



SASTRA

ENGINEERING · MANAGEMENT · LAW · SCIENCES · HUMANITIES · EDUCATION

DEEMED TO BE UNIVERSITY

(U/S 3 of the UGC Act, 1956)



THINK MERIT | THINK TRANSPARENCY | THINK SASTRA

T H A N J A V U R | K U M B A K O N A M | C H E N N A I

Students' Performance Prediction using Machine Learning algorithms

Under the guidance of :

Mr. EASHWAR K B

Asst. Professor – II

Department of CSE

Presented by :


Suwetha S

224058033

II-M.Sc. Computer Science

BASE PAPER

Trends in Neuroscience and Education 33 (2023) 100214




ELSEVIER

Contents lists available at [ScienceDirect](https://www.sciencedirect.com)

Trends in Neuroscience and Education

journal homepage: www.elsevier.com/locate/tine




Research paper

Student course grade prediction using the random forest algorithm: Analysis of predictors' importance

Mirna Nachouki^{*}, Elfadil A. Mohamed, Riyadh Mehdi, Mahmoud Abou Naaj

Artificial Intelligence Research Centre, Department of Information Technology, Ajman University, UAE



Check for updates

OBJECTIVE

Students' performance prediction uses analytical methods, including machine learning and statistics, to estimate students' future academic performance. It predicts student performance by considering the parameters that past-grades, attendance records, socio-economic background, and study attitudes and health records. By analyzing different datasets derived from that student transcripts

ABSTRACT

This proposed project is aimed to predict the students' performance using the following machine learning algorithms: Random Forest, Support Vector Machine (SVM) and Gradient Boosting. By analyzing different dataset derived from student transcripts, it is aimed to understand how useful each algorithm is in forecasting the students' academic performance. Through a detailed validation the strengths and limitations of these predictive abilities of Random Forest, Gradient Boosting and Support Vector Machine (SVM) within the environment of student academic performance. This study increases to the conversation on improving predictive methodologies in education, providing valuable understanding for institutions looking to increase student success through knowledgeable decision making and personalized interventions.

LITERATURE REVIEW

| Year | Title | Journal Name | Techniques used | Limitations |
|------|---|-----------------------------|--|---|
| 2018 | Early segmentation of students according to their academic performance: A predictive modelling approach | Decision Support Systems | decision trees, support vector machines, naive Bayes, bagged trees and boosted trees | The data is from 2003 to 2015, so the study may not be very helpful for understanding the current education situation since things may have changed |
| 2019 | Prediction of academic performance associated with internet usage behaviors using machine learning algorithms | Computers in Human Behavior | decision tree, neural network and support vector machine | None of the problems or difficulties related to this approach are discussed in this paper |

CONT.,

| Year | Title | Journal Name | Techniques used | Limitations |
|------|---|--------------|---|--|
| 2020 | Multiple Features Fusion Attention Mechanism Enhanced Deep Knowledge Tracing for Student Performance Prediction | IEEE Access | Recurrent Neural Network (RNN) | What should be done next or how to improve the study is not discussed in this paper |
| 2021 | Multiclass Prediction Model for Student Grade Prediction Using Machine Learning | IEEE Access | Decision Tree, Support Vector Machine (SVM), Naïve Bayes, K- Nearest Neighbor (kNN), Logistic Regression and Random Forest | This paper discusses two approaches for selecting particular characteristics, but it provides insufficient information about these approaches |

CONT.,

| Year | Title | Journal Name | Techniques used | Limitations |
|------|---|--------------------------------------|--|---|
| 2022 | A machine learning prediction of academic performance of secondary school students using radial basis function neural network | Trends in Neuroscience and Education | Bayes net, decision tree, k-nearest neighbors, logistic regression, naive Bayes, random forest and random tree | The data is only from students in Bangladesh, making it unsuitable for use in any other educational setting or location |
| 2023 | Machine Learning Algorithms based Student Performance Prediction based on Previous Records | IEEE Access | Bayesian classification | The Bayesian classifier technique is only used in this paper |

WORKFLOW



DATASET DESCRIPTION

Sample Dataset

Source: Real-time Data from Students'

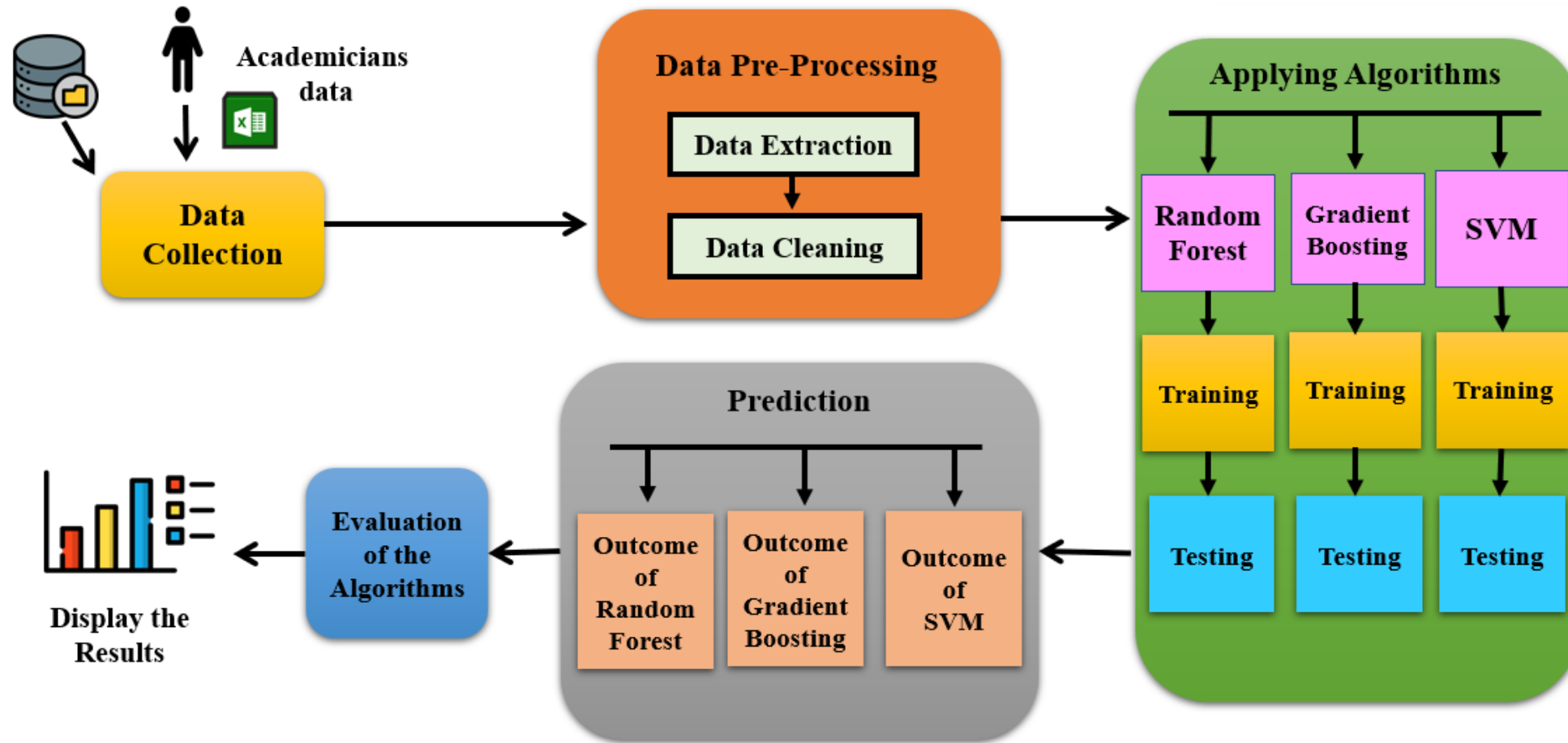
| A | B | C | D | E | F | G | H | I | J | K |
|------------|-------------|-----|--------|------------------------------|-------------------------|-------------------------|------------------------------------|---------------|--|---|
| Name | Register No | Age | Gender | Your Current Course of Study | 10th Board of Education | 12th Board of Education | 12th Specialization | Mode of Study | If you are a Day Scholar and traveling some km (kilometres) on every day, What are the problems you face ? | What is your current attendance percentage(%) ? |
| Aarthika S | 224026002 | 19 | Female | 3rd B.Com | State Board | State Board | Commerce with Computer application | Day Scholer | Headache | 80 - 90 |
| Abinaya K | 224026004 | 19 | Female | 3rd B.Com | State Board | State Board | Commerce with Accountancy | Day Scholer | Weight loss because of traveling | 90 above |
| Abinaya.M | 225026005 | 20 | Female | 3rd B.Com | State Board | State Board | Commerce with Accountancy | Day Scholer | Headache | 80 - 90 |
| Abinaya R | 224026009 | 20 | Female | 3rd B.Com | State Board | State Board | Commerce with Accountancy | Day Scholer | Weight loss because of traveling | 70 - 80 |

