# Student Performance Prediction using Support Vector Machine

**Prathmesh Jadhav, Abhishek Toshniwal, Yash Golawar, Pratik Khomhane**
Sinhgad College of Engineering, Pune, Maharashtra, India

**Abstract:** *Predicting students' academic progress is critical in a university setting for early diagnosis of at-risk pupils. The goal of this work is to offer data mining models that use classification methods based on Support Vector Machine (SVM) algorithms to forecast students' academic accomplishment after a preparatory year, as well as to select the optimal approach. Based on graduation CGPA, students' academic achievement is classified as High, Average, or Below Average, and these classifications are applied to a newly produced dataset.*

**Keywords:** Support Vector Machine (SVM), computer vision, neural networks, Data Mining, Error Measurement, Accuracy.

## I. INTRODUCTION

Education is a field that generates and accumulates a lot of data. During the typical educational process, thousands of hours are spent at school and performing countless chores at home. Students' interactions with instructional materials generate a large amount of data. Education management systems and online educational platforms collect data on student participation in the online system. Statistics on students' progress and the outcomes of assignments and exercises, as well as involvement in group projects and dialogues, are examples. The university has accumulated data on its candidates throughout time, including their gender, age, and grades in numerous subjects upon high school graduation. Following that, data is collected on the same individuals, but time as students. Attendance, grades in various subjects, the success of scientific activities, the types of tasks assigned, the teachers who taught the course, and so forth.

Data mining, also known as knowledge discovery, is growing more popular since it aids in the study of data from multiple perspectives and the summary of that data into useful information. In educational data mining, decision trees, neural networks, k-nearest neighbour, naive bayes, support vector machines, and other methods are used.

On a csv dataset, the purpose of this study is to use the Support Vector Machine approach to predict and analyse student performance.

## II. RELATED WORK

To predict student achievement, numerous research have been undertaken. The Decision Tree (DT) approach was used to predict engineering student achievement in [8]. The researchers gathered information on 340 students in order to forecast how well they will fare on their first-year exam. In their training set, the model they constructed was 60% correct. [9] estimated final student grades using parameters from two different datasets using WEKA data mining methods. Each dataset includes statistics on a group of students from a specific college course during the last four semesters. Among the other classifiers, the IBK has the highest accuracy. [10] summarised previous studies that used various data mining tools to forecast, analyse, and evaluate student performance. The authors of [11] employed DT classification methods to measure student performance and created classifier models using artificial neural networks. The work was analysed using a number of criteria in attempt to predict student achievement. Analyzing the learner's strengths and weaknesses in order to improve future performance. This conclusion emphasises the importance of data mining approaches in course evaluation and higher education data mining. The writers of [12] provide a study that will aid students and teachers in improving the grades of at-risk kids. Data from the student's previous database, including attendance, seminar attendance, and assignment grades, was utilised to forecast performance at the end of the semester. When compared to other categorization algorithms, the authors employed the Nave Bayes algorithm, which had the highest accuracy. Multiple decision tree algorithms were tested on an educational dataset to classify students' educational performance, according to the researchers in [13]. The study's main goal

is to find the best algorithm among the most extensively used decision tree algorithms and to provide a standard for each one. On the studied dataset, the Classification and Regression Tree strategy outperformed the other methods based on the accuracy and precision obtained by 10-fold cross validations.

[14] An overview of data mining approaches for forecasting student performance was provided, with a focus on how to use the prediction algorithm to uncover the most significant components of a student's data. To predict student achievement, researchers frequently employ neural networks and decision trees as categorization algorithms. [15] Future outcomes, as well as the factors that influence them, were found and assessed using data mining approaches. The analysis was carried out by applying the FP Growth Algorithm to determine the Association rules for the same, which were then ranked by Lift Metric. The Rule Based Induction Method was then used to classify the data.

## III. DATA SET AND DESCRIPTION

Using various classification approaches, the student's The performance of students was predicted using a college database. Each dataset contains 396 instances, each with 12 characteristics. Table 1 lists the characteristics, along with their descriptions and values. The following attributes were taken into account for prediction:

Travel time:
  Home to school or college travel time.
Type numeric: 1- <15minutes, 2- 15-30minutes, 3- 30-1hour, 4- >1hour.

