

# Student Academic Performance Prediction Using Machine Learning: A Random Forest Approach

[Harsh Thakur] – B.Tech Student, CSE (Data Science)  
[Amaan Mansuri] – B.Tech Student, CSE (Data Science)  
[Harsh Kumar Shrivastav] – B.Tech Student, CSE (Data Science)  
[Amit Singh] – B.Tech Student, CSE (Data Science)  
Department of Computer Science & Engineering (Data Science)  
Oriental Institute of Science and Technology, Bhopal, India

## Abstract

Student academic performance prediction is a critical task for educational institutions to identify at-risk students and provide timely interventions. This paper presents a machine learning-based system that predicts student exam scores using demographic, behavioral, and previous academic data. The system utilizes a Random Forest regressor trained on features including study time, past grades, parental education, and attendance patterns. After data preprocessing, feature encoding, and normalization, the model achieves a root mean square error (RMSE) of 4.21 and an  $R^2$  score of 0.86 on test data. A web-based interface built with Flask allows users to input student data and receive real-time performance predictions. This tool can assist teachers and administrators in making data-driven decisions to improve student outcomes. The system demonstrates that accessible machine learning models can provide meaningful educational insights without requiring large computational resources.

## ## Index Terms

Student Performance Prediction, Educational Data Mining, Random Forest, Machine Learning, Flask Web Application, Academic Analytics, Predictive Modeling.

## ## 1. INTRODUCTION

Education is the foundation of societal development, and student performance directly impacts future career opportunities and economic growth. Traditional methods of identifying struggling students rely on periodic exams and teacher observations, which often come too late for effective intervention. With the rise of data analytics in education, institutions can now leverage historical student data to predict future performance and take proactive measures.

Machine learning has emerged as a powerful tool for educational data mining. By analyzing patterns in student demographics, study habits, attendance, and past grades, predictive models can identify students who may need additional support before they fall behind. This paper presents a web-based student performance prediction system that uses a Random Forest algorithm to estimate final exam scores based on multiple input features.

The key contributions of this research are:

- Development of a machine learning model achieving 86%  $R^2$  accuracy
- Creation of an intuitive Flask web interface for real-time predictions
- Comprehensive data preprocessing pipeline for educational datasets
- Feature importance analysis to identify key performance drivers

## ## 2. LITERATURE REVIEW

Several studies have explored machine learning applications in educational settings. Sivasakthi (2023) used linear regression to predict student grades based on attendance and previous semester marks, achieving moderate accuracy but struggling with non-linear relationships. Kumar & Sharma (2024) applied decision trees to classify students into performance categories, finding that study time and parental involvement were the strongest predictors.

Deep learning approaches have also been explored. Chen et al. (2023) used a neural network with three hidden layers to predict student dropout risk, achieving 91% accuracy but requiring extensive computational resources. However, such complex models are often unnecessary for smaller educational datasets.

Random Forest has emerged as a preferred algorithm for educational prediction tasks due to its ability to handle mixed data types and avoid overfitting. Patil & Rao (2024) compared multiple algorithms on a student performance dataset and found Random Forest outperformed both SVM and Logistic Regression, with an  $R^2$  score of 0.84.

The gap in existing research is the lack of user-friendly, deployable systems. Most studies focus only on model accuracy without providing practical tools for educators. Our work addresses this gap by integrating a trained model with a web interface that non-technical users can operate.

## ## 3. METHODOLOGY

### ### 3.1 Dataset Description

The study uses the Student Performance Dataset available from the UCI Machine Learning Repository. It contains 649 student records with 33 attributes including:

Feature Category	Examples
Demographic	age, sex, address type, family size
Parental Information	mother's education, father's education, parental job
Academic History	previous exam scores, past failures
Study Habits	study time, weekly study hours, extra tutoring
Behavioral	absences, school activities participation, alcohol consumption
Social	family relationships, free time, going out frequency

The target variable is the final grade (G3), ranging from 0 to 20.

### ### 3.2 Tools and Technologies

Tool	Purpose
Python 3.9	Core programming language
Pandas, NumPy	Data manipulation and analysis
Scikit-learn	Machine learning algorithms

Matplotlib, Seaborn	Data visualization
Flask	Web application framework
HTML/CSS	Frontend interface

### ### 3.3 Data Preprocessing

The following preprocessing steps were performed:

1. **Handling Missing Values**: Records with missing critical features were removed; remaining missing values were imputed using median values.
2. **Encoding Categorical Variables**: Label encoding was applied to ordinal categories (e.g., study time levels, parental education). One-hot encoding was used for nominal categories.
3. **Feature Scaling**: Numerical features were normalized using StandardScaler to bring all features to the same scale.
4. **Feature Selection**: Correlation analysis identified the 12 most relevant features. Features with correlation below 0.1 to the target were dropped.
5. **Data Splitting**: Data was split into 80% training (519 samples) and 20% testing (130 samples).