CONT.,

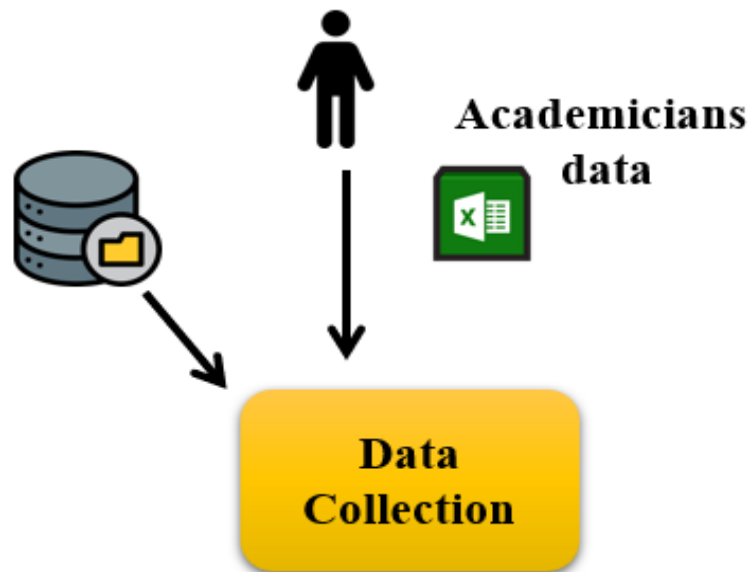
| L | M | N | O | P | Q | R | S |
|--|---|--|--------------------------------------|---|---|--|--|
| How is your Grade Point Average calculated ? | Does family situations affect your studies? | If you said "Yes", which factor influences you ? | What is your family's annual income? | If you get a low grade in a course, what do you think could be the reason ? | If you are less focused on your studies and your friends score higher than you, does that discourage you? | What kind of course do you struggle with the most? | Why are you struggling with that course? |
| Test in classes | Maybe | Education Investment | 10000 upto 50000 | Hopelessness | No | Managing Business Process | Lack of interest in course |
| Internal Marks | No | Nothing | 10000 upto 50000 | Anxiety | No | Business Law | Difficulty Level |
| Internal Marks | No | Nothing | 10000 upto 50000 | Loneliness | No | Insurance | Difficulty Level |
| Assignments | No | Education Investment | 50000 upto 1,00,000 | Fear about Future | No | Business Law | Difficulty Level |

| T | U | V | W | X | Y | Z | AA | AB | AC |
|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| 1st Sem SGPA | 1st Sem CGPA | 2nd Sem SGPA | 2nd Sem CGPA | 3rd Sem SGPA | 3rd Sem CGPA | 4th Sem SGPA | 4th Sem CGPA | 5th Sem SGPA | 5th Sem CGPA |
| 6.1539 | 6.1539 | 8 | 7.1273 | 8 | 7.4146 | 8.3929 | 7.6636 | 8.7037 | 7.8686 |
| 4.6539 | 4.6539 | 6.931 | 6.4 | 7.7037 | 6.8293 | 7.4643 | 6.9909 | 7.7037 | 7.1314 |
| 7.6539 | 7.6539 | 8.2414 | 7.9636 | 8.0741 | 8 | 7.8571 | 7.9636 | 8.6667 | 8.1022 |
| 7.2692 | 7.2692 | 7.8966 | 7.6 | 7.8519 | 7.6829 | 7.8929 | 7.7364 | 8 | 7.7883 |

PROPOSED ARCHITECTURE



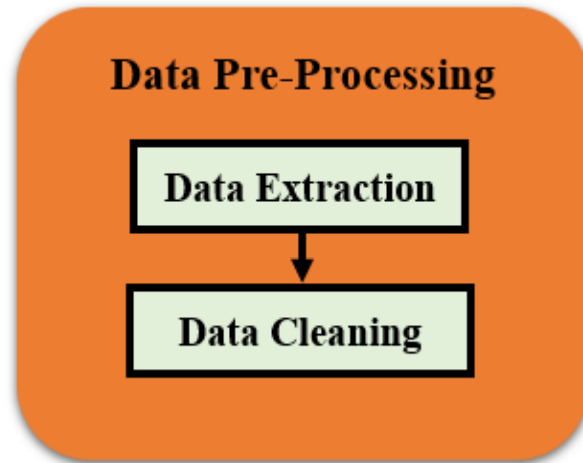
MODULE 1 : DATA COLLECTION



Data Collection

- The information was gathered from students using a questionnaire with some questions
- Additionally, the information was collected from the student academic repository. This was the way the data were collected

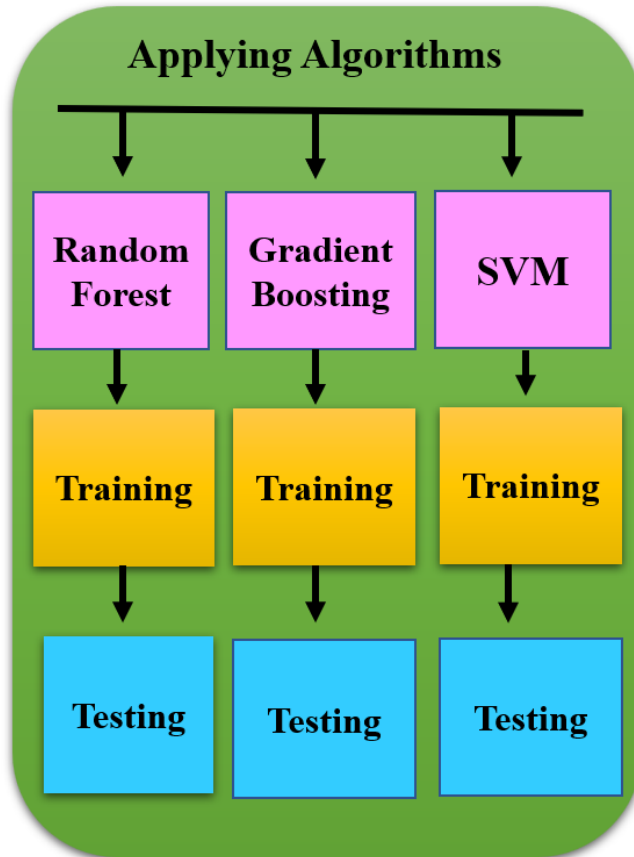
MODULE 2 : DATA PRE-PROCESSING



In this stage, using the following two methods, Data Extraction and Data Cleaning, the Data Pre-Processing was performed.

