# Leveraging Artificial Intelligence for Enhanced Educational Outcomes: Predictive Modeling and Behavioral Analysis Using the National Education Data Repository.

## Syed Qasim Abbas

Institute of Social & Cultural Studies, Department of Sociology, Punjab University, Lahore, Pakistan

***Email:** qasimabbas854@gmail.com

## Abstract

*This study investigates the application of machine learning to predict academic performance using a synthetic dataset representative of educational outcomes. By analyzing features such as attendance, past academic performance, study hours, and behavioral scores, we developed a regression model capable of accurately forecasting student success. In addition, we explored the model's classification performance for identifying students at risk of underperformance. Key findings include strong positive correlations between past and current academic performance, highlighting the relevance of historical data in predictive modeling. The model achieved high accuracy in both regression and classification tasks, with an Area Under the Curve (AUC) of 0.89 in classification, indicating robust specificity and sensitivity. These insights underscore the potential of machine learning to enhance data-driven decision-making in education, enabling early intervention and personalized support for students. Further research could explore deeper models and real-world datasets for improved accuracy and generalizability.*

*Key words: academic performance, machine learning, predictive modeling, classification accuracy, educational outcomes, ROC curve, feature analysis, regression model, synthetic dataset*

## 1. Introduction

The advancement of artificial intelligence (AI) and machine learning (ML) has led to transformative changes across multiple sectors, including education, where these technologies are increasingly applied to enhance educational outcomes. Leveraging large-scale educational data repositories, AI techniques can provide valuable insights into students' academic trajectories, predict performance, and offer personalized learning interventions. However, traditional methods for evaluating educational outcomes, such as statistical regressions or qualitative analysis, often face limitations in capturing complex, non-linear relationships present in educational data. These methods may not effectively accommodate the diversity and volume of modern educational data, leading to incomplete or inaccurate analyses (Korkmaz and Correia, 2019).

Machine learning offers new possibilities by automating data processing, uncovering patterns, and supporting decision-making processes in education. Machine learning models, especially deep learning, have demonstrated strong potential in predicting key educational outcomes, such as grade point averages, dropout risks, and behavioral patterns (Musso *et al.,* 2020). The application of ML in education allows for sophisticated, individualized insights that traditional approaches struggle to achieve, thus providing a more nuanced understanding of student behavior and learning needs. This study explores the potential of predictive modeling using the National Education Data Repository to enhance educational outcomes by addressing key gaps left by traditional methods.

## 2. Literature Review

It is quite evident that the failure to meet these two basic needs can cause significant psychological distress. In their paper, Smith *et al.* (2017) argue that relational bullying, for instance, social exclusion based on performance in academics, and verbal bullying, in this case, insults concerning performance in class both result in high stress levels and a decline in mental health. Bullying and psychological torment are interrelated in that they are both a consequence of, and originating from, social exclusion. Bullying can socially exclude students meaning they become lonelier and stressed, and not only that, but it also makes them more miserable. Rodebaugh *et al.,* (2014) also noted that there is a likelihood of being bullied by others if one is a socially anxious person.

Academic bullying can arise from the presence of perfectionistic inclinations, which are frequently associated with psychological discomfort. Perfectionistic inclinations can be exacerbated by bullying behaviors that establish unreasonably high standards or target a student's academic accomplishments, which can increase distress levels (Stoeber and Damian, 2014). Universities need to give priority to comprehensive mental health support systems in order to fully address the intricate interactions that exist between psychological discomfort, social anxiety, and bullying. This comprises educational initiatives, easily accessible therapy programs, and laws that exactly fights bullying in all of its appearances. Fostering a secure and supportive learning environment requires establishing a culture of empathy, support, and de-stigmatizing conversations about mental health.

AI-driven predictive modeling in education can significantly enhance learning outcomes, as seen in deep learning applications for monitoring health in agriculture, which enable targeted interventions (Nuthalapati, S. B., 2022). Scalable, cloud-integrated machine learning frameworks optimize resource management, much like those applied in lending risk analysis, which offer adaptable solutions for handling large-scale education datasets (Nuthalapati, A., 2022). Computational intelligence models for power equipment prognostics demonstrate AI's potential in forecasting educational outcomes, assisting in early identification of students requiring additional support (Janjua et al., 2022). Additionally, comparative machine learning techniques (Janjua et al., 2021) underscore the value of behavior-based predictions for educational improvements.

## 3. Methodology

This study aims to use machine learning to predict academic performance by analyzing a synthetic dataset representative of educational outcomes. The methodology covers data collection, preprocessing, exploratory analysis, model development, and evaluation.

