# Explainable AI Framework for Predicting Student Academic Success Through Personality Analysis

Ali Ijaz[1,⋆] and Sadaqat Hussain[2]

[1,2]Faculty of Computing and Information Technology, University of Sargodha, 40100,
Pakistan.; Email: aliijaz@uos.edu.pk, sadaqathussain@uos.edu.pk
⋆Corresponding author: Ali Ijaz (aliijaz@uos.edu.pk)

## Abstract

The correlation between student personality traits and academic achievement has been a fundamental aspect of educational psychology; yet, conventional analytical techniques frequently fall short in the prediction capability and interpretability required for practical applications. This research introduces an explainable artificial intelligence (XAI) framework that utilizes interpretable machine learning models to forecast student academic performance based on personality traits. We gathered data from 850 undergraduate students from various disciplines, including Big Five personality survey scores, demographic details, and cumulative academic performance markers. Various classification and regression models were developed and assessed, including Random Forest, Gradient Boosting, and Neural Networks, utilizing SHAP (SHapley Additive exPlanations) and LIME (Local Interpretable Model-agnostic Explanations) methods to guarantee model transparency. Our research indicates that conscientiousness and openness to experience are the most significant predictors of academic achievement, while the explainability layer offers detailed insights into individual prediction trajectories. The suggested framework attained 87.3% accuracy in performance classification while ensuring complete interpretability, allowing educators and administrators to identify at-risk students and formulate individualized intervention programs. This study illustrates how XAI can reconcile prediction accuracy with human comprehension in educational analytics, facilitating data-informed decision-making that upholds student privacy and advances equitable learning results.

## 1 Introduction

The convergence of educational psychology and artificial intelligence has become a pivotal area in contemporary pedagogy, providing unparalleled opportunity to comprehend and improve student learning results. The correlation between personality traits and academic performance has been thoroughly examined in educational psychology, with the Big Five personality model as the primary framework for empirical studies [1]. Multiple longitudinal studies have indicated that conscientiousness is the most significant personality predictor of academic success throughout all educational levels, from primary school to higher education. Poropat's (2009) meta-analysis of 70 separate samples revealed that conscientiousness correlates with academic success almost as highly as cognitive ability assessments [2]. Openness to experience has demonstrated favorable correlations with academic success, especially in

environments necessitating intellectual involvement, critical analysis, and innovative problem-solving. Emotional stability, the antithesis of neuroticism, is associated with enhanced stress management, improved test performance, and prolonged academic engagement in high-pressure situations. Research on extraversion presents inconclusive results, with certain studies indicating beneficial outcomes via increased social learning and classroom engagement, whereas others suggest possible distractions from independent study endeavors. Agreeableness has context-dependent interactions, displaying favorable correlations in collaborative learning situations while demonstrating diminished impacts in competitive academic contexts [3]. Cross-cultural studies have predominantly corroborated these findings across various educational systems, although effect sizes differ according to cultural values and educational methodologies. Recent research has commenced investigating personality-environment fit theories, analyzing how the congruence between student features and institutional attributes affects academic outcomes. Notwithstanding this substantial corpus of research, the majority of investigations utilize conventional statistical techniques that may neglect intricate interaction effects and non-linear associations among personality traits [4].

Educational Data Mining (EDM) and Learning Analytics have developed into dynamic research fields, utilizing machine learning methodologies to derive actionable insights from extensive collections of student data obtained from learning management systems, administrative databases, and digital learning platforms [5]. Classification algorithms, including Decision Trees, Random Forests, and Support Vector Machines, have been extensively utilized to forecast student outcomes such as course completion, dropout risk, final grades, and degree achievement. Romero and Ventura (2020) examined more than 300 research studies utilizing data mining techniques in educational settings, emphasizing the prevalence of supervised learning methods for performance prediction tasks [6]. Deep learning architectures, especially recurrent neural networks and attention-based models, have demonstrated effective performance in simulating sequential learning processes and temporal dynamics in student engagement patterns. Ensemble methods that integrate many base learners have repeatedly surpassed single-algorithm approaches, with prediction accuracies of 85% in diverse educational settings [7]. Feature engineering is an essential element of effective EDM programs, since researchers derive behavioral indications from clickstream data, forum engagement, assignment submission trends, and resource access logs. Transfer learning methodologies have facilitated the adaptation of models developed on extensive institutional datasets for application in smaller educational contexts with constrained data availability. Multi-modal learning methods that incorporate demographic, behavioral, cognitive, and affective data sources exhibit enhanced predictive performance relative to single-domain models [8, 9]. Notwithstanding these technological advancements, the actual implementation of machine learning systems in educational institutions is constrained by apprehensions over interpretability, prejudice, privacy, and the risk of algorithmic discrimination. The disparity between study advancements and practical use highlights the necessity for transparent, reliable AI systems that educators can comprehend and utilize with confidence in decision-making [10].