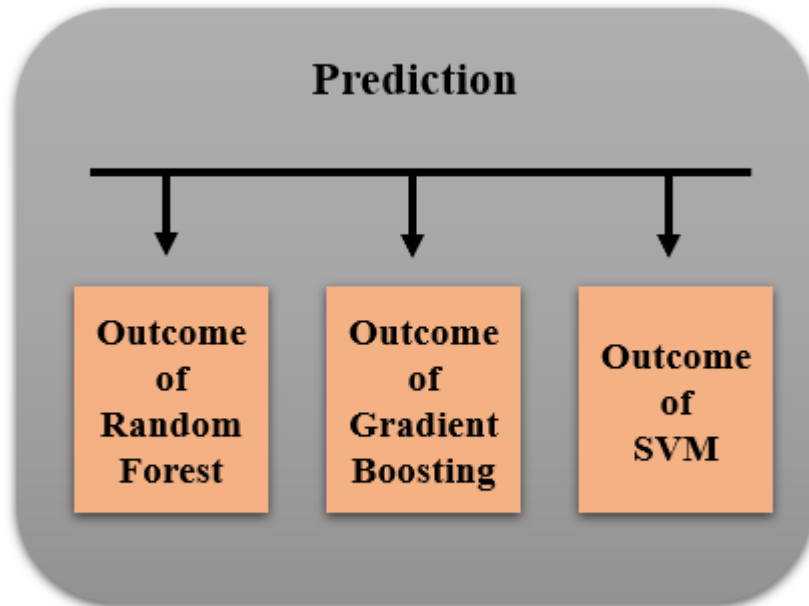
- **Data Extraction:** It was the process of collecting data samples from different sources. For example, Student academic repository and Spreadsheets documents
- **Data Cleaning:** A dataset sometimes contained entries with incomplete information. Those datasets or entries were not regarded for further processing. There existed specific techniques to impute fill in the missing data for these incomplete entries

MODULE 3 : APPLYING ALGORITHMS



- In this stage, various machine learning algorithms were applied: Random Forest, Gradient Boosting, and Support Vector Machine (SVM)
- Additionally, the dataset was divided into two subsets: one for training the models and another for testing the models, both derived from the same dataset

MODULE 4 : PREDICTION

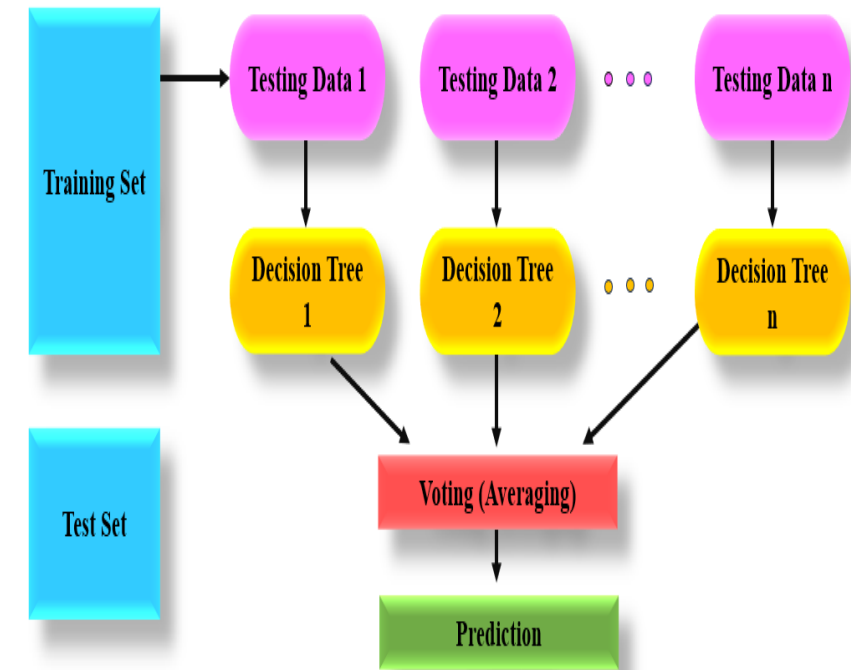


- In this stage, the outcomes of the Random Forest algorithm, Gradient Boosting algorithm, and Support Vector Machine (SVM) algorithm were obtained separately
- Predictions of results were obtained through various evaluation techniques

ALGORITHM USED

1) Random Forest

- Random Forest is an ensemble learning algorithm used for both classification and regression tasks
- It builds multiple decision trees during training and combines their predictions to make more accurate and robust predictions
- Each tree in the forest is trained on a random subset of the data and makes independent predictions and the final outputs are determined by a majority vote or averaging depending on the task



CONT.,

2) Gradient Boosting

- Gradient Boosting is a popular boosting algorithm in machine learning used for classification and regression tasks
- Gradient Boosting is also a powerful algorithm that utilizes boosting, where it constructs an ensemble of decision trees is acceptable for handling large datasets

3) Support Vector Machine

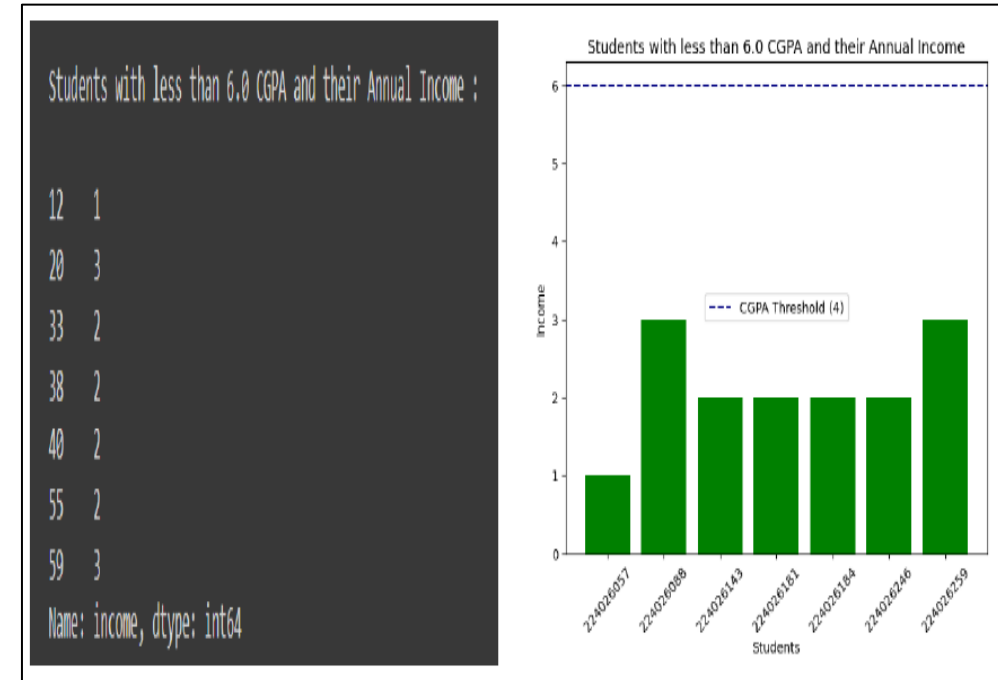
- A Support Vector Machine(SVM) is a supervised machine learning algorithm used for classification and regression tasks
- It works by finding the hyperplane that best separates data points into different classes, with a maximum margin between the classes

OUTPUT FOR 3rd YEAR

For 3rd Year Students':