Study time:
  Weekly avg. study time of student.
    Type numeric: 1- <2hour, 2- 2-5hour, 3- 5-10hour,
    4- >10hour

Attended Extra-class:
  Whether student has attended any extra class.
    Type Boolean: Yes or No

Extra-curricular:
  Whether student has taken part in any Extra-curricular        activities.
    Type Boolean: Yes or No

Tuition:
  Whether student has taken any tuition other than college.
    Type Boolean: Yes or No

Health Condition:
  Students' health condition throughout semester.

Absent days:
  For how many days student was absent.
    Type numeric:0-90 days

Test 1 Mark:
  Marks Scored by student in test 1.
  Type numeric: 0-100
Test 2 Mark:
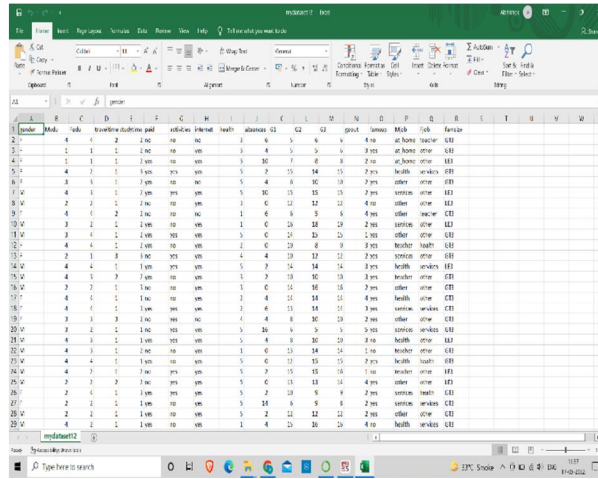  Marks Scored by Student in test 2.
    Type numeric: 0-100

**Figure 1:** Data description and values

## IV. SVM CLASSIFICATION

In this step, the data was classified using the classification method (SVM), and the students' status was then supplied. The preceding phase's test and training data records were collected and used to classify the data using the SVM algorithm. A database was used to compile the data. Several of the training record sets were tested. After applying the SVM classifier function to the data sets, the classification information were obtained. The accuracy of the prognosis for overall functioning was established after these classification details were analysed to acquire the results. Finally, based on the conditions, the students were divided into groups. Internal Marks for certain theoretical courses, Internal Lab Marks, whether he or she engages in sports, department activities, or is experiencing any psychological stress as a result of the college or family environment, and so on.

It's a promising data classification approach for both linear and non-linear data. It uses non-linear mapping to translate the original training data into higher dimensions. Students are referred to as multidimensional things since they are defined by a variety of aspects. Within this additional dimension, it finds the linear optimal separation hyperplane (that is, a "decision boundary" separating pupils from one class from another). A hyperplane can be used to split data into two groups at any time (H1 and H2). The SVM uses support vectors ("critical" preliminary information samples) and edges to locate this hyperplane (Large edge and small edge, which is characterised by help vectors).
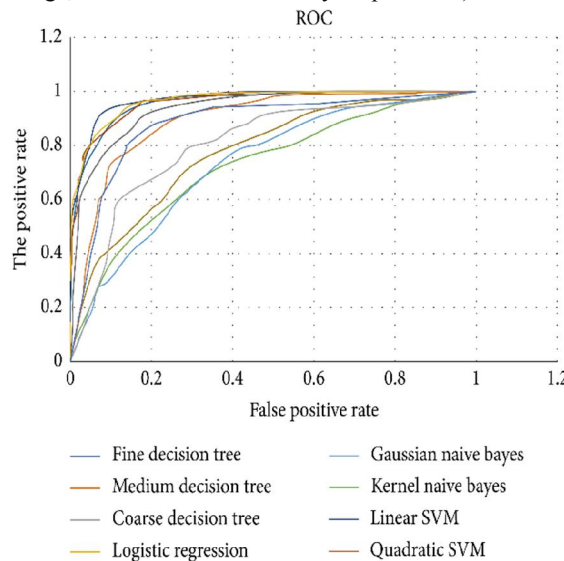


**Figure 2:** FP rate

382

When compared to other classifiers, SVM has a higher True positive rate, as shown in the graph above.