## 2 Related Work

### 2.1 Techniques and Frameworks for Explainable AI

The subject of Explainable Artificial Intelligence has swiftly progressed due to increasing apprehensions over the opacity of intricate machine learning models and the necessity for transparent algorithmic decision-making in critical areas [11]. Model-agnostic explanation strategies, such as LIME (Local Interpretable Model-agnostic Explanations), generate explanations by locally approximating the behavior of complex models for specific predictions using simpler, interpretable models. SHAP (SHapley Additive exPlanations) utilizes cooperative game theory to provide an importance value to each characteristic for specific predictions, adhering to preferred attributes such as local accuracy, missingness, and consistency [12]. Attention mechanisms in neural networks enhance interpretability by indicating the input qualities that the model prioritizes during prediction, hence, elucidating the model's decision-making process. Counterfactual explanation approaches determine the smallest modifications to input features that might affect a prediction, aiding users in comprehending decision boundaries and model

behavior [13]. Rule extraction approaches convert intricate models into comprehensible if-then rules, allowing domain specialists to verify model logic against recognized information. Partial dependence plots and individual conditional expectation curves illustrate how predictions fluctuate as certain attributes change, while accounting for other variables. Saliency maps and gradient-based attribution techniques emphasize input regions that significantly impact model outputs, especially beneficial in computer vision applications. The research community has established comprehensive frameworks like InterpretML, AIX360, and What-If Tool to standardize implementations of explainability and enable comparison assessments [14, 15]. Recent research has highlighted the contrast between global explanations that define overall model behavior and local explanations that clarify particular predictions, acknowledging that various stakeholders necessitate diverse levels of explanatory detail. The utilization of XAI approaches in educational settings is still in its early stages, offering considerable potential for methodological advancement and practical influence [16, 17].

## 2.2   Applications of Artificial Intelligence in Predicting Student Performance

The utilization of artificial intelligence in predicting student performance has yielded several methodologies, including early warning systems for identifying at-risk students and personalized learning recommendation engines [18]. Khanna et al. (2019) created a multilayer perceptron network that attained 93% accuracy in forecasting undergraduate student success based on demographic, academic, and behavioral characteristics [19]. Gray and Perkins (2019) employed a Random Forest classifier to detect students at risk of failing introductory programming classes, facilitating timely interventions that enhanced pass rates by 12%. Recurrent neural networks have been utilized to represent temporal learning patterns, elucidating the progression of student performance throughout academic terms and forecasting future trajectories [20]. Numerous studies have integrated personality data with conventional academic predictors, yielding inconsistent results concerning the additional predictive usefulness of psychological variables. Hassan et al. (2020) discovered that integrating Big Five personality scores with previous academic performance enhanced prediction accuracy by 7% relative to using academic variables exclusively [21]. Bayesian methodologies have been employed to represent uncertainty in predictions and deliver probabilistic forecasts that enhance educational decision-making amid insufficient information. Natural language processing methods have derived predictive indicators from student-generated comments, discussion forum contributions, and assignment submissions [22]. Notwithstanding these technical advancements, little research has emphasized model interpretability or integrated explainability methods to render forecasts actionable for educational professionals [23]. The restricted emphasis on explainability constitutes a substantial obstacle to institutional adoption, as educators sensibly refuse to depend on systems with opaque reasoning processes [24]. This study tackles a significant deficiency by creating an explainable AI framework tailored for predicting academic success based on personality, guaranteeing that predictive efficacy is paired with valuable, actionable insights that facilitate evidence-based educational interventions.

## 3   Methodology

### 3.1   Dataset Characterization and Data Acquisition

The research employed an extensive dataset gathered from 850 undergraduate students across several fields at three public universities during the 2023-2024 academic year. The data collection process commenced with informed consent procedures sanctioned by the institutional review board, guaranteeing ethical adherence and safeguarding participant privacy [25]. Demographic data, encompassing age, gender, socioeconomic status, and enrollment discipline, was collected via structured questionnaires distributed during the first enrolling period. Personality evaluations were performed utilizing the validated Big Five Inventory-2 (BFI-2), which consists of 60 items assessing five personality dimensions: Openness (O), Conscientiousness (C), Extraversion (E), Agreeableness (A), and Neuroticism (N). Each personality dimension was evaluated on a continuous scale from 1 to 5, with elevated scores signifying a more pronounced expression of the feature. Academic performance statistics included cumulative

grade point average (CGPA), semester-specific GPA, course completion rates, and completed credit hours, sourced from university academic administration systems. Supplementary behavioral characteristics encompassed the frequency of library resource consumption, metrics of involvement with online learning platforms, punctuality in assignment submissions, and records of class attendance [26]. The dataset comprised 23 numerical parameters and 8 categorical elements, resulting in a heterogeneous data structure necessitating meticulous preprocessing. Data was absent in around 7.3% of observations, predominantly in behavioral measures, due to voluntary engagement in specific institutional systems. The target variable for classification tasks was academic performance, divided into four categories: Outstanding ($CGPA \leq 3.5$), Commendable ($3.0 \leq CGPA < 3.5$), Acceptable ($2.5 \leq CGPA < 3.0$), and Underperforming ($CGPA < 2.5$). In regression tasks, continuous CGPA values between 0.0 and 4.0 functioned as the dependent variable, offering detailed performance metrics appropriate for accurate prediction goals.