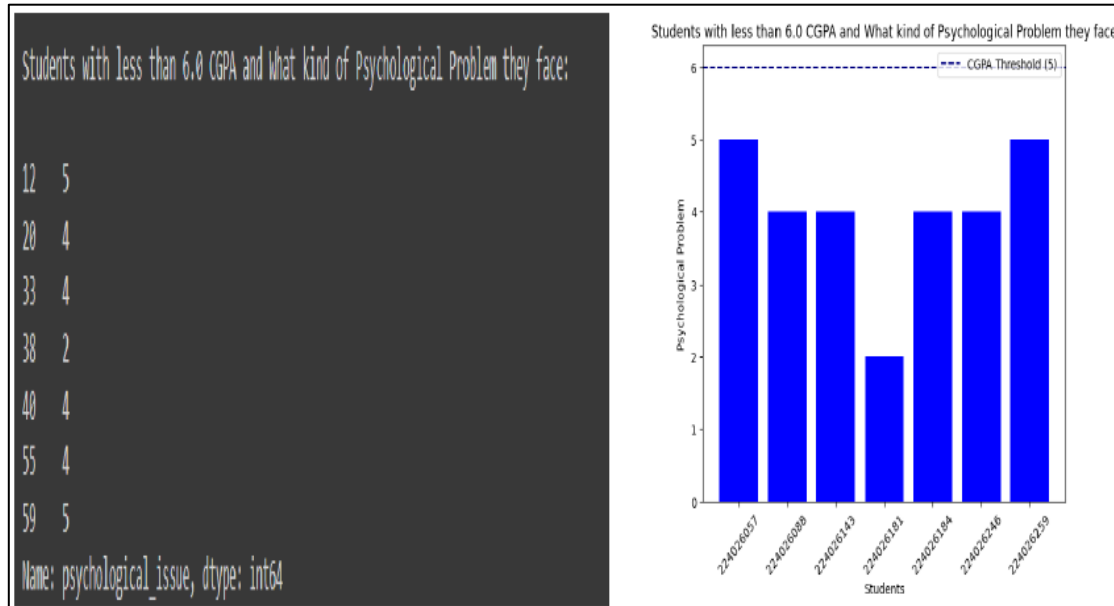


3rd year students' with CGPA below 6.0

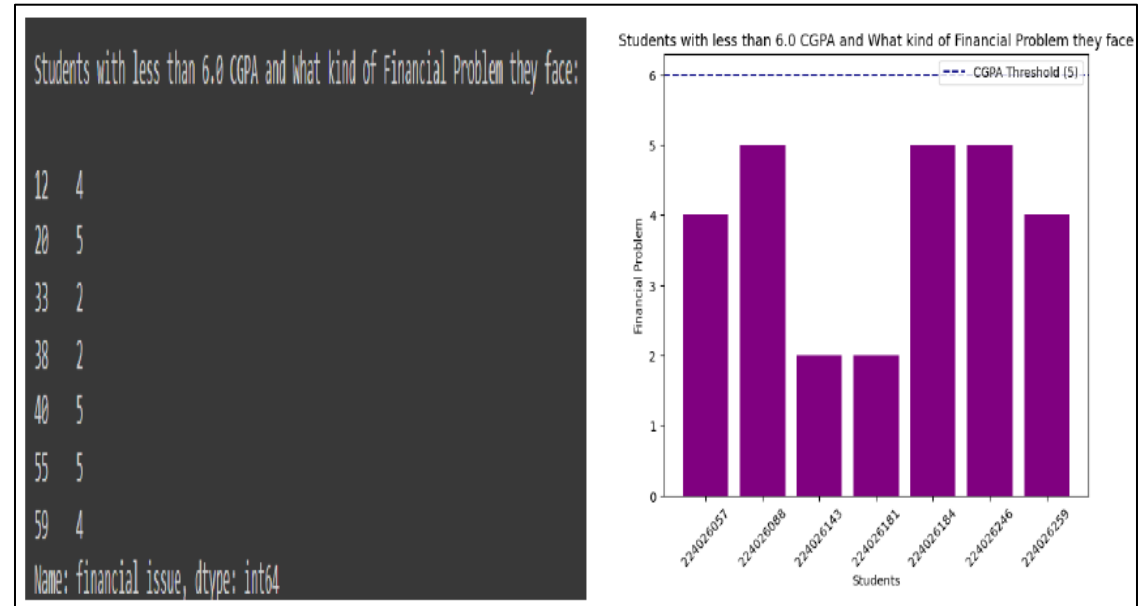


3rd year students' with CGPA below 6.0 and their annual income

CONT.,



3rd year students' with CGPA below 6.0 and what kind of Psychological problems they may face

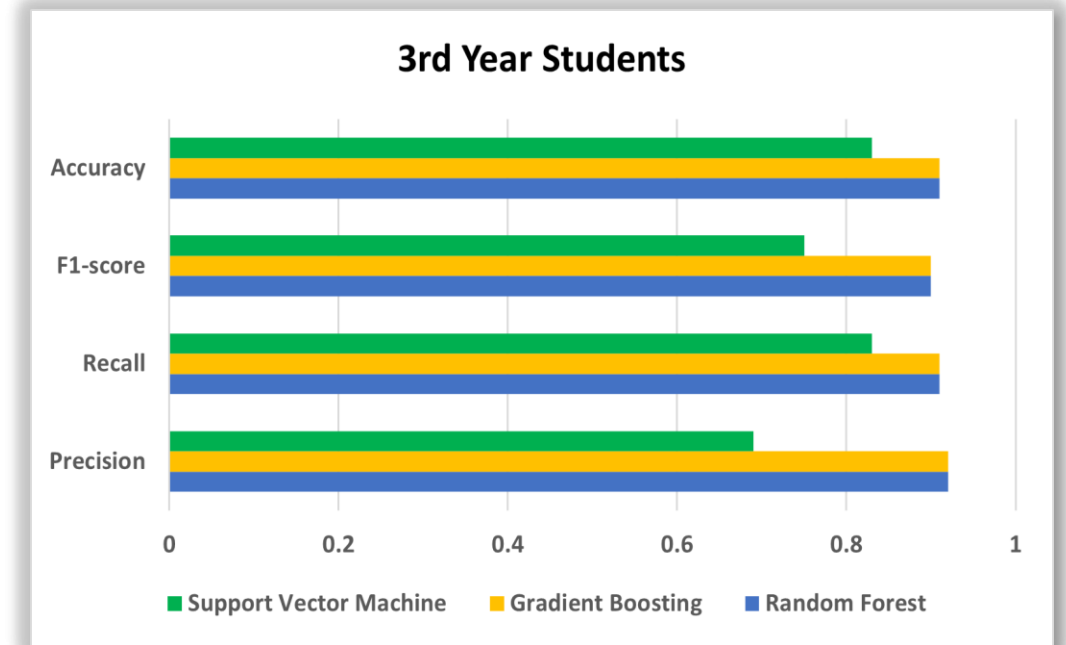


3rd year students' with CGPA below 6.0 and what kind of Financial problems they may face

