

# Intelligent System for Student Performance Prediction Using AI

A. Gopi Chandrika<sup>1</sup>, J. Sampath Eswar Sai Kumar<sup>2</sup>, G. Seshagiri<sup>3</sup>, A. Shyam<sup>4</sup>, V. V. Vidya Sagar<sup>5</sup>  
<sup>1,2,3,4</sup>*Department of Computer Science and Engineering (Data Science) Raghu Engineering College  
(Autonomous), JNTU GV, Dakamarri, Visakhapatnam*

<sup>5</sup>*Assistant Professor, Department of Computer Science and Engineering (Data Science)  
Raghu Engineering College (Autonomous), JNTU GV, Dakamarri, Visakhapatnam*

**Abstract**—Student performance prediction is an important area in education, as it helps in identifying students who may need additional support and enables timely academic intervention. In this study, an intelligent system is developed to predict student performance in mathematics using machine learning along with explainable artificial intelligence techniques.

The system is built using a dataset that includes both demographic and academic factors such as gender, race/ethnicity, parental education level, type of lunch, participation in test preparation courses, and scores in reading and writing. A Linear Regression model is used to understand how these factors influence the mathematics score. To improve the model's effectiveness, preprocessing methods like one-hot encoding and feature scaling are applied.

The developed model is deployed through a FastAPI backend and connected to an interactive web interface, allowing users to input data and obtain predictions instantly. To ensure the model is transparent and understandable, SHAP (SHapley Additive exPlanations) is used to interpret the contribution of each feature to the prediction.

The results show that the model performs well when important academic features are included, while its performance decreases significantly when these features are removed. This emphasizes the importance of selecting relevant features in predictive modelling.

Overall, the proposed system not only provides accurate predictions but also offers meaningful insights, making it useful for educators and institutions in monitoring student performance and making informed decisions.

**Index Terms**—Student Performance Prediction, Machine Learning, Linear Regression, Explainable AI, SHAP, Educational Data Mining.

## I. INTRODUCTION

### 1.1 Context and Motivation

In recent years, the use of data-driven techniques in the education sector has gained significant attention. Educational institutions generate large volumes of student-related data, which can be effectively utilized to analyse learning patterns and predict academic outcomes. Predicting student performance at an early stage can help educators identify students who may require additional support and guidance.

Traditional evaluation methods primarily rely on periodic examinations and manual assessment, which often fail to provide timely insights into student progress. With the advancement of machine learning, it has become possible to build predictive systems that can analyse multiple factors influencing academic performance. These systems can assist in improving teaching strategies and enhancing overall student success rates.

The motivation behind this work is to develop an intelligent and automated system that not only predicts student performance accurately but also provides meaningful insights into the contributing factors.

### 1.2 Problem Statement

Despite the availability of educational data, accurately predicting student performance remains a challenging task. Many existing approaches rely on basic statistical techniques or traditional machine learning models that do not fully utilize the relationships between different features.

Additionally, most prediction systems operate as black-box models, providing results without

explaining the reasoning behind them. This lack of transparency reduces trust and limits their practical application in real-world educational environments.

Therefore, there is a need for a system that can:

- Accurately predict student performance using relevant features
- Handle both categorical and numerical data effectively
- Provide interpretable results to understand the influence of each factor

### 1.3 Objectives

The main objectives of this study are:

- To develop a machine learning-based model for predicting student performance
- To analyse and compare multiple models such as Linear Regression, Decision Tree, Random Forest, XGBoost, and ANN
- To implement a real-time prediction system using a web-based interface
- To incorporate explainable AI techniques (SHAP) for understanding feature contributions
- To evaluate the impact of key features on model performance

## II. LITERATURE REVIEW

The field of student performance prediction has gained significant attention with the advancement of machine learning and explainable artificial intelligence techniques.

- Khan et al. (2021) proposed an artificial intelligence-based system to monitor student performance and recommend preventive measures. Their work focused on the early identification of academically at-risk students. However, the approach lacked interpretability, making it difficult to understand how different features influenced predictions.
- Luo et al. (2024) developed a machine learning framework integrated with SHAP for predicting and analysing student performance. Their study emphasized the importance of explainable AI by providing feature-level insights. Despite improved transparency, the system did not focus on real-time deployment or user interaction.
- Villegas et al. (2025) investigated the use of multiple machine learning models combined with

SHAP analysis for academic performance prediction. Their results demonstrated that ensemble models achieve high accuracy while maintaining interpretability. However, the work mainly focused on model evaluation rather than practical system implementation.