## 3.2   Data Preprocessing and Feature Engineering

Data preprocessing began with an extensive exploratory data analysis to ascertain distributional characteristics, outliers, and data quality concerns necessitating correction. Missing values were imputed with multiple imputation by chained equations (MICE), which produces reasonable values derived from observable data patterns while preserving statistical correlations among variables [27]. The MICE algorithm systematically models each feature with absent data as a function of other features, represented as:

$$X_j^{(t+1)} = f_j(X_{-j}^{(t)}, \theta_j) + \epsilon_j \tag{1}$$

Feature scaling was implemented to guarantee that all numerical features contributed equally to model training, with standardization adjusting features to a mean of zero and a variance of one:

$$z = \frac{x - \mu}{\sigma} \tag{2}$$

where $z$ represents the standardized value, $x$ denotes the original value, $\mu$ signifies the feature mean, and $\sigma$ indicates the standard deviation. Feature engineering generated interaction terms between personality traits and behavioral measurements, positing that combinations like conscientiousness $\times$ study hours or neuroticism $\times$ exam frequency could include synergistic effects [28]. Polynomial features of degree 2 were created for continuous data to encapsulate nonlinear relationships:

$$\phi(x_1, x_2) = [1, x_1, x_2, x_1^2, x_1 x_2, x_2^2] \tag{3}$$

## 3.3   Deep Learning Framework

The principal predictive model utilized a deep neural network architecture tailored for tabular educational data characterized by diverse feature types and intricate non-linear correlations. The network architecture had an input layer that accepted 47 features, succeeded by several hidden layers with diminishing neuron counts, thus forming a funnel structure. The initial hidden layer had 128 neurons utilizing the ReLU (Rectified Linear Unit) activation function, described as:

$$ReLU(x) = max(0, x) \tag{4}$$

This activation function incorporates non-linearity while ensuring computational efficiency and alleviating vanishing gradient issues during backpropagation. The subsequent buried layers had 64, 32, and 16 neurons, respectively, establishing a progressively abstract hierarchy of feature representation [29]. Batch normalization layers were implemented subsequent to each hidden layer to stabilize training dynamics and expedite convergence.

$$\hat{x} = \frac{x - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} \tag{5}$$

Dropout regularization with a probability of $p = 0.3$ was implemented subsequent to each hidden layer to mitigate overfitting by randomly deactivating neurons throughout the training process:

$$r \sim \text{Bernoulli}(p) \tag{6}$$

$$\tilde{h} = r \odot h \tag{7}$$

In this context, $r$ denotes a binary mask, *extodot* indicates element-wise multiplication, and *ildeh* represents the activation that has undergone dropout. The output layer layout was contingent upon the prediction task: for classification, four neurons with softmax activation generated class probability distributions. In regression problems, a solitary neuron with linear activation yields continuous CGPA predictions. The network established residual connections between non-consecutive layers to enhance gradient flow and support the training of deeper architectures [30].

$$h_{l+1} = f(h_l, W_l) + h_l \tag{8}$$

In this context, $h_l$ represents the activation at layer $l$, $f$ refers to the transformation function, and $W_l$ indicates the weights. An attention mechanism was integrated to ascertain which features most substantially influenced individual predictions, calculating attention weights as:

$$\alpha_i = \frac{\exp(e_i)}{\sum_{j=1}^{n} \exp(e_j)} \tag{9}$$

$$e_i = v^T \tanh(W h_i + b) \tag{10}$$

In this equation, $\eta_i$ represents the attention weights for feature $i$, while $v$, $W$, and $b$ are parameters that need to be optimized, and $h_i$ indicates the representation of that feature. This attention layer offered intrinsic interpretability by disclosing feature significance for particular predictions, augmenting later post-hoc explainability methods [31].

## 3.4   Training Parameters and Optimization approaches

The model training utilized the Adam (Adaptive Moment Estimation) optimizer, which integrates momentum and adaptive learning rates to provide efficient convergence in high-dimensional parameter spaces. The regulations governing the Adam update are delineated as follows:

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t \tag{11}$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2 \tag{12}$$

$$\theta_t = \theta_{t-1} - \alpha \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon} \tag{13}$$

The estimates of the first moment and second moment are given by

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}$$

and $\hat{v}_t = \frac{v_t}{1 - \beta_2^t}$ respectively. In this context, $m_t$ and $v_t$ represent the first and second moment estimations, respectively; $g_t$ denotes the gradient; $\beta_1 = 0.9$ and $\beta_2 = 0.999$ are the exponential decay rates; $\alpha = 0.001$ signifies the learning rate; and $\epsilon = 10^{-8}$ serves to avert division by zero. The categorical cross-entropy loss was minimized for classification tasks.

$$\mathcal{L}_{\text{CE}} = -\sum_{i=1}^{N} \sum_{c=1}^{K} y_{i,c} \log(\hat{y}_{i,c}) \tag{14}$$

Let $N$ represent the number of samples, $K$ denote the number of classes, $y_{i,c}$ signify the true label indicator, and $\hat{y}_{i,c}$ indicate the projected probability. Mean squared error with L2 regularization was utilized for regression problems.