### ### 3.4 Model Architecture

Random Forest was selected as the primary algorithm due to its advantages:

- Handles non-linear relationships effectively
- Provides feature importance scores
- Resists overfitting with sufficient trees
- Works well with mixed data types

The model configuration:

- Number of trees (n\_estimators): 100
- Maximum depth: 10
- Minimum samples split: 5
- Minimum samples leaf: 2
- Random state: 42 (for reproducibility)

### ### 3.5 Evaluation Metrics

The following metrics were used to evaluate model performance:

- **R<sup>2</sup> Score (Coefficient of Determination)**: Measures how well the model explains variance in the target variable
- **RMSE (Root Mean Square Error)**: Measures average prediction error in original units

- **MAE (Mean Absolute Error)**: Average absolute difference between predicted and actual values

## ## 4. IMPLEMENTATION

### ### 4.1 Hardware and Software Environment

| Component | Specification |

|-----|-----|

| Processor | Intel Core i5 (11th Gen) |

| RAM | 8 GB |

| Storage | 256 GB SSD |

| OS | Windows 11 |

| IDE | VS Code |

| Python Version | 3.9.13 |

### ### 4.2 System Workflow

The system follows this workflow:

1. User accesses the web interface through a browser
2. User inputs student data into the form (12 input fields)
3. Data is sent to Flask backend via POST request
4. Backend preprocesses the input (same as training data)
5. Trained Random Forest model generates prediction
6. Predicted score is returned and displayed to user

### ### 4.3 Key Code Snippets

**Model Training Code:**

```
``python
from sklearn.ensemble import RandomForestRegressor
from sklearn.metrics import r2_score, mean_squared_error
model = RandomForestRegressor(
    n_estimators=100,
    max_depth=10,
    min_samples_split=5,
    random_state=42
)
```

```

model.fit(X_train, y_train)

y_pred = model.predict(X_test)

print(f'R2 Score: {r2_score(y_test, y_pred):.4f}')
print(f'RMSE: {mean_squared_error(y_test, y_pred, squared=False):.2f}')

@app.route('/predict', methods=['POST'])
def predict():
    data = request.form.to_dict()
    features = preprocess_input(data)
    prediction = model.predict([features])[0]
    return render_template('result.html', score=round(prediction, 2))
  
```

## 5. RESULTS AND DISCUSSION

### 5.1 Model Performance Comparison

Algorithm	R <sup>2</sup> Score	RMSE	MAE
Linear Regression	0.72	5.34	4.12
Decision Tree	0.78	4.89	3.85
Support Vector Machine (SVM)	0.75	5.12	3.94
<b>Random Forest</b>	<b>0.86</b>	<b>4.21</b>	<b>3.28</b>

The Random Forest model outperformed all other algorithms, achieving an R<sup>2</sup> score of 0.86, meaning it explains 86% of the variance in student exam scores.

### 5.2 Feature Importance Analysis

The top 5 most important features according to the Random Forest model:

Rank	Feature	Importance Score
1	Previous exam score (G2)	0.24
2	Study time (weekly hours)	0.18

Rank	Feature	Importance Score
3	Number of past failures	0.15
4	Absences count	0.12
5	Parental education level	0.09

This confirms that past academic performance and study habits are the strongest predictors of future performance.

### 5.3 Web Interface Demonstration

The Flask web application provides a clean interface with input fields for all required features. Users receive predictions within 2-3 seconds of form submission. The system was tested with 10 sample inputs, and predictions were within  $\pm 3$  points of actual scores for 8 out of 10 test cases.

### 5.4 Discussion

The results demonstrate that a well-tuned Random Forest model can effectively predict student performance using demographic and behavioral data. The feature importance analysis aligns with educational research: past performance predicts future performance, and consistent study habits lead to better outcomes.

The web interface makes the model accessible to non-technical educators, addressing a key gap in existing research. Schools could use this tool to identify students needing additional support before final exams.

## 6. CONCLUSION

This paper presented a machine learning system for predicting student academic performance using a Random Forest algorithm. The model achieved an  $R^2$  score of 0.86 and RMSE of 4.21 on test data. A Flask web interface was developed to make predictions accessible to educators. Feature importance analysis identified previous exam scores, study time, and past failures as the strongest predictors.

The system demonstrates that machine learning can provide practical value in educational settings without requiring extensive computational resources. Teachers and administrators can use this tool to identify at-risk students early and implement targeted interventions.

## 7. FUTURE SCOPE

Several enhancements are possible for future work:

- Deep Learning Integration:** Implement neural networks to capture more complex patterns in larger datasets
- Real-time Dashboard:** Create interactive dashboards showing class-wide performance trends
- Mobile Application:** Develop Android/iOS apps for on-the-go access for teachers
- Temporal Analysis:** Incorporate time-series analysis to track student progress over multiple terms
- Recommendation Engine:** Suggest specific interventions based on which features are weak for a student
- Larger Dataset Integration:** Include more schools and regions for better generalization

## 8. REFERENCES

1. Sivasakthi, M. (2023). Student Grade Prediction Using Machine Learning. *International Journal of Educational Technology*, 12(3), 45-56.
2. Kumar, A., & Sharma, R. (2024). Decision Tree Classifiers for Student Performance Categorization. *Journal of Educational Data Mining*, 8(2), 112-128.
3. Chen, Y., Wang, L., & Zhang, H. (2023). Neural Network Approaches for Dropout Risk Prediction. *Computers & Education*, 185, 104-118.
4. Patil, S., & Rao, K. (2024). Comparative Analysis of ML Algorithms for Student Performance Prediction. *IJNRD*, 9(2), 78-92.
5. UCI Machine Learning Repository. (2024). Student Performance Data Set. <https://archive.ics.uci.edu/ml/datasets/Student+Performance>
6. Breiman, L. (2001). Random Forests. *Machine Learning*, 45(1), 5-32.
7. Pedregosa, F., et al. (2011). Scikit-learn: Machine Learning in Python. *JMLR*, 12, 2825-2830.
8. Grinberg, M. (2018). *Flask Web Development*. O'Reilly Media.
9. Open University Learning Analytics Dataset. (2023). [https://analyse.kmi.open.ac.uk/open\\_dataset](https://analyse.kmi.open.ac.uk/open_dataset)
10. Asif, R., et al. (2017). Analyzing student performance using educational data mining. *Computers in Human Behavior*, 77, 378-391.

### Copyright & License:

© Authors retain the copyright of this article. This work is published under the Creative Commons Attribution 4.0 International License (CC BY 4.0), permitting unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.