- Ahmed et al. (2024) conducted a comparative study of various machine learning algorithms for predicting student outcomes. The research showed that models such as Random Forest and Linear Regression perform effectively on educational datasets. However, the study did not incorporate explainable AI techniques to justify model predictions.
- Kesgin et al. (2025) focused on fairness and bias in student performance prediction models. Their work ensured that machine learning systems produce unbiased results across different groups. While addressing ethical concerns, the study lacked an interactive and deployable system for real-time usage.

From the above studies, it is evident that while existing approaches focus on prediction accuracy and fairness, there is a need for a system that combines accurate prediction, explainability, and real-time deployment. This work addresses these gaps by integrating machine learning with SHAP-based explainability and a web-based application.

## III. METHODOLOGY

### 3.1 Dataset Description

The dataset used in this study contains student-related demographic and academic attributes that influence performance in mathematics. It includes both categorical and numerical features, allowing the model to capture various factors affecting student outcomes. The input features considered are gender, race/ethnicity group, parental level of education, lunch type, test preparation course status, reading score, and writing score. These features represent both academic ability and socio-economic background of students. The target variable is the mathematics score, which is treated as a continuous numerical value. The objective is to predict this score based on the given input features.

### 3.2 Data Preprocessing

To prepare the dataset for training and prediction, several preprocessing steps were performed.

Categorical features such as gender, race/ethnicity, parental education level, and lunch type were converted into numerical form using One-Hot Encoding. Numerical features like reading and writing scores were standardized using StandardScaler to bring all values to a similar range.

The test preparation course feature was converted into a binary format to indicate whether the course was completed or not.

Finally, the dataset was split into training and testing sets. In addition, k-fold cross-validation was used to evaluate the model and ensure reliable performance.

### 3.3 Models Used

Multiple machine learning models were implemented and evaluated to identify the most suitable approach for predicting student performance.

Linear Regression was used as the baseline model due to its simplicity and ability to model linear relationships. Decision Tree was applied to capture non-linear patterns in the data. Random Forest, an ensemble method, was used to improve prediction accuracy and reduce overfitting by combining multiple decision trees. XGBoost was also implemented as an advanced boosting technique to enhance model performance.

In addition, an Artificial Neural Network (ANN) was used to model complex non-linear relationships. The ANN learns feature representations through multiple layers, enabling it to capture deeper patterns in the dataset.

Among all models, Linear Regression achieved the best performance, while ANN produced comparable results with only a marginal difference. This indicates that the dataset primarily follows linear relationships, making simpler models more effective. Therefore, Linear Regression was selected as the final model due to its accuracy, interpretability, and computational efficiency.

### 3.4 Explainable AI

To improve transparency and interpretability, SHAP (SHapley Additive exPlanations) was integrated into the system.

SHAP provides a quantitative measure of how each feature contributes to the final prediction. It helps in

identifying the most influential features such as reading score and writing score, which have a significant impact on the predicted mathematics score. By providing feature-level explanations, SHAP makes the model more understandable and reliable, which is especially important in educational applications where decision-making should be transparent.

### 3.5 System Architecture

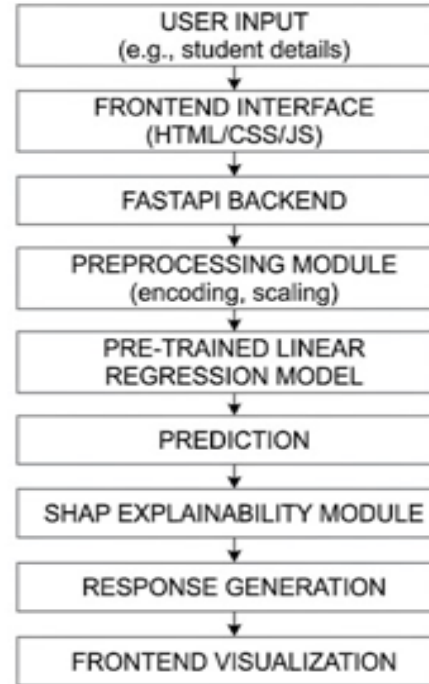


Figure 1 Single Column System Architecture Flow

## IV. SYSTEM DESIGN AND IMPLEMENTATION

### 4.1 System Modules

The proposed system is divided into several functional modules, each responsible for a specific task in the prediction pipeline.

- **Data Processing Module:** Performs preprocessing using encoding and scaling to prepare data for the model.
- **Model Training Module:** Trains multiple models (Linear Regression, Decision Tree, Random Forest, XGBoost) and selects the best-performing model.
- **API Module (FastAPI):** Handles communication between frontend and backend, processes input, and returns predictions.