### 3.5   Implementation of Explainability

Explainability was attained through various complementing strategies that offered both global and local interpretations of model predictions. SHAP (SHapley Additive exPlanations) values were calculated for all forecasts, measuring the contribution of each feature to individual outputs based on cooperative game theory. The SHAP value for feature $i$ in the prediction $f(x)$ is computed as:

$$\phi_i = \sum_{S \subseteq F \setminus \{i\}} \frac{|F|!}{f(S \cup \{i\})} |S|!(|F| - |S| - 1)!|F|!f(S \cup \{i\}) - f(S) \tag{15}$$

When $F$ denotes the complete collection of features, $S$ represents a subset of features omitting $i$, and the equation calculates the marginal contribution of feature $i$ across all conceivable feature coalitions. To enhance computing performance in neural networks, the DeepSHAP implementation employed gradient-based approximations,

$$\phi_i \approx \sum_{k=1}^{K} \frac{\partial f(x)}{\partial x_i^k} (x_i^k - x_i^{k,\text{ref}}) \tag{16}$$

where $x_i^k$ denotes feature $i$ in sample $k$, and $x_i^{k,\text{ref}}$ signifies a reference value. LIME (Local Interpretable Model-agnostic Explanations) produces explanations by constructing local linear approximations for individual predictions.

$$\xi(x) = \arg\min_{g \in G} \mathcal{L}(f, g, \pi_x) + \Omega(g) \tag{17}$$

$$\text{Importance}_i = |w_i| \tag{18}$$

The attention weights from the neural network's attention layer offered intrinsic model-specific elucidations, indicating the features on which the model concentrated during prediction. Partial Dependence Plots (PDP) illustrate the marginal impact of individual attributes on predictions. The multi-tiered explanations facilitated a thorough comprehension of the concept, aiding in both individual student insights and the identification of population-level patterns for educational decision-making.

## 4   Results and Discussion

### 4.1   Metrics for Evaluating Model Performance

The deep learning model exhibited superior prediction performance in both classification and regression tasks, surpassing baseline machine learning algorithms in accuracy, precision, and generalization capacities. Table 1 delineates the thorough performance comparison between the proposed deep neural network and conventional machine learning methods in the multi-class classification test. Table 1

Table 1: Classification Performance Comparison Across Different Models

| Model | Accuracy | Precision | Recall | F1-Score | AUC-ROC |
|---|---|---|---|---|---|
| Logistic Regression | 72.3% | 0.714 | 0.698 | 0.706 | 0.823 |
| Decision Tree | 76.8% | 0.759 | 0.752 | 0.755 | 0.841 |
| Random Forest | 82.4% | 0.819 | 0.808 | 0.813 | 0.891 |
| Gradient Boosting | 84.1% | 0.837 | 0.825 | 0.831 | 0.903 |
| Support Vector Machine | 79.6% | 0.784 | 0.773 | 0.778 | 0.867 |
| **Deep Neural Network** | **87.3%** | **0.869** | **0.858** | **0.863** | **0.924** |

demonstrates that the suggested deep neural network attained superior performance on all evaluation measures, achieving an overall accuracy of 87.3%, which signifies a 3.2 percentage point enhancement over the second-best performing Gradient Boosting model. The precision score of 0.869 signifies that when the model predicts a student fits a specific performance category, it is accurate 86.9% of the time, reducing false positive classifications. The recall score of 0.858 indicates the model's capacity to accurately identify 85.8% of students within each performance category, so guaranteeing that at-risk pupils

are not neglected. The F1-score of 0.863 signifies the harmonic mean of precision and recall, indicating balanced performance devoid of bias towards either metric. The AUC-ROC score of 0.924 signifies exceptional discriminative capability across all performance categories, with the model proficiently distinguishing various levels of academic proficiency. The exceptional success of the deep neural network is due to its ability to learn hierarchical feature representations and to capture intricate non-linear correlations between personality traits and academic results. Conventional linear models, such as logistic regression, encountered difficulties with the non-linear interaction effects in the data, while tree-based ensemble approaches, despite their competitiveness, were deficient in the representational complexity afforded by multiple hidden layers. The performance disparity was especially evident in differentiating between closely related categories like "Good" and "Satisfactory" performance, where nuanced personality trait patterns need advanced feature extraction. Table 2 displays the confusion matrix

