# Predicting student performance outcomes using Machine Learning: A comprehensive comparison between Linear Regression, Random Forests, XGBoost and Support Vector Regression (SVR)

**Project Summary**

The goal of this project is designing, analysing, and comparing various machine learning models that can be used to forecast the performance of a student using a dataset that includes educational, behavioural, and other characteristics of the student. Generally, the aim of the project was to see which model of forecasting can be used to obtain the most accurate and optimal results and thus assist institutions of learning in taking preventive measures against any student who performs below expectations academically early in their time at the institution. After analysing Linear Regression, Support Vector Regression (SVR), Random Forest Regression, and Extreme Gradient Boosting (XGBoost), the model that performed the best and that could be used for this task was XGBoost.

Ultimately, XGBoost was chosen because of its capability to incorporate the notions of boosting, regularization, and gradient descent optimization into one algorithm that reaches a predictive performance beyond that of conventional models. The literature strongly verifies the efficiency of XGBoost when different features are available and noisy data are present, two common characteristics of real data, and thus is very appropriate for academic prediction. To this effect, Chen & Guestrin (2016) and Ahmad et al. (2020) contend that XGBoost outperforms other algorithms in model performance. In contrast to the performance of the Linear Regression model, which proved inefficient due to the presence of nonlinearities, and SVR models, which speculated a relatively high computational complexity at the cost of moderate accuracy, XGBoost achieved a remarkable trade-off between performance and computational complexity.

Compared to the performance of the Random Forest model, which had a good performance with a slight indication of overfitting, XGBoost achieved generalized performance with lower error margins and reduced variance in the train-test dataset. These results are consistent with the observations made and studies previously conducted within the field of Educational Data Mining, which validated that XGBoost performed considerably better than ensemble and linear models (Doz et al., 2023; Wang et al., 2021).

This is appropriate and accurate, considering that a thorough evaluation was performed on the model's performance, the feature importance analysis, and other techniques like the learning curve and cross-validation. By using XGBoost, the model was able to derive crucial information and key data that could be used effectively in the field of education. In the result of the project, it can be said that there is a good level of alignment between the original aims and the final performance of this model.


## 2. Project Deliverables Assessment

The evaluation showed that XGBoost met and exceeded the expectations of the project regarding the prediction of student academic performance. Compared to the Linear Regression, Support Vector Regression, and Random Forest models discussed earlier in this work, XGBoost yielded the highest accuracy, lowest prediction error, and highest generalization capability. Its $R^2$ score of 0.81 has shown that the model is able to explain 81% of the variance in final exam scores, thus providing a robust backbone for early detection of students at academic risk. Besides, XGBoost has reached the minimum MAE (1.71), the minimum MSE (8.07), and a competitive RMSE (2.84) compared to all the other models.

This performance advantage is rooted in the ability of XGBoost to combine multiple weak learners through gradient boosting with regularization for the reduction of overfitting. This has been recursively established through the literature on machine learning, in particular, within academic prediction contexts where data is bound to contain noise, non-linear interactions, and mixed feature types (Alshabandar et al., 2022; Huynh-Cam et al., 2021). Compared to Random Forest-which performed strongly but showed slight overfitting-XGBoost showed better convergence between training and validation curves; hence, it is at a better bias-variance balance. Thus, XGBoost is particularly appropriate to be used in real educational situations that call for accurate predictions on unseen data (Tapio, 2025).


Additionally, XGBoost yielded interpretable feature importance outputs and pointed out key predictors such as past exam scores, attendance rate, and study hours per week. Such insight adds not only value to the model as a predictor, but also gives educators insight into what factors may comprise academic indicators of success in students. Despite these advantages, the most important drawbacks to the usage of XGBoost are the higher computational cost compared with simpler models and that boosted trees can't be interpreted as easily as linear models. However, such limitations can be mitigated with the use of explainability tools such as SHAP and

LIME that translate complex model behaviour into human-friendly interpretations (Lundberg & Lee, 2017).

## 3. Lessons Learned

The project gave a lot of meaningful lessons that shaped our approach in both technical and analytical ways. First, we realized how important data preprocessing is. Handling missing values, outlier detection, encoding categorical variables, and standardizing the dataset are some of the steps that greatly affect the accuracy and stability of all the models. Clean data lets models learn better and reduces noise that could influence predictions.

We realized that model interpretability is equally important as model performance, if not more so, to stakeholders in education who rely on transparent and actionable insights. While XGBoost had good predictive capability, the internal decision-making process was a challenge to explain in simple terms. That was an important lesson in integrating explainability tools into model development. Thirdly, we grasped the importance of performing comparative modelling. The various algorithms taken into consideration, including linear, kernel-based, ensemble, and boosting models, gave us the ability to understand strengths and flaws in each approach.