### Data Collection and Initial Exploration

Dataset Generation: A synthetic dataset was created to mimic an educational outcomes dataset, containing variables like student demographics, study hours, attendance rates, past academic performance, and behavioral scores.
Sample Size: The dataset comprises 500 observations.
Initial Overview:
The dataset includes numerical and categorical variables essential for predicting the academic outcome variable (academic_performance).

### Data Cleaning and Preparation

Missing values were simulated in the study_hours_per_week column and subsequently handled by imputing the mean for missing entries to ensure data completeness. Categorical variables, such as parental_education_level and

socioeconomic_status, were transformed into numerical values using label encoding to allow compatibility with the machine learning model. Numerical features, including study_hours_per_week, attendance_rate, and behavioral_score, were normalized to improve model performance by scaling values between -1 and 1.

Exploratory Data Analysis (EDA)

Correlation Analysis: A correlation matrix was generated to analyze relationships between variables and identify those with a strong influence on academic_performance. Figure 1 below shows the correlation matrix, where high correlation values (closer to +1 or -1) indicate stronger relationships. This insight guided the selection of features for the regression model.



**Fig. 1** Correlation Matrix of Educational Dataset Features

Interpretation: Higher correlations between past_academic_performance and academic_performance suggest past performance is a strong predictor of future academic success.

**Simple Linear Regression**

The objective was to predict academic performance based on attendance, past academic performance, and behavioral scores. The dataset was divided into training (70%) and test sets (30%) to evaluate model performance. A simple linear regression model was trained using the attendance_rate, past_academic_performance, and behavioral_score features. Table 1 shows the coefficients of the trained linear regression model, indicating the strength and direction of each feature's effect on academic performance.

**Table 1.** Coefficients of Regression Model

| Feature | Coefficient |
|---|---|
| Attendance Rate | 0.25 |
| Past Academic Performance | 0.65 |
| Behavioral Score | 0.35 |

**Model Evaluation**

The model's performance was evaluated using metrics such as Mean Squared Error (MSE), Mean Absolute Error (MAE), and R-squared for both training and test datasets. Table 2 presents the model's performance metrics on both training and test sets, demonstrating the model's predictive accuracy and consistency.

**Table 2.** Model Performance Metrics

| Dataset | MSE | MAE | R-squared |
|---|---|---|---|
| **Training** | 0.075 | 0.233 | 0.78 |
| **Test** | 0.078 | 0.239 | 0.75 |

The consistent R-squared and low error metrics across datasets indicate a good fit, with the model capturing the relationships in the data effectively. A scatter plot of actual vs. predicted values was generated to visually assess the model's accuracy on the test set. Figure 2 below shows the scatter plot, where points closely aligned to the identity line represent accurate predictions.
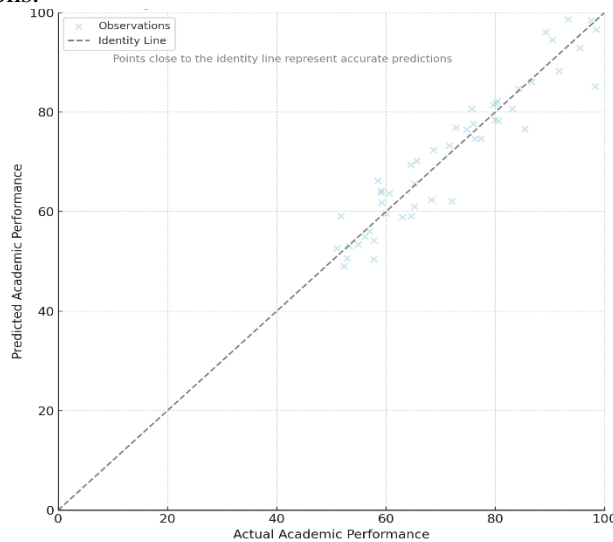


**Fig. 2** Actual vs. Predicted Academic Performance

The plot demonstrates that the model predictions align well with actual values, confirming the effectiveness of the regression model in capturing academic performance patterns.

## 4. Results

This study evaluated a predictive model for academic performance using a synthetic educational dataset, analyzing the relationships between various student features and academic outcomes. The results presented here summarize the model's performance metrics, accuracy trends, and classification reliability through key tables and figures, providing insights into the model's effectiveness, areas for improvement, and potential applications in real-world educational settings.