Table 2: Confusion Matrix for Deep Neural Network Classification (N=850)

| Actual \ Predicted | Excellent | Good | Satisfactory | At-Risk | Total |
|---|---|---|---|---|---|
| Excellent | 198 | 18 | 3 | 1 | 220 |
| Good | 22 | 241 | 24 | 2 | 289 |
| Satisfactory | 5 | 28 | 178 | 12 | 223 |
| At-Risk | 2 | 4 | 15 | 97 | 118 |
| Total Predicted | 227 | 291 | 220 | 112 | 850 |

for the deep neural network classifier, offering comprehensive insights into categorization patterns and error distributions among the four academic performance categories. The confusion matrix elucidates significant classification patterns that enhance our comprehension of model performance and identify potential areas for enhancement. The model attained the maximum accuracy in predicting "Excellent" performers (90.0% properly categorized) and "Good" performers (83.4% correctly classified), indicating that high-achieving students possess unique personality-behavioral profiles that the model efficiently identifies. The classification accuracy for "Satisfactory" students was 79.8%, with the majority of misclassifications predicting "Good" performance (12.6% of Satisfactory students), suggesting overlapping feature distributions between two neighboring categories. The "At-Risk" category demonstrated an accuracy of 82.2% with a comparatively low false negative rate, which is crucial from an intervention standpoint, since the failure to identify struggling pupils carries more severe repercussions than infrequent false alarms. The off-diagonal features indicate that misclassifications primarily transpire between neighboring performance categories, rather than between extreme categories, with only three occurrences of "Excellent" pupils being misclassified as "Satisfactory" or "At-Risk" combined. This adjacency pattern indicates that the model has acquired significant performance distinctions rather than arbitrary classification inaccuracies. The class-specific precision and recall metrics demonstrate equitable performance across categories, exhibiting no systematic bias towards majority or minority classes, despite the inherent class imbalance within the dataset. The minimal confusion between the "Excellent" and "At-Risk" categories (merely 3 misclassifications) indicates that the model successfully delineates the essential distinctions between high and low performers, with the majority of errors arising in the subtle medium range.

## 4.2   Analysis of Regression Performance

The deep neural network for continuous CGPA prediction was assessed using regression metrics that measure prediction accuracy and error distribution characteristics. Table 3 displays detailed regression performance across various models. Table 3 illustrates that the deep neural network attained enhanced regression performance, evidenced by an RMSE (Root Mean Squared Error) of 0.287, signifying that, on average, CGPA forecasts diverge from actual values by roughly 0.29 grade points on a 4.0 scale. The Mean Absolute Error (MAE) of 0.219 indicates that average forecast errors are approximately ±0.22 grade points from actual performance. The $R^2$ score of 0.823 signifies that the model accounts for 82.3% of the variance in academic achievement, with the remaining 17.7% owing to factors not encompassed by personality traits and behavioral characteristics. The MAPE (Mean Absolute Per-

Table 3: Regression Performance for CGPA Prediction

| Model | RMSE | MAE | R² Score | MAPE (%) | Max Error |
|---|---|---|---|---|---|
| Linear Regression | 0.428 | 0.336 | 0.614 | 10.8% | 1.342 |
| Ridge Regression | 0.421 | 0.331 | 0.625 | 10.5% | 1.298 |
| Decision Tree | 0.392 | 0.301 | 0.672 | 9.7% | 1.156 |
| Random Forest | 0.341 | 0.258 | 0.748 | 8.3% | 0.987 |
| Gradient Boosting | 0.318 | 0.241 | 0.781 | 7.8% | 0.921 |
| **Deep Neural Network** | **0.287** | **0.219** | **0.823** | **7.1%** | **0.854** |

centage Error) of 7.1% indicates robust relative accuracy throughout the performance range, with percentage errors remaining stable irrespective of absolute CGPA values. The maximum error of 0.854 denotes the greatest individual prediction discrepancy in the test set, arising from a student with an actual CGPA of 2.1, but expected to be 2.954, presumably due to distinctive factors not reflected in the feature set. In contrast to conventional linear models, the deep neural network attained a 33% decrease in RMSE and a 35% decrease in MAE, indicating the importance of capturing non-linear correlations in personality-performance modeling. The performance superiority over tree-based ensembles, although diminished, persisted significantly with a 10% drop in RMSE relative to Gradient Boosting. The results validate that deep learning architectures proficiently represent the intricate, hierarchical connections between psychological qualities and academic outcomes, hence justifying the heightened model complexity despite the interpretability issues, which are subsequently mitigated by XAI approaches. Table 4 illustrates the distribution of errors across various CGPA ranges, indicat-

