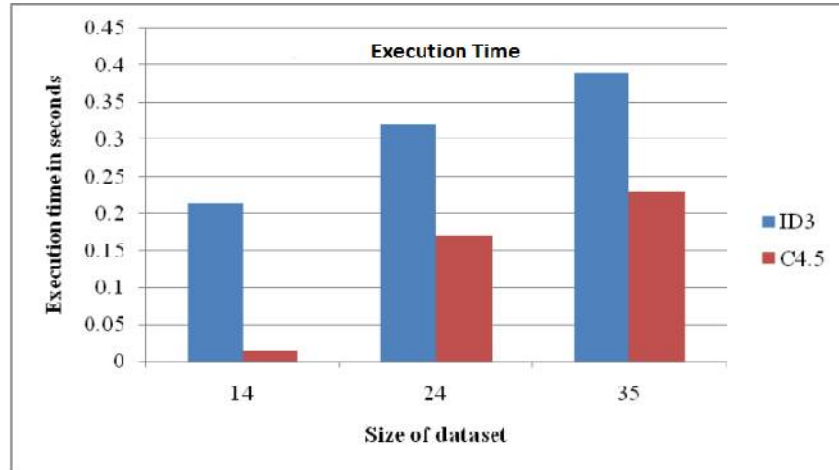


Fig.3. Comparison of Execution Time for ID3 & C4.5 Algorithm

In result student getting the all records of given test is show in Fig. 4. Figure shows the test given by student marks obtains, total marks, date oaf test and provide suggestion to study for improving performance.



Student Performance System....				
Try To Learn Every thing				
Home Result LogOut				
Result				
Test	Marks Obtain	Total Marks	Test Date	Study
C	6	20	22/04/2017	Array
Fundamental	5	20	23/04/17	Computer Media
Fundamental	3	20	24/04/17	Computer Media
C	8	20	25/04/17	loop

Fig.4. Student Record

In Fig.5 shows the student performance record of test given by that particular student which involve subject, total marks, marks obtain, date of conducting test, weak concept shoes the category in which student is weak.

International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirccce.com

Vol. 5, Issue 4, April 2017

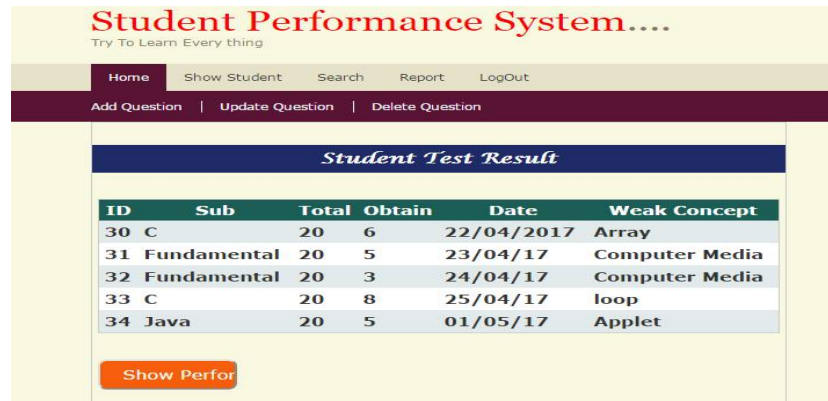


Fig.5. Student performance report

In Fig.6 indicate performance of student in graphical to teacher.

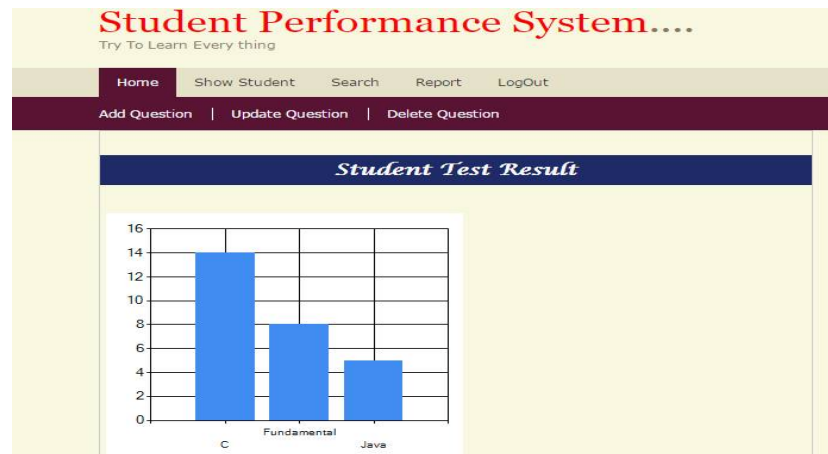


Fig.6. Performance result

VII. CONCLUSION

Academic success of students of any professional Institution has become the major issue for the management. An early analysis of students at risk of poor performance helps the management take timely action to improve their performance through extra coaching and counselling. The result of this study indicates that data mining techniques capabilities provided effective improving tools for analysis student performance. In this paper, data mining is utilized to analyse course evaluation questionnaires. Here, the most important variables that separate “satisfactory” and “not satisfactory” student performances and there weakness’ in particular subject or field. Hopefully, these can help instructors to improve their performances. Tree-based methods classify instances by sorting the instances down the tree from the root to some leaf node, which provides the classification of a particular instance. Each node in the tree specifies a test of some attribute of the instance and each branch descending from that node corresponds to one of the possible values for this attribute. This paper focuses on analysis student academic performance by using advantage of data mining techniques model.