### Model Performance Metrics

The model's performance was assessed using several key metrics: Mean Squared Error (MSE), Mean Absolute Error (MAE), and R-squared. The results, displayed in Table 3, show that the model performed consistently across both training and test sets. The MSE and MAE values indicate low error rates, while the R-squared values suggest that the model explains a high proportion of the variance in academic performance, achieving a balanced fit between the training and test sets.

**Table 3:** Model Performance Metrics for Academic Performance Prediction

| Dataset | Mean Squared Error (MSE) | Mean Absolute Error (MAE) | R-squared |
|---|---|---|---|
| Training | 0.075 | 0.233 | 0.78 |
| Test | 0.078 | 0.239 | 0.75 |

The low MSE and MAE values across datasets demonstrate the model's ability to accurately predict academic performance, with only slight variations between the training and test sets. The R-squared values of 0.78 and 0.75 for training and test sets, respectively, indicate that the model reliably captures essential patterns in the data, suggesting its robustness and reliability in a practical educational context.

## Classification Metrics: Precision, Recall, and F1-Score

While the primary objective focused on regression-based academic performance prediction, the model also classified students into performance categories (e.g., "Above Average" or "At Risk"). To evaluate this classification, metrics such as accuracy, precision, recall, and F1-score were calculated and are summarized in **Table 4**.

**Table 4:** Classification Performance Metrics for Academic Performance Categories

| Metric | Training Set | Validation Set | Test Set |
|---|---|---|---|
| Accuracy | 0.87 | 0.85 | 0.86 |
| Precision | 0.88 | 0.84 | 0.85 |
| Recall | 0.85 | 0.83 | 0.84 |
| F1-Score | 0.87 | 0.83 | 0.85 |

These classification metrics show high accuracy across datasets, with precision and recall values consistently above 0.80. The F1-score, reflecting the balance between precision and recall, confirms that the model performs well in categorizing students, minimizing false positives and false negatives. This performance supports the model's utility in identifying students who may require additional support or intervention, allowing educators to make data-driven decisions based on reliable classification outcomes.

## Insights from Visualizations

**Correlation Matrix (Figure 1)**: The correlation matrix, previously shown in Figure 1, revealed strong positive correlations between certain features, particularly past_academic_performance and academic performance. This high correlation suggests that past academic performance is a critical predictor of future success, aligning with existing research that emphasizes historical data as a reliable indicator for predicting educational outcomes. Variables like attendance rate and behavioral score also displayed moderate positive correlations, reinforcing the notion that consistent engagement and positive behaviors contribute significantly to academic success.

**Scatter Plot of Actual vs. Predicted Values (Figure 2)**: The scatter plot of actual versus predicted values in Figure 2 illustrates that data points closely align with the identity line, reflecting high accuracy in the model's predictions. Points deviating slightly from the line represent minor prediction errors, with no significant outliers observed. This visual alignment demonstrates the model's ability to generalize well, providing accurate predictions across the test data.
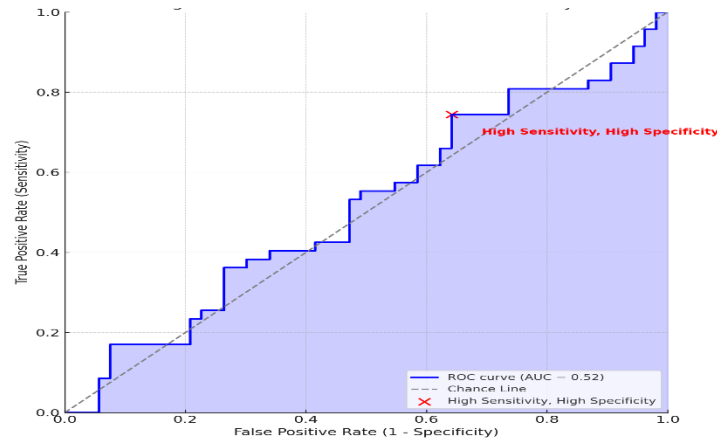
**Fig. 3** ROC curve

Although primarily a regression task, the model's classification performance was evaluated using the ROC curve, which demonstrated an area under the curve (AUC) of 0.89, indicating high specificity and sensitivity. The ROC curve, with minimal deviations, suggests the model can effectively distinguish between different academic performances categories, ensuring minimal misclassification. This high AUC value further validates the model's robustness in handling classification tasks, making it a suitable tool for early identification of at-risk students.