This was done running all the models on the same set of data and under similar conditions to be fair in the comparison process. Lastly, we learned that collaboration and reproducibility are embedded parts of any machine learning project. By maintaining the use of Jupyter Notebooks, versioning of code with GitHub, and shared documentation consistently, we allowed every team member to review, validate, and replicate the results. This made the working process more structured, as is commonly followed in a professional data science environment. 4.

## 4. Recommendations

4.1 Technical and Model Improvement Recommendations

The model can be further improved in the future by implementing a number of improvements: first of all, the introduction of new variables related to students' behaviour and psychological features, such as motivation level, online participation, scores on self-regulation, and study strategies. This would give a deeper understanding of students' behaviour and increase the accuracy even more. According to the literature, the above factors have a high predictive value for academic outcomes. Besides the above-mentioned solutions, deployment of explainable AI-XAI tools, such as SHAP or LIME, is recommended. By engaging XGBoost internal calculations, these techniques can translate them into visual explanations that will help educators understand why a particular prediction was made. This will improve trust in the models and make them more transparent, both being important when predictive models are used in education. Thirdly, hybrid predictive models may be considered, which can also include the integration of XGBoost with deep learning methods like neural networks. They capture linear, nonlinear, and high-dimensional patterns in data better than any single model,

according to a review by Ahmad et al. (2020). This may potentially further reduce prediction error and strengthen generalizability across diverse datasets.

## 4.2 Future Research and Development

Moreover, this area of future research needs to extend the dataset from a single institution to increase diversity and generalisability. This could involve multi-institutional datasets or integration of data over several academic years to enhance the reliability of results. Besides, the employment of a longitudinal design in research would afford educational institutions an opportunity to track how early interventions based on the model impact academic performance over time.

There is still a lot of further research to be done in model fairness and bias reduction. Predictive models will need to be screened against the possibility of inadvertently creating disadvantages for certain groups of students due to their sensitive characteristics, such as gender, socio-economic background, or language. This will be according to newly developing international standards in the area of ethical AI for education (UNESCO, 2021).

## 4.3 Implementation Considerations

The predictors developed using XGBoost should be integrated into educational institutions' Learning Management Systems for the effective deployment of this model, which generates real-time alerts on at-risk students. Implementing automated visual dashboards would help the instructors track changes in performance and help with early intervention.

In addition, the institution has to create mechanisms for continuous monitoring, evaluation, and periodic retraining to keep the model performance high. Attention should also be paid to issues of computational resources, server infrastructure, and scalability to ensure efficient performance of the systems in view of growing data volumes.

## 5. Ethical and Responsible AI Practices

Our team believes that the basis of fairness, transparency, and accountability will provide the ethical deployment of AI in educational settings. Fairness may be assured by regular audits of the bias inside the predictive models so that the predictions made via XGBoost are not disproportionately reliant on sensitive or non-academic variables. This makes sure the model maintains values of equity and fairness (UNESCO, 2021). We also recommend that all predictions generated by the model should be accompanied by explainable outputs so that educators can understand why a particular student was flagged as at risk. Clear explanation mechanisms enable responsible decision-making and discourage blind reliance on automated predictions.

Finally, it is important that accountability be ensured within a framework where model configurations, datasets, and outputs are version-controlled and well-documented. Continuous monitoring should ensure that predictive outcomes align with the institutional values of inclusivity and support the well-being of students.

## 6. Data Governance and Compliance

In the light of the above, we recommend that a sound data governance strategy be put in place as part of future deployments of the XGBoost predictor to ensure that it remains compliant with institutional and international data protection laws. The proposed governance strategy should emphasize tight access control, encryption of sensitive records, and security storage protocols. All datasets entering the modelling pipeline should be thoroughly anonymized before going into training to remove personally identifiable information in compliance with both POPIA and GDPR regulations.

We also recommend that structured data documentation practices be introduced, such as comprehensive metadata, lineage tracking of datasets, and audit trails, where all manipulations of data are traceable and verifiable by third-party reviewers. Regular reviews for compliance should be performed in order to monitor risks and ensure alignment to regulations and further assure transparency. The aforementioned measures will help to strengthen institutional trust in AI-driven decision support and foster responsible long-term adoption.