CONT.,

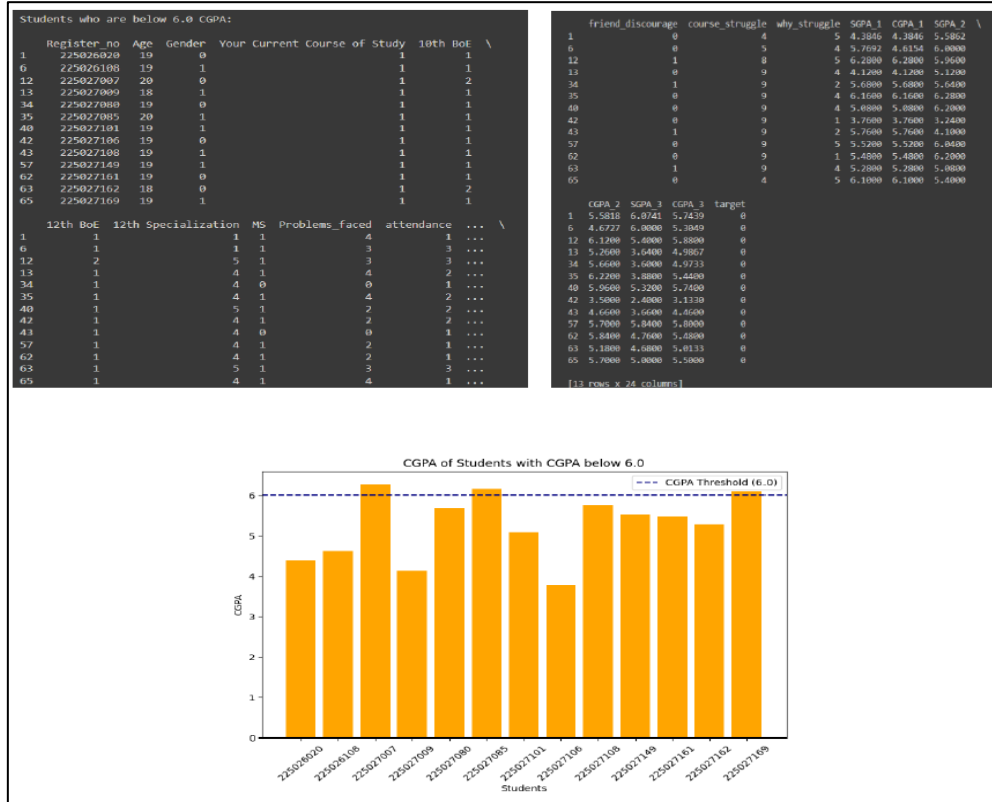
3rd Year Students' Results:

| 3 rd Year Students' | | | | |
|--------------------------------|-----------|--------|----------|----------|
| Algorithms | Precision | Recall | F1-Score | Accuracy |
| Random Forest | 0.92 | 0.91 | 0.9 | 0.91 |
| Gradient Boosting | 0.92 | 0.91 | 0.9 | 0.91 |
| Support Vector Machines | 0.69 | 0.83 | 0.75 | 0.83 |

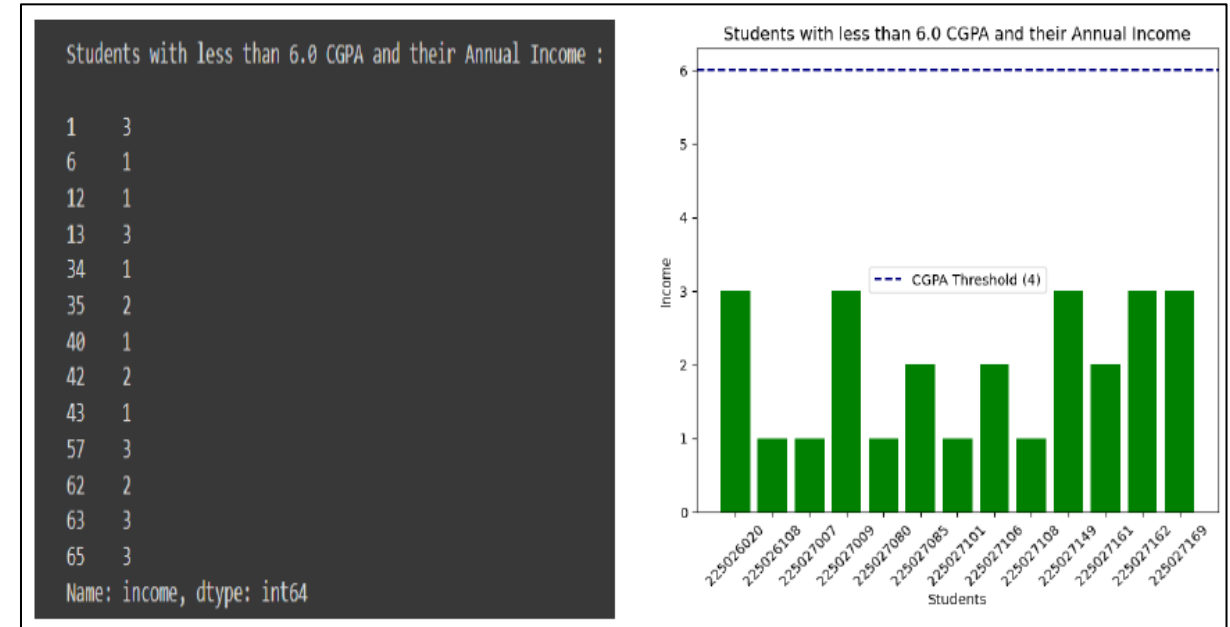


OUTPUT FOR 2nd YEAR

For 2nd Year Students':