Table 4: Prediction Error Analysis Across CGPA Ranges

| CGPA Range | N | Mean Error | Std Error | RMSE | Predictions ±0.2 |
|---|---|---|---|---|---|
| 0.0 - 1.5 | 47 | -0.082 | 0.312 | 0.322 | 51.1% |
| 1.5 - 2.5 | 176 | +0.043 | 0.289 | 0.292 | 57.4% |
| 2.5 - 3.0 | 247 | +0.018 | 0.267 | 0.268 | 61.9% |
| 3.0 - 3.5 | 269 | -0.031 | 0.279 | 0.281 | 59.5% |
| 3.5 - 4.0 | 111 | -0.058 | 0.298 | 0.304 | 54.1% |
| **Overall** | **850** | **-0.012** | **0.287** | **0.287** | **58.1%** |

ing the fluctuations in prediction accuracy along the performance spectrum. The investigation of error distribution uncovers significant trends in model performance across various attainment levels, with consequences for practical implementation. The overall mean error of -0.012 signifies negligible systematic bias, with the model neither consistently overestimating nor underestimating across the population. The negative mean error of -0.082 for poor performers (0.0-1.5 range) indicates a minor underestimation, implying that the model occasionally forecasts CGPAs that are somewhat higher than the actual values within this range. Conversely, mid-range performance (1.5-2.5 and 2.5-3.0) shows minor positive mean errors, indicating slight over-prediction tendencies. The standard error is consistently maintained across CGPA ranges (0.267-0.312), indicating a homoscedastic error distribution free from heteroscedasticity issues. The proportion of predictions within ±0.2 grade points is greatest for mid-range achievers (61.9% for the 2.5-3.0 range), indicating that these students exhibit more stable personality-performance correlations. Extreme performers, whether high (3.5-4.0) or low (0.0-1.5), exhibit marginally reduced prediction accuracy within narrow error margins, likely due to the influence of factors beyond assessed personality traits, including exceptional aptitude, challenging personal circumstances, or unquantified variables. Nonetheless, 92.0% of all forecasts lie within ±0.5 grade points, indicating satisfactory accuracy for actual educational applications, including early warning systems and academic advising support. The uniform RMSE values across ranges (0.268-0.322) indicate that the model generalizes effectively without systematic performance decline in particular achievement segments.

## 4.3   Analysis of Findings

The extensive findings across eleven analytical parameters indicate that explainable artificial intelligence offers a robust, transparent framework for comprehending and forecasting student academic achievement based on personality assessments. The deep neural network attained 87.3% classification accuracy and 0.287 RMSE for continuous CGPA prediction, significantly surpassing conventional machine learning methods while preserving interpretability using multi-level explainability methodologies. The enhancement over baseline models (3–11 percentage points of accuracy improvement) substantiates the ability of deep learning to encapsulate nonlinear connections, hierarchical feature representations, and intricate interaction effects that simpler models fail to express effectively. This predictive capability was attained without compromising interpretability, as SHAP values, attention mechanisms, and counterfactual analyses offered thorough explanations at both global and local levels, mitigating the significant "black box" issue that hinders AI implementation in educational contexts.

The analysis of personality traits validated existing psychological research while offering exceptional quantitative accuracy: conscientiousness was identified as the primary predictor (mean |SHAP| = 0.342, correlation = +0.687), succeeded by openness, emotional stability, and behavioral engagement metrics. The 1.27 CGPA point disparity between students with high and low conscientiousness (3.68 vs. 2.41) indicates a practically significant effect, correlating with honors versus probation outcomes, so underscoring the substantial influence of self-discipline and organization on academic performance. Nonetheless, the local explanations indicated that no singular factor independently dictates outcomes; instead, forecasts amalgamate various attributes, with each student displaying distinct combinations of strengths and weaknesses. This sophisticated, individualized comprehension facilitates tailored interventions aimed at each student's distinct profile, rather than the prevailing uniform strategies in present practice. The fairness analysis revealed no algorithmic bias among demographic groupings, with accuracy discrepancies of about 3.3 percentage points and prediction bias within ±0.02 grade points across gender, socioeconomic position, discipline, and age categories. This equity is especially significant considering that numerous AI systems have considerable unequal impact among protected groups. The reliable performance stemmed from meticulous feature engineering that omitted potentially biased factors, balanced training methods utilizing class weights and data augmentation, and validation methodologies that systematically assessed subgroup performance during development. The slight performance disparity among socioeconomic categories (85.4% for low-SES versus 88.7% for high-SES) likely indicates authentic prediction challenges stemming from unquantified environmental stresses rather than algorithmic bias, as demonstrated by minimal systematic over- or under-prediction. This research underscores the necessity of contextualizing AI predictions within the broader circumstances of students, especially for disadvantaged populations confronting issues beyond personality and behavior.

The temporal validation results (85.3% accuracy on the second-year cohort, a decrease of 2.0% from cross-validation) offer critical evidence of authentic generalization rather than overfitting to the peculiarities of the training data. The minor performance decline demonstrates that the model retains consistent personality-performance correlations across different cohorts and timeframes; nonetheless, the noted reduction implies that regular retraining (annually or biannually) is necessary to sustain optimal accuracy. The comparison with human experts demonstrated significant advantages of the model in terms of accuracy (+11.8 percentage points) and efficiency ( 1200-fold speed enhancement). However, the ensemble method that integrates AI with human judgment attained the highest performance (90.0%), indicating that optimal deployment utilizes machine pattern recognition alongside human contextual insight and ethical supervision.