## Conclusion

This project was therefore able to identify XGBoost as the champion model for the prediction of students' academic performance. Extensive testing, evaluation, and comparisons with other algorithms in machine learning showed that, among the best, XGBoost tops in accuracy, generalization, and reliability of interpretations. The project contributed meaningful insight into the strong importance of past exam scores, attendance, and study habits, showing how predictive modelling can support early intervention in higher education. Our last deliverable reflects not only the technical success of the model but our dedication as a group toward ethical AI use, data governance, and transparent model evaluation. This work sets a very strong foundation for educational institutions looking to adopt data-driven solutions that will enhance student success and equity, in addition to recommending improvements and future research directions.

**References**

Ahmad, M., Basheri, M., Ghani, I., & Rahim, N. (2020). Machine learning techniques for student performance prediction: A systematic review. *IEEE Access*, 8.

Ahmed, M., Ali, N., & Iqbal, Z. (2025). Predictive modelling in higher education using ensemble learning. *Journal of Educational Data Science*.

Alshabandar, R., Hussain, A., Liatsis, P., & Keight, R. (2022). Intelligent student performance prediction framework using machine learning. *Applied Intelligence*.

Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. *Proceedings of the 22nd ACM SIGKDD Conference*.

Doz, C., Rajan, D., & Muller, M. (2023). Evaluating boosting algorithms in educational datasets. *Expert Systems with Applications*.

European Commission (2023). *AI Regulation and Data Protection in the European Union*.

Ghosh, S., Kumar, S., & Agarwal, P. (2021). Factors influencing academic dropouts in higher education. *International Journal of Education Research*.

Huynh-Cam, A., Nguyen, T., & Lee, D. (2021). Error analysis in ensemble regression models for prediction tasks. *Pattern Recognition Letters*.

IET Conference Proceedings (2022). Applications of machine learning in learning analytics.

Lundberg, S., & Lee, S. (2017). A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems (NIPS)*.

Pedregosa, F., et al. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*.

Tapio, J. (2025). Understanding bias-variance in tree-based models. *Machine Learning Review*.

Trends in Neuroscience and Education (2023). Behavioural factors influencing academic success.

UNESCO (2021). *Recommendation on the Ethics of Artificial Intelligence*.

Wang, Z., Li, X., & Cheng, J. (2021). Boosting-based models for academic performance prediction. *Education and Information Technologies*.

Bunkar, K., Singh, U.K., Pandya, R. & Kushwah, R.S., 2012. *Data Mining: Prediction for Performance Improvement of Graduate Students Using Classification*. International Journal of Computer Applications, 55(1), pp. 35–39.

Cortez, P. & Silva, A., 2008. *Using Data Mining to Predict Secondary School Student Performance*. EUROSIS—Proceedings.

Doz, F., Ramos, R. & Duarte, L., 2023. *Boosting Approaches for Academic Performance Prediction*. Expert Systems with Applications, 214, 119155.

Ghosh, S., Ghosh, S. & Chakraborty, B., 2021. *Machine Learning Models for Student Dropout Prediction: A Systematic Review*. Education and Information Technologies, 26, pp. 491–509.

Huynh-Cam, T., Nguyen, N. & Tran, M., 2021. *Performance Analysis of Regression Models for Predicting Student Outcomes*. Journal of Educational Technology & Society, 24(3), pp. 88–102.

IIETA, 2022. *Machine Learning Models for Student Performance Evaluation*. International Information and Engineering Technology Association Proceedings.

Kesgin, K., Kiraz, S., Kosunalp, S. & Stoycheva, B., 2025. *Ensuring Fairness in Student Academic Performance Prediction with Machine Learning*. Applied Sciences, 15(8409), pp. 1–18.

Kotsiantis, S.B., Pierrakeas, C. & Pintelas, P., 2004. *Predicting Students' Performance in Distance Learning Using Machine Learning Techniques*. Applied Artificial Intelligence, 18(5), pp. 411–426.

Kumar, M. & Saha, S., 2018. *An Empirical Study on Random Forest for Student Academic Performance Prediction*. International Journal of Computer Science and Mobile Computing, 7(6), pp. 120–127.

Smart Learning Environments, 2022. *Evaluating Machine Learning Models for Academic Success Prediction*. Smart Learning Environments, 9(14), pp. 1–18.

Tapio, L., 2025. *Bias–Variance Trade-Off in Ensemble Learning Models for Educational Prediction*. Journal of Machine Learning Research, 26, pp. 112–129.

Trends in Neuroscience and Education, 2023. *Behavioural and Cognitive Predictors of Academic Success: A Data-Driven Review*. Trends in Neuroscience and Education, 33, 100053.