## V. EXPERIMENTAL RESULTS AND DISCUSSION

The 396 instances were the focus of the experiment and debate. Our mission is to discover students who are at the top of their class in any graduating programme. For this, we obtained data from both first-year and second-year or final-year students. Because SVM gives the best results for small numbers of vectors, the dataset was likewise constrained to 396 students. The students' personal information was taken from their database. The students' attendance status was also retrieved from the same database. Because it has an impact on a student's grade

Every institution/university requires that each of its students complete the bare minimum of tutoring hours for each topic. Assume a student does not regularly attend classes due to a high level of involvement in other academic endeavours. The poor attendance rate of a student will almost likely affect his or her academic progress. As a result, their professors should issue a warning to keep their attendance within the "safe" range. Questionnaires were utilised to collect data from students. The following is the format of the questionnaire. This information is fed into the system.

**Model of a Questionnaire:**
- Your Name:
- Father's Name:
- Father's Occupation:
- Mother's Name:
- Mother's Occupation:
- ADDRESS:
- DOB:
- Email id:
- Gender:
- Age:
- Course:
- Course Completion:
- Joined any Tuition:
- Attended Extra classes:
- Travel Time:
- Study Time(weekly):
- Health Condition:
- Co-Curricular Activities:
- Absent days:
- Marks Scored in Test1:
- Marks Scored in Test 2:

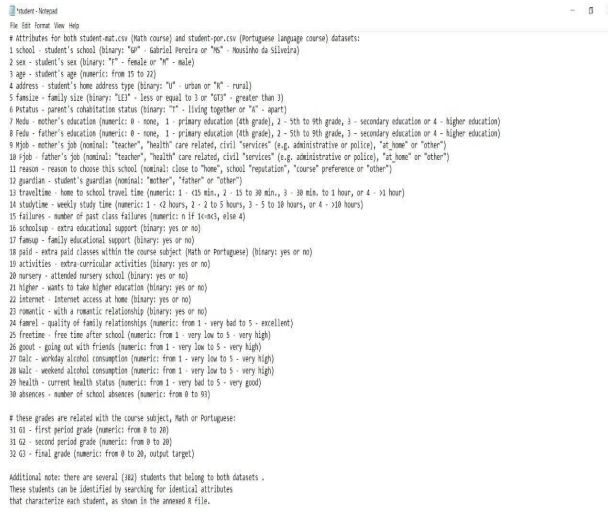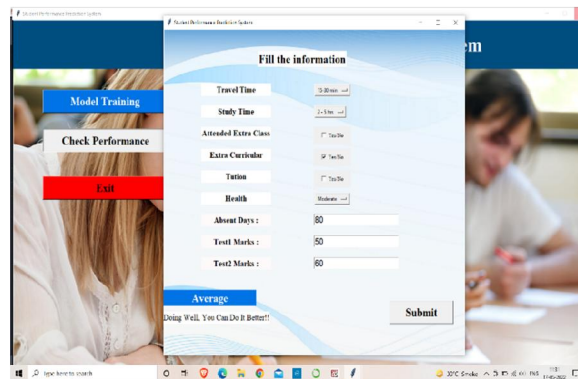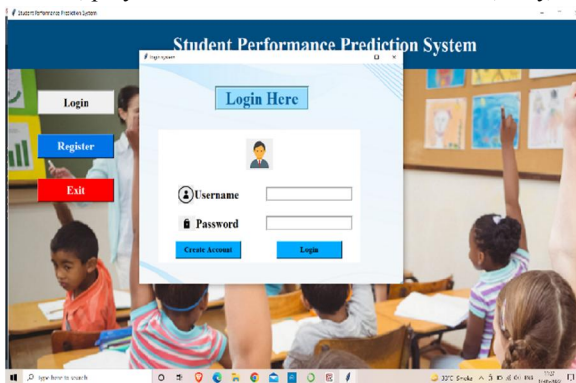Figure 3 shows a number of other similar questions that were posed.