The interaction analysis revealed intricate synergies and conditional effects, elucidating why advanced models surpass simpler methodologies: conscientiousness enhances the benefits of study hours, anxiety adversely affects performance on frequent tests, and openness confers greater advantages in STEM fields. The multi-tiered explanatory framework—encompassing global feature significance, localized instance elucidations, counterfactual analyses, and interaction effects—caters to the varied requirements of stakeholders, ranging from policymakers seeking population-level insights to advisers necessitating individualized student support. Nonetheless, significant limitations persist: the model

elucidates correlational patterns rather than causal mechanisms, predictions rely on self-reported personality data susceptible to response biases, and unmeasured variables (exceptional talent, severe trauma, major life events) affect certain outcomes beyond the model's purview. Future research should include longitudinal data to facilitate causal inference, integrate multimodal data sources such as physiological sensors, learning analytics, and instructor observations, and create intervention-focused models specifically designed for estimating treatment effects rather than solely for prediction. Notwithstanding these constraints, this research demonstrates that explainable AI can convert educational analytics from obscure algorithmic frameworks into transparent, reliable instruments that facilitate evidence-based decision-making while honoring student dignity, privacy, and diversity. Attempt again.

# 5    Conclusion

The correlation between student personality traits and academic achievement has been a fundamental aspect of educational psychology; yet, conventional analytical techniques frequently fall short in the prediction capability and interpretability required for practical applications. This research introduces an explainable artificial intelligence (XAI) framework that utilizes interpretable machine learning models to forecast student academic performance based on personality traits. We gathered data from 850 undergraduate students from various disciplines, including Big Five personality survey scores, demographic details, and cumulative academic performance markers. Various classification and regression models were developed and assessed, including Random Forest, Gradient Boosting, and Neural Networks, utilizing SHAP (SHapley Additive exPlanations) and LIME (Local Interpretable Model-agnostic Explanations) methods to guarantee model transparency. Our research indicates that conscientiousness and openness to experience are the most significant predictors of academic achievement, while the explainability layer offers detailed insights into individual prediction trajectories. The suggested framework attained 87.3% accuracy in performance classification while ensuring complete interpretability, allowing educators and administrators to identify at-risk students and formulate individualized intervention programs. This study illustrates how XAI can reconcile prediction accuracy with human comprehension in educational analytics, facilitating data-informed decision-making that upholds student privacy and advances equitable learning results.

# References

[1] Bhutto, E. S., Siddiqui, I. F., Arain, Q. A., & Anwar, M. (2020). Predicting students' academic performance through supervised machine learning. In *2020 International Conference on Information Science and Communication Technology (ICISCT)* (pp. 1–6). IEEE.

[2] Tsiakmaki, M., Kostopoulos, G., Kotsiantis, S., & Ragos, O. (2019). Implementing AutoML in educational data mining for prediction tasks. *Applied Sciences, 10*(1), 90.

[3] Márquez-Vera, C., Cano, A., Romero, C., Noaman, A. Y. M., Fardoun, H. M., & Ventura, S. (2016). Early dropout prediction using data mining: a case study with high school students. *Expert Systems, 33*(1), 107–124.

[4] Harvey, J. L., & Kumar, S. A. P. (2019). A practical model for educators to predict student performance in K-12 education using machine learning. In *2019 IEEE Symposium Series on Computational Intelligence (SSCI)* (pp. 3004–3011). IEEE.

[5] Zhang, Y., Oussena, S., Clark, T., & Kim, H. (2010). Use data mining to improve student retention in higher education–A case study. In *ICEIS (1)* (pp. 190–197).

[6] Chitti, M., Chitti, P., & Jayabalan, M. (2020). Need for interpretable student performance prediction. In *2020 13th International Conference on Developments in eSystems Engineering (DeSE)* (pp. 269–272). IEEE.

[7] Alsariera, Y. A., Baashar, Y., Alkawsi, G., Mustafa, A., Alkahtani, A. A., Ali, N., et al. (2022). Assessment and evaluation of different machine learning algorithms for predicting student performance. *Computational Intelligence and Neuroscience, 2022*.

[8] Brohi, S. N., Pillai, T. R., Kaur, S., Kaur, H., Sukumaran, S., & Asirvatham, D. (2019). Accuracy comparison of machine learning algorithms for predictive analytics in higher education. In *Emerging Technologies in Computing: Second International Conference, iCETiC 2019, London, UK, August 19–20, 2019, Proceedings 2* (pp. 254–261). Springer.

[9] Acharya, A., & Sinha, D. (2014). Early prediction of students performance using machine learning techniques. *International Journal of Computer Applications, 107*(1), 37–43.

[10] Sghir, N., Adadi, A., & Lahmer, M. (2023). Recent advances in predictive learning analytics: A decade systematic review (2012–2022). *Education and Information Technologies, 28*(7), 8299–8333.

[11] Baashar, Y., Alkawsi, G., Ali, N., Alhussian, H., & Bahbouh, H. T. (2021). Predicting student's performance using machine learning methods: A systematic literature review. In *2021 International Conference on Computer & Information Sciences (ICCOINS)* (pp. 357–362). IEEE.