- Frontend Interface Module: Provides a user interface for entering inputs and displaying results.
- Explainability Module: Uses SHAP to show feature contributions for each prediction.

#### 4.2 Web Application

The system is implemented as a web-based application for real-time prediction.

Users input student details through the interface, and the system generates the predicted math score along with pass/fail status and performance level. Feature contributions are displayed visually to improve interpretability.

The application ensures smooth interaction between frontend and backend, providing a simple and responsive user experience.

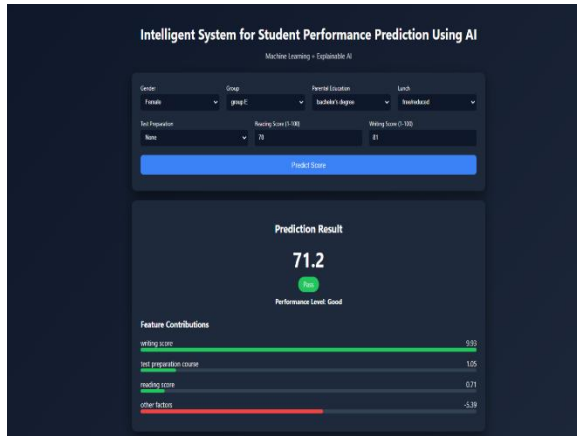


Figure 2: Web-based interface showing input, prediction results, and feature contribution visualization

### V. RESULTS AND DISCUSSION

#### 5.1 Model Comparison

The performance of different models was evaluated using cross-validation.

Model	R <sup>2</sup>	MAE	MSE	RMSE
Linear Regression	0.8718	4.301	29.14	5.395
Decision Tree	0.6953	6.644	69.65	8.340
Random Forest	0.8355	4.880	37.52	6.120

Model	R <sup>2</sup>	MAE	MSE	RMSE
XGBOOST	0.8097	5.271	43.23	6.570
ANN	0.8661	4.387	30.64	5.531

Table 1: Model Performance Comparison

Linear Regression achieved the best performance, indicating strong linear relationships in the dataset.

#### 5.2 Feature Analysis

Model performance was significantly higher when reading and writing scores were included. Removing these features resulted in a sharp drop in accuracy, showing their importance in predicting math scores.

#### 5.3 Explainability Analysis

SHAP analysis revealed that reading and writing scores have the highest impact on predictions, followed by test preparation and other factors. This improves model transparency and helps interpret predictions effectively.

### VI. CONCLUSION

This study presents an intelligent system for predicting student performance using machine learning techniques. Multiple models were evaluated, and Linear Regression achieved the best performance, while ANN produced comparable results with only a marginal difference. This indicates that the dataset primarily follows linear relationships, making simpler models more effective.

The integration of SHAP improved interpretability by providing clear insights into feature contributions, especially highlighting the importance of reading and writing scores. The system also supports real-time prediction through a web-based interface, making it practical for educational use.

### VII. FUTURE WORK

The proposed system can be further enhanced by incorporating additional features such as attendance, behavioral data, and socio-economic factors. Future work may also include applying deep learning models on larger datasets to capture more complex patterns. Additionally, the system can be extended to mobile platforms for better accessibility and integrated with

institutional systems for real-time monitoring. Further improvements can focus on enhancing explainability and ensuring fairness in predictions.

#### REFERENCES

- [1] I. Khan, A. R. Ahmad, N. Jabeur, and M. N. Mahdi, "An artificial intelligence approach to monitor student performance and devise preventive measures," *Smart Learning Environments*, vol. 8, no. 17, 2021.
- [2] E. Ahmed et al., "Student performance prediction using machine learning techniques," *Computational Intelligence and Neuroscience*, 2024.
- [3] Z. Luo et al., "A method for prediction and analysis of student performance using machine learning and SHAP," *Mathematics*, MDPI, 2024.
- [4] W. Villegas et al., "Machine learning models for academic performance prediction using SHAP analysis," *Frontiers in Education*, 2025.
- [5] K. Kesgin et al., "Explaining and ensuring fairness in student academic performance prediction with machine learning," *Applied Sciences*, 2025.
- [6] S. M. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," *NeurIPS*, 2017.
- [7] M. T. Ribeiro, S. Singh, and C. Guestrin, "Why should I trust you? Explaining the predictions of any classifier," *KDD*, 2016.
- [8] C. M. Bishop, "Pattern Recognition and Machine Learning," Springer, 2006.
- [9] P. Cortez and A. Silva, "Using data mining to predict secondary school student performance," 2008.