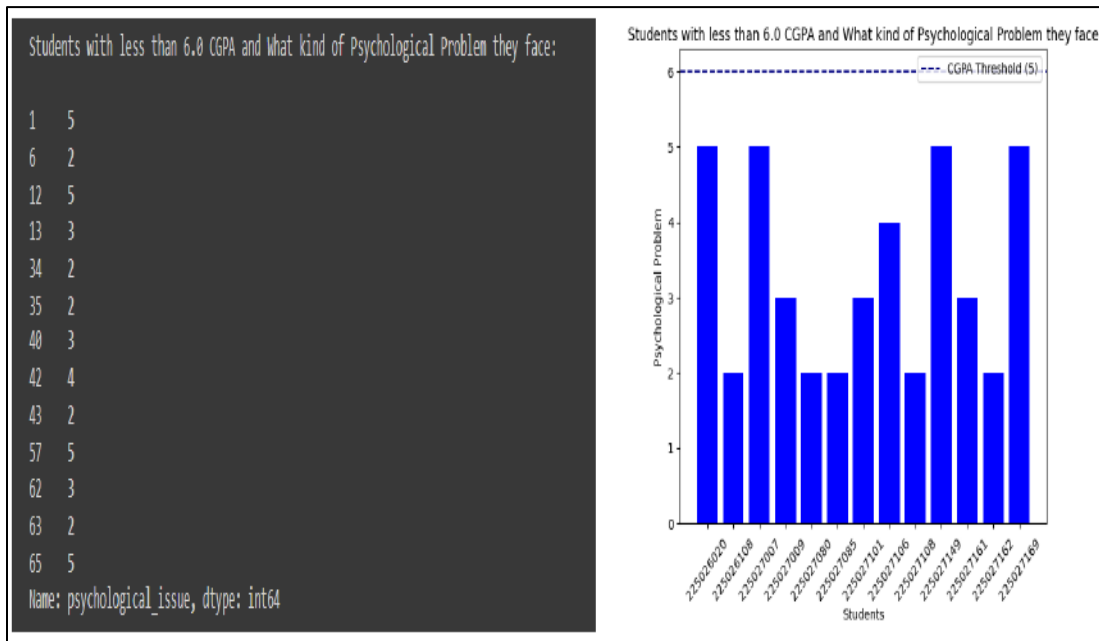


2nd year students' with CGPA below 6.0

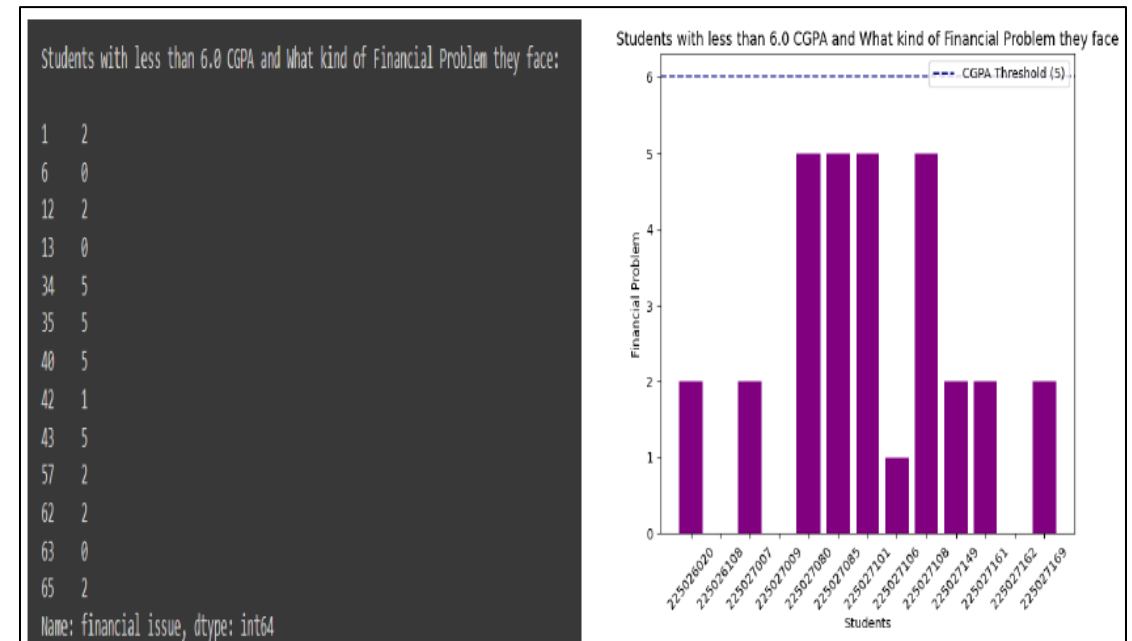


2nd year students' with CGPA below 6.0 and their annual income

CONT.,



2nd year students' with CGPA below 6.0 and what kind of Psychological problems they may face

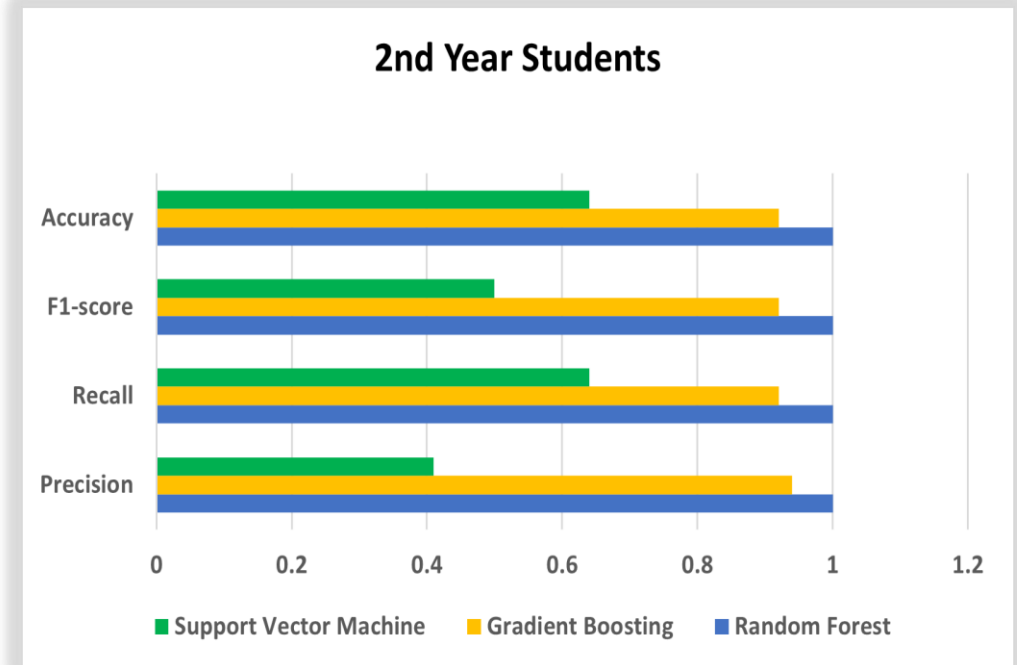


2nd year students' with CGPA below 6.0 and what kind of Financial problems they may face

CONT.,

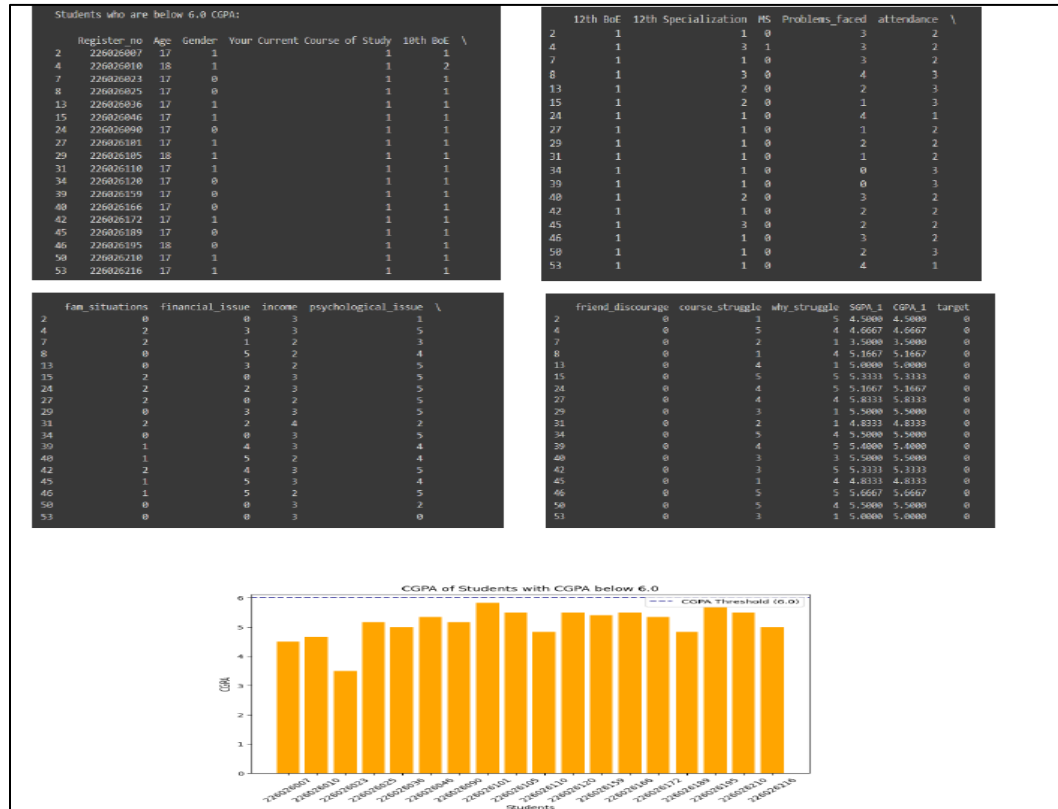
2nd Year Students' Results:

| 2 nd Year Students' | | | | |
|--------------------------------|-----------|--------|----------|----------|
| Algorithms | Precision | Recall | F1-Score | Accuracy |
| Random Forest | 1.0 | 1.0 | 1.0 | 1.0 |
| Gradient Boosting | 0.94 | 0.92 | 0.92 | 0.92 |
| Support Vector Machines | 0.41 | 0.64 | 0.5 | 0.64 |