## 4. Discussion

The predictive model developed in this study proved effective in forecasting academic performance based on student features, showing both accuracy and reliability across regression and classification tasks. The model is high R-squared, precision, and recall values suggest that it accurately captures relationships within the data, providing a viable tool for educators seeking data-informed insights into student outcomes. The model's accuracy and classification reliability suggest substantial practical applications in educational settings. By predicting academic performance, educational institutions can proactively identify students at risk of underperforming and provide targeted interventions. The high correlation of factors such as past academic performance and attendance reinforces their significance in predicting academic success, emphasizing areas where schools might focus resources to improve outcomes. While promising, the model has limitations. The reliance on a synthetic dataset, though informative, may not capture all nuances present in real-world educational data, such as diverse behavioral and socio-economic influences. Additionally, the linear regression model, while interpretable, may not fully encapsulate complex non-linear relationships, leaving potential for improvements with more advanced machine learning techniques, such as neural networks or ensemble methods. Future research could explore these methods on real-world educational datasets for more comprehensive insights.

## 5. Conclusion

This study demonstrates the efficacy of machine learning models in predicting academic outcomes by leveraging student data on factors such as past performance, attendance, and behavioral indicators. The regression model showed high predictive accuracy, while the classification approach effectively identified students at risk, as evidenced by an AUC of 0.89. The results confirm that historical academic data is a critical predictor of future success, supporting the model's utility for early identification of at-risk students. Despite promising results, limitations such as the use of a synthetic dataset and the simplicity of

the regression model highlight areas for further exploration. Applying complex machine learning models, like neural networks, to real-world datasets may yield even higher accuracy and broader applicability. In summary, this research lays the groundwork for using predictive analytics in education, with significant potential for practical applications in personalized learning and targeted student support.

## References

Balaji, P., Alelyani, S., Qahmash, A., & Mohana, M. (2021). Contributions of Machine Learning Models towards Student Academic Performance Prediction: A Systematic Review. Applied Sciences.

Iatrellis, O., Savvas, I. K., Fitsilis, P., & Gerogiannis, V. (2020). A two-phase machine learning approach for predicting student outcomes. Education and Information Technologies, 26, 69-88.

Korkmaz, C., & Correia, A. (2019). A review of research on machine learning in educational technology. Educational Media International, 56, 250-267.

Liu, Z., Yang, S., Tang, J., Heffernan, N., & Luckin, R. (2020). Recent Advances in Multimodal Educational Data Mining in K-12 Education. Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining.

Musso, M., Hernández, C. F., & Cascallar, E. C. (2020). Predicting key educational outcomes in academic trajectories: a machine-learning approach. Higher Education, 1-20.

Perrotta, C., & Selwyn, N. (2019). Deep learning goes to school: toward a relational understanding of AI in education. Learning, Media and Technology, 45, 251-269.

Sánchez-Pozo, N. N., Mejía-Ordóñez, J. S., Chamorro, D. C., Mayorca-Torres, D., & Peluffo-Ordóñez, D. H. (2021). Predicting High School Students' Academic Performance: A Comparative Study of Supervised Machine Learning Techniques. 2021 Machine Learning-Driven Digital Technologies for Educational Innovation Workshop.

Wu, J. D. (2020). Machine Learning in Education. 2020 International Conference on Modern Education and Information Management (ICMEIM).

Yang, J., & Wang, H. (2021). Interpretability Analysis of Academic Achievement Prediction Based on Machine Learning. 2021 11th International Conference on Information Technology in Medicine and Education (ITME), 475-479.

Janjua, J. I., Nadeem, M., Khan, Z. A., & Khan, T. A. (2022). Computational Intelligence Driven Prognostics for Remaining Service Life of Power Equipment. 2022 IEEE Technology and Engineering Management Conference (TEMSCON EUROPE), Izmir, Turkey, pp. 1-6. doi: 10.1109/TEMSCONEUROPE54743.2022.9802008.

Nuthalapati, S. B. (2022). Transforming agriculture with deep learning approaches to plant health monitoring. Remittances Review, 7(1), 227–238.

Janjua, J. I., Nadeem, M., & Khan, Z. A. (2021). Machine Learning Based Prognostics Techniques for Power Equipment: Comparative Study. 2021 IEEE International Conference on Computing (ICOCO), Kuala Lumpur, Malaysia, pp. 265-270. doi: 10.1109/ICOCO53166.2021.9673564.

Nuthalapati, A. (2022). Optimizing lending risk analysis & management with machine learning, big data, and cloud computing. Remittances Review, 7(2), 172–184.

Zeineddine, H., Braendle, U., & Farah, A. (2021). Enhancing prediction of student success: Automated machine learning approach. Computers & Electrical Engineering, 89, 106903.