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirccce.com

Vol. 5, Issue 4, April 2017

REFERENCES

1. Mustafa Agaoglu, "Predicting Instructor Performance Using Data Mining Techniques in Higher Education," IEEE Access , Volume: 4 ,2016.
2. Tripti Mishra,Dr. Dharminder Kumar,Dr. Sangeeta Gupta,"Mining Students' Data for Performance Prediction," in fourth International Conference on Advanced Computing & Communication Technologies,2014.
3. Keno C. Piad, Menchita Dumlao, Melvin A. Ballera, Shaneth C. Ambat," Predicting IT Employability Using Data Mining Techniques," in third International Conference on Digital Information Processing, Data Mining, and Wireless Communications (DIPDMWC), 2016.
4. Bipin Bihari Jayasingh,"A Data Mining Approach to Inquiry Based Inductive Learning Practice In Engineering Education," in IEEE 6th International Conference on Advanced Computing,2016.
5. S. M. Merchán,"Analysis of Data Mining Techniques for Constructing a Predictive Model for Academic,"IEEE Latin America Transactions, vol. 14, no. 6, June 2016.
6. Konstantina Chrysafiadi and Maria Virvou," Fuzzy Logic for adaptive instruction in an e-learning environment for computer programming," IEEE Transactions on Fuzzy Systems ,Volume: 23, Issue: 1, Feb. 2015.
7. M. Mayilvaganan,D. Kalpanadevi , " Comparison of Classification Techniques for predicting the performance of Students Academic Environment," in International Conference on Communication and Network Technologies (ICCNT), 2014.
8. Cristóbal Romero," 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, November 2010.
9. Priyanka Anandrao Patil, R. V. Mane," Prediction of Students Performance Using Frequent Pattern Tree," Sixth International Conference on Computational Intelligence and Communication Networks,2014.
10. Behrouz Minaei-Bidgoli, Deborah A. Kashy , Gerd Kortemeyer, William F. Punch, "Predicting Student Performance: An Application Of Data Mining Methods With An Educational Web-Based System," in 33'd ASEE/IEEE Frontiers in Education Conference, T2A-13,November 5-4,2003.
11. Peter Brusilovsky, Sibel Somyürek , Julio Guerra , Roya Hosseini , Vladimir Zadorozhny , Paula J. Durlach,"Open Social Student Modeling for Personalized Learning," IEEE Transactions on Emerging Topics in Computing, Volume: 4, Issue: 3, July-Sept. 2016.
12. Pedro G. Espejo, Sebasti´an Ventura, and Francisco Herrera," A Survey on the Application of Genetic Programming to Classification," IEEE Transactions On Systems, Man, And Cybernetics—Part C: Applications And Reviews, Vol. 40, No. 2, March 2010.
13. Carlos Márquez Vera, Cristóbal Romero Morales and Sebastián Ventura Soto,"Predicting of school failure and dropout by using data mining techniques", The IEEE Journal of Latin-American Learning Technologies (IEEE-RITA) , Vol. 8, No. 1, pp 7-14, Feb 2013.
14. A. Peña-Ayala, "Review: Educational data mining: A survey and a data mining-based analysis of recent works," Expert Systems with Applications, vol. 41, no. 4, pp. 1432-1462, 2014.
15. R.S.J.D Baker and K.Yacef, "The State of Educational Data Mining in 2009: A Review and Future Visions" , Journal of Educational Data Mining, 1, Vol 1, No 1, 2009.
16. V. Rus, M. Lintean, and R. Azevedo, "Automatic detection of student mental models during prior knowledge activation in MetaTutor," in Proc. Int. Conf. Educ. Data Mining, Cordoba, Spain, 2009, pp. 161-170.

BIOGRAPHY

Ankita A Nichat is a Student of M.E Computer Science & Engineering Department, H.V.P.M'S College of Engineering & Technology, Amravati, Maharashtra , India. She received Bachelor of Engineering Degree in 2015 from SGBAU Amravati, Maharashtra, India. Her research interests are Education technology and Data Mining.

Dr.Anjali B Raut is a HOD of Department of Computer Science & Engineering, H.V.P.M'S College of Engineering & Technology, Amravati , Maharashtra , India.