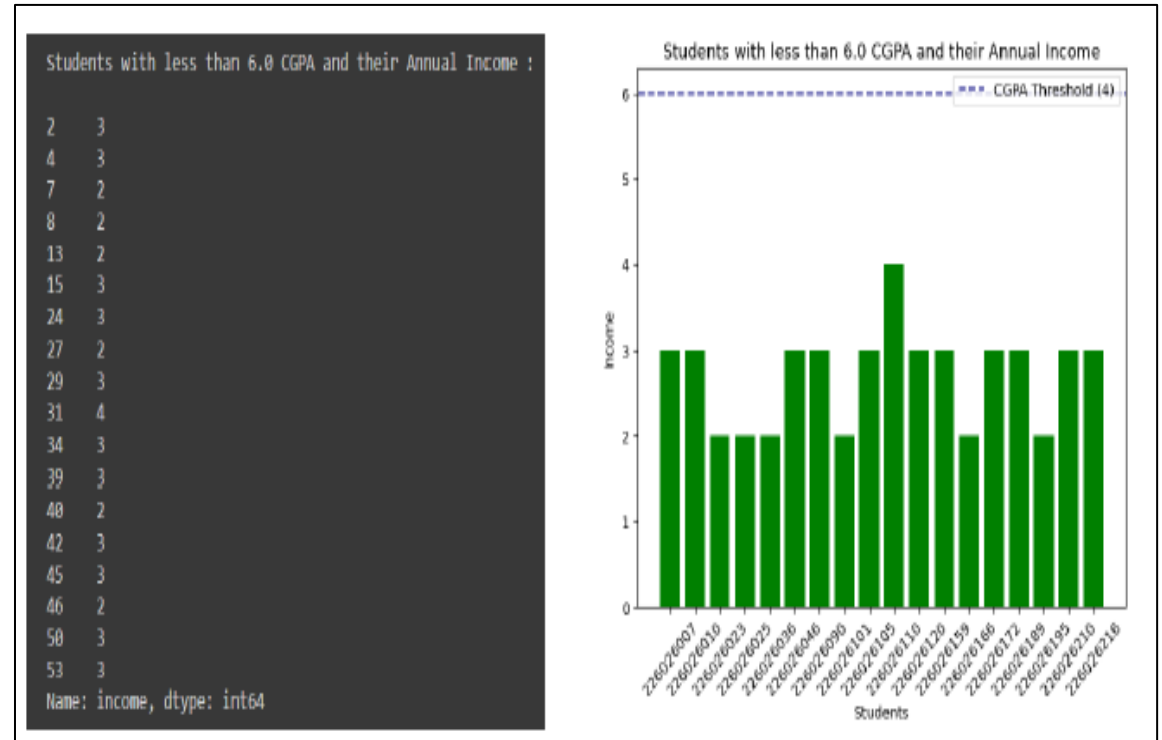


OUTPUT FOR 1st YEAR

For 1st Year Students':

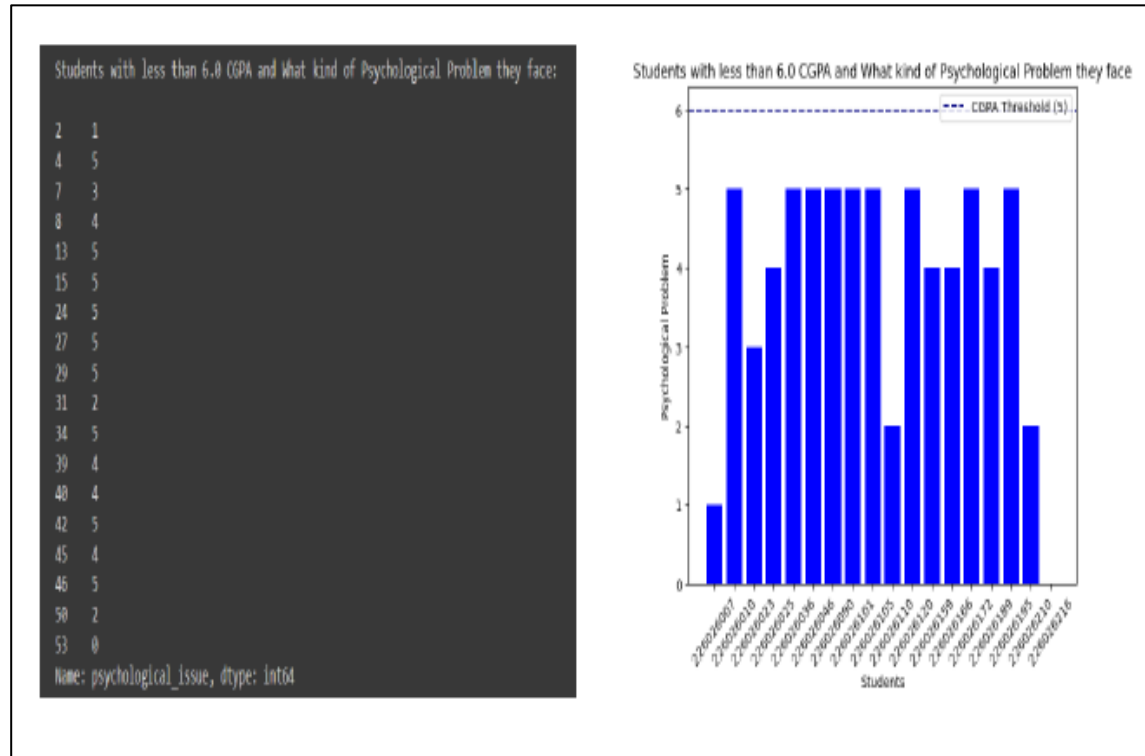


1st year students' with CGPA below 6.0

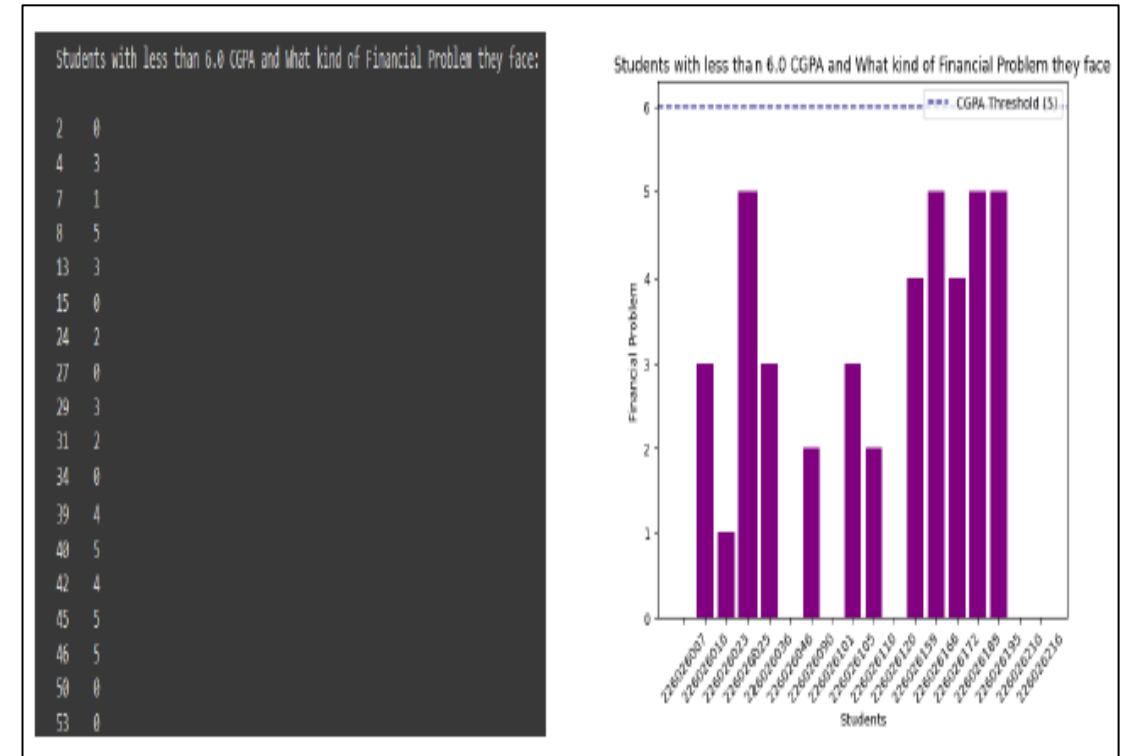


1st year students' with CGPA below 6.0 and their annual income

CONT.,



1st year students with CGPA below 6.0 and what kind of Psychological problems they may face