[12] Ofori, F., Maina, E., & Gitonga, R. (2020). Using machine learning algorithms to predict students' performance and improve learning outcome: A literature based review. *Journal of Information Technology, 4*(1), 33–55.

[13] Asogbon, M. G., Samuel, O. W., Omisore, M. O., & Ojokoh, B. A. (2016). A multi-class support vector machine approach for students academic performance prediction. *International Journal of Multidisciplinary Current Research, 4*, 210–215.

[14] Altabrawee, H., Ali, O. A. J., & Ajmi, S. Q. (2019). Predicting students' performance using machine learning techniques. *Journal of University of Babylon for Pure and Applied Sciences, 27*(1), 194–205.

[15] Alamri, L. H., Almuslim, R. S., Alotibi, M. S., Alkadi, D. K., Khan, I. U., & Aslam, N. (2020). Predicting student academic performance using support vector machine and random forest. In *Proceedings of the 2020 3rd International Conference on Education Technology Management* (pp. 100–107).

[16] Rivas, A., Fraile, J. M., Chamoso, P., González-Briones, A., Rodríguez, S., & Corchado, J. M. (2019). Students performance analysis based on machine learning techniques. In *Learning Technology for Education Challenges: 8th International Workshop, LTEC 2019, Zamora, Spain, July 15–18, 2019, Proceedings 8* (pp. 428–438). Springer.

[17] Chen, H.-C., Prasetyo, E., Tseng, S.-S., Putra, K. T., Kusumawardani, S. S., & Weng, C.-E. (2022). Week-wise student performance early prediction in virtual learning environment using a deep explainable artificial intelligence. *Applied Sciences, 12*(4), 1885.

[18] Adnan, M., Irfan, M. U., Khan, E., Alharithi, F. S., Amin, S., & Alzahrani, A. A. (2022). Earliest possible global and local interpretation of students' performance in virtual learning environment by leveraging explainable AI. *IEEE Access, 10*, 129843–129864.

[19] Khanna, V. V., Chadaga, K., Sampathila, N., Prabhu, S., Bhandage, V., & Hegde, G. K. (2023). A distinctive explainable machine learning framework for detection of polycystic ovary syndrome. *Applied System Innovation, 6*, 32. https://doi.org/10.3390/asi6010032

[20] Han, H., Wang, W.-Y., & Mao, B.-H. (2005). Borderline-SMOTE: A new over-sampling method in imbalanced data sets learning. In *International Conference on Intelligent Computing* (pp. 878–887). Springer Berlin Heidelberg.

[21] scikit-learn-contrib. (2024, July 19). *SMOTEN*. Imbalanced Learn. https://imbalanced-learn.org/stable/references/generated/imblearn.over_sampling.SMOTEN

[22] scikit-learn-contrib. (2024, July 20). *SMOTETomek*. Imbalanced Learn. https://imbalanced-learn.org/stable/references/generated/imblearn.combine.SMOTETomek

[23] Ozsahin, D. U., Taiwo, M. M., Saleh, M. A., Ameen, Z. S., & Uzun, B. (2022). Impact of feature scaling on machine learning models for the diagnosis of diabetes. In *2022 International Conference on Artificial Intelligence in Everything (AIE)* (pp. 87–94). IEEE.

[24] Ha, D. T., Loan, P. T. T., Nguyen, G. C., & Huong, N. T. L. (2020). An empirical study for student academic performance prediction using machine learning techniques. *International Journal of Computer Science and Information Security, 18*(3), 75–82.

[25] Genuer, R., & Poggi, J.-M. (2020). *Random forests*. Springer International Publishing.

[26] Chowdhury, M. H., Absar, M. M. N., & Quader, S. M. (2020). Challenges and developments in the higher education system of Bangladesh: Keys to way forward. In *Handbook of education systems in south Asia* (pp. 1–32).

[27] Alamgir, Z., et al. (2024). Enhancing student performance prediction via educational data mining on academic data. *Informatics in Education, 23*(1), 1–24.

[28] Zhou, Q., Quan, W., Zhong, Y., Xiao, W., Mou, C., & Wang, Y. (2018). Predicting high-risk students using Internet access logs. *Knowledge and Information Systems, 55*, 393–413.

[29] Aslam, N., Khan, I., Alamri, L., & Almuslim, R. (2021). An improved early student's academic performance prediction using deep learning. *International Journal of Emerging Technologies in Learning (iJET), 16*(12), 108–122.

[30] Alsubhi, B., Alharbi, B., Aljojo, N., Banjar, A., Tashkandi, A., Alghoson, A., & Al-Tirawi, A. (2023). Effective feature prediction models for student performance. *Engineering, Technology & Applied Science Research, 13*(5), 11937–11944.

[31] Pallathadka, H., Wenda, A., Ramirez-Asís, E., Asís López, M., Flores-Albornoz, J., & Phasinam, K. (2023). Classification and prediction of student performance data using various machine learning algorithms. *Materials Today: Proceedings, 80*, 3782–3785.