**Figure 3:** Questionnaire Model with Description

During the pre-processing step, only a few samples were chosen from the data input. In most cases, variables that are unrelated to the forecast are ignored. DOB and E-mail ID, for example, are not included in any partitioning or categorization algorithms. Occasionally, the data values obtained are insufficient or unavailable. These kinds of samples should all be ignored. A range of factors, such as parents' education and career, may indirectly improve a student's performance. They do, however, play a minor role in the film. As a result, they, too, omitted.

## VI. CONCLUSION

The effectiveness of DM approaches to predict students' development after a year of preparation in higher education was investigated in this model. A student's academic success is measured using the Grade Point Average (GPA) (excellent, average, or poor). Throughout the experiment, we used SVM classifiers on the student dataset to predict the student's academic year achievement. The SVM classifier outperformed other strategies in predicting student achievement, according to the findings. Additionally, grades on the first and second prelim tests, as well as the final exam, Skills course, and extracurricular skills course, all had a part in predicting students' academic progress. By giving out timely warnings to students, the findings will assist predict students' final grades early enough for suitable interventions to be undertaken. With competent counselling, the number of low-achieving pupils can be reduced.

## REFERENCES

[1]. M. Goyal and R. Vohra, IJCSI International Journal of Computer Science Issues, Vol. 9, Issue 2, No. 1, March 2012, "Applications of Data Mining in Higher Education."

[2]. R. Huebner, Studies in Higher Education Journal, "A study of educational data mining research," 2012.

[3]. M.S. Mythili and A.R. Mohamed Shanavas, "An Analysis of Students' Performance Using Classification Algorithms," IOSR, Journal of Computer Engineering, Volume 16, Number 1, January 2014.

[4]. S. Lakshmi Prabha and A.R.Mohamed Shanavas, "Educational data mining applications," ORAJ, Vol. 1, No. 1, August 2014.

[5]. "Data mining in course management systems: Moodle case study and tutorial," Computers & Education, vol. 51, no. 1, 2008, pp. 368-384.

[6]. S. Ayesha, T. Mustafa, A. Sattar, and M. Khan, "Data mining model for higher education system," European Journal of Scientific Research, vol.43, no.1, pp.24-29, 2010.

[7]. Z. J. Kovacic, "Early prediction of student success: Mining student enrollment data," InSITE 2010 Proceedings.

[8]. I. Milos, S. Petar, V. Mladen, and A. Wejdan, INFOTEH- JAHORINA Vol. 15, March 2016, 684.

[9]. P.Kavipriya, A Review on Predicting Students' Academic Performance Using Data Mining Techniques Earlier, International Journal of Advanced Research in Computer Science and Software Engineering, Volume 6, Issue 12, December 2016 ISSN: 2277 128X.

[10]. N. Ankita and R. Anjali, Student Performance Analysis Using Data Mining Technique, International Journal of Innovative Research in Computer and Communication Engineering, Vol. 5, Issue 1, January 2017.

[11]. P. Shruthi and B. Chaitra, Student Performance Prediction in the Education Sector Using Data Mining, Volume 6, Issue 3, March 2016, International Journal of Advanced Research in Computer Science and Software Engineering.

[12]. 2012, S.K. Yadav, B. Bharadwaj, and S. Pal. A Comparative Study of Data Mining Applications for Predicting Student Performance Vol. 1, No. 12 of the International Journal of Innovative Technology and Creative Engineering (ISSN: 2045-711),A. Mohamed Shahiria, W. Husaina, N. Abdul Rashida, "A Review on Predicting Student's Performance using Data Mining Techniques" Procedia Computer Science 72 (2015) 414 – 422, ELSEVIER. December.

[13]. K. Kohli and S. Birla, " Data Mining on Student Database to Improve Future Performance", International Journal of Computer Applications (0975 – 8887), Volume 146 – No.15, July 2016.

[14]. Hilal Almarabeh, "Analysis of Students' Performance by Using Different Data Mining Classifiers", International Journal of Modern Education and Computer Science (IJMECS), Vol.9, No.8, pp.9-15, 2017.DOI: 10.5815/ijmecs.2017.08.02

[15]. K. B. Eashwar, R. Venkatesan and D. Ganesh, Student Performance Prediction Using Svm, International Journal of Mechanical Engineering and Technology 8(11), 2017, pp. 649–662.