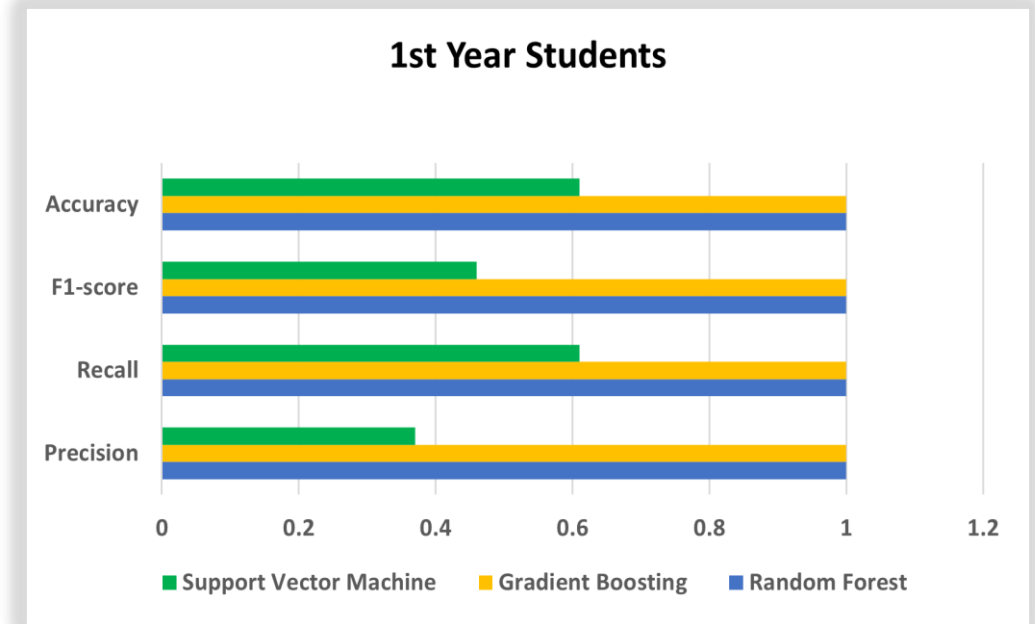


1st year students with CGPA below 6.0 and what kind of Financial problems they may face

CONT.,

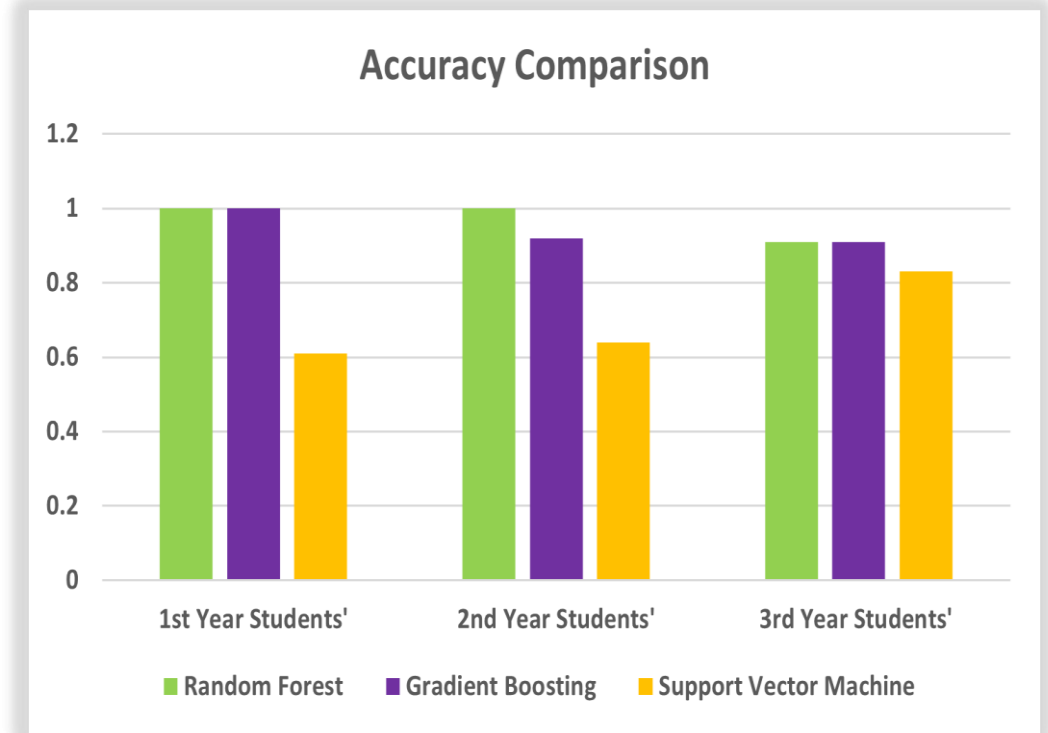
1st Year Students' Results:

| 1 st Year Students' | | | | |
|--------------------------------|-----------|--------|----------|----------|
| Algorithms | Precision | Recall | F1-Score | Accuracy |
| Random Forest | 1.0 | 1.0 | 1.0 | 1.0 |
| Gradient Boosting | 1.0 | 1.0 | 1.0 | 1.0 |
| Support Vector Machines | 0.37 | 0.61 | 0.46 | 0.61 |



ACCURACY COMPARISON

| Model | 1 st Year Students' | 2 nd Year Students' | 3 rd Year Students' |
|------------------------|--------------------------------|--------------------------------|--------------------------------|
| Random Forest | 1.0 | 1.0 | 0.91 |
| Gradient Boosting | 1.0 | 0.92 | 0.91 |
| Support Vector Machine | 0.61 | 0.64 | 0.83 |



CONCLUSION

In conclusion, for the prediction of student performance using machine learning algorithms, we gathered information from students to construct our dataset. We employed three machine learning algorithms: Random Forest, Gradient Boosting, and Support Vector Machine. Among these algorithms, Random Forest exhibited the highest accuracy. Our dataset was divided into three partitions: 1st year, 2nd year, and 3rd year. For 1st-year students, we achieved comparable levels of accuracy with both Random Forest and Gradient Boosting. In the case of 2nd-year students, Random Forest yielded the highest accuracy. Similarly, for 3rd-year students, Random Forest and Gradient Boosting demonstrated equal levels of accuracy. Overall, Random Forest consistently outperformed the other algorithms in terms of accuracy. Furthermore, this analysis can assist education management in identifying student performance based on various factors such as family situation, psychological aspects, and financial issues.

Requirements

Hardware Specifications :

| | |
|-----------|---|
| Processor | : 12th Gen Intel(R) Core(TM) i51235U 1.30 GHz |
| Hard Disk | : 512 GB |
| RAM | : 8.00 GB |

Software Specification :

| | |
|----------------------|------------------|
| OS | : Windows 11 |
| Programming Language | : Python |
| Dataset | : Real Time Data |