



OPERATIONS
RESEARCH
CENTER



Student Performance Analytics

15.072: Advanced Analytics Edge

Fall 2024

Nidhish Nerur, Brimar Ólafsson, James Pinter, Angel Xie, Yuki Yu

1 Problem Motivation

Education is crucial to shaping future opportunities, making it vital to understand factors influencing student performance. Our project focuses on predicting academic success for the early identification of struggling students. With these predictions, targeted interventions, such as tutoring or counseling, can improve outcomes and help students reach their potential. By identifying key factors for academic success and challenges, we aim to guide resource allocation and educational policies. Ultimately, this project seeks to provide data-driven recommendations to foster equitable and successful learning environments.

Historically, educators have utilized descriptive statistics and regression analysis to examine factors affecting student performance. These traditional methods effectively identify correlations between variables such as attendance and socioeconomic status with academic outcomes, but fail to capture non-linear relationships. A recent study on forecasting student academic performance demonstrated that machine learning algorithms, including Naïve Bayes, K-Nearest Neighbor, Support Vector Machines, and Neural Networks, have significantly improved prediction accuracy when applied to educational datasets, outperforming conventional models (Dawar et al., 2024). Our project seeks to leverage similar machine learning methodologies to uncover deeper insights into the key determinants of student academic performance.

2 Dataset Description

Kaggle Dataset We have a dataset from Kaggle, comprising 6,607 student records with 20 features. These features encompass various aspects of student life, including classroom discipline (e.g., Hours_Studied, Attendance), home environment (e.g., Parental_Involvement, Family_Income, Internet_Access), and extracurricular activities (e.g., Extracurricular_Activities, Physical_Activity). The target variable for our predictive models is the student’s final exam score, providing a quantitative measure of academic performance.

Dataset Source: <https://www.kaggle.com/datasets/lainguyn123/student-performance-factors>

2.1 Data Augmentation

Synthetic Data Generation Given the relatively small size of our combined dataset (6,607 observations), we supplemented it with synthetic data. We utilized Synthetic Minority Oversampling Technique (SMOTE) to duplicate instances of minority classes in classification. Additionally, we utilized Gaussian Mixture Models to estimate numeric features in a way that resembled the underlying distribution of data. Our augmentation approach aims to provide a more diverse and balanced dataset with additional observations, potentially enhancing our models’ predictive capabilities. However, we carefully evaluate the benefits of synthetic data against the risk of overfitting. The results for the augmented models are found in the Appendix.

2.2 Data Processing

Our preprocessing approach follows the data science pipeline, encompassing exploratory data analysis, cleaning, and feature engineering. We view the correlation matrix between numeric features and do not see significant risk of multi-collinearity:

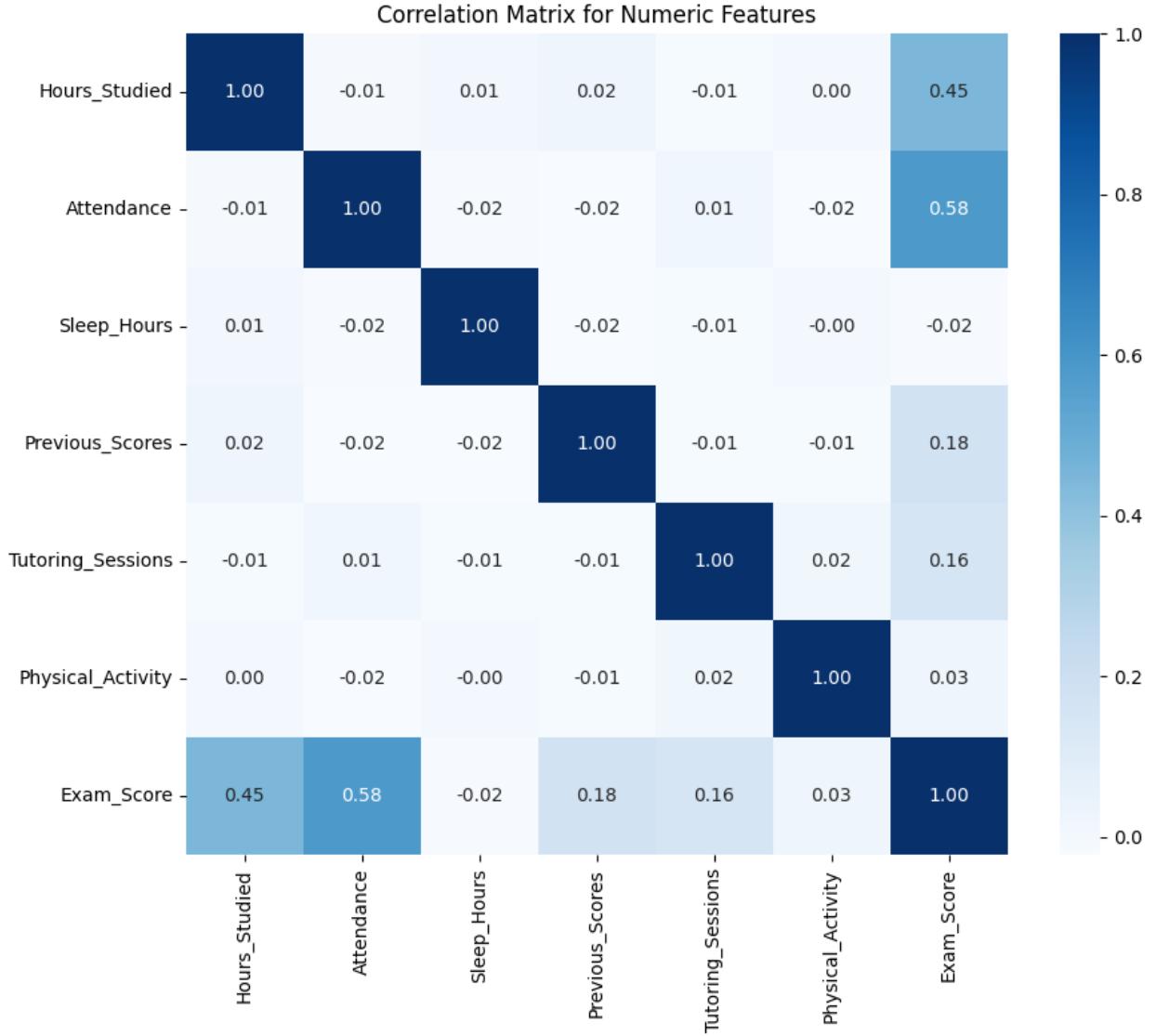


Figure 1: Correlation Matrix

Then, we found linear relationships between attendance rate and exam score as well as the hours studied and exam score, which are shown in the appendix. We created new interactions between hours studied and attendance, learning engagement, home support levels, study to sleep ratio, and balance of physical activity with sleep. Finally, we handled missing values, encoded categorical variables, and normalized numeric features to use machine learning algorithms.

3 Regression Task

For the regression task, we tested linear regression models, decision trees, and ensemble methods such as random forest, gradient boosting, and neural networks. We built baseline models and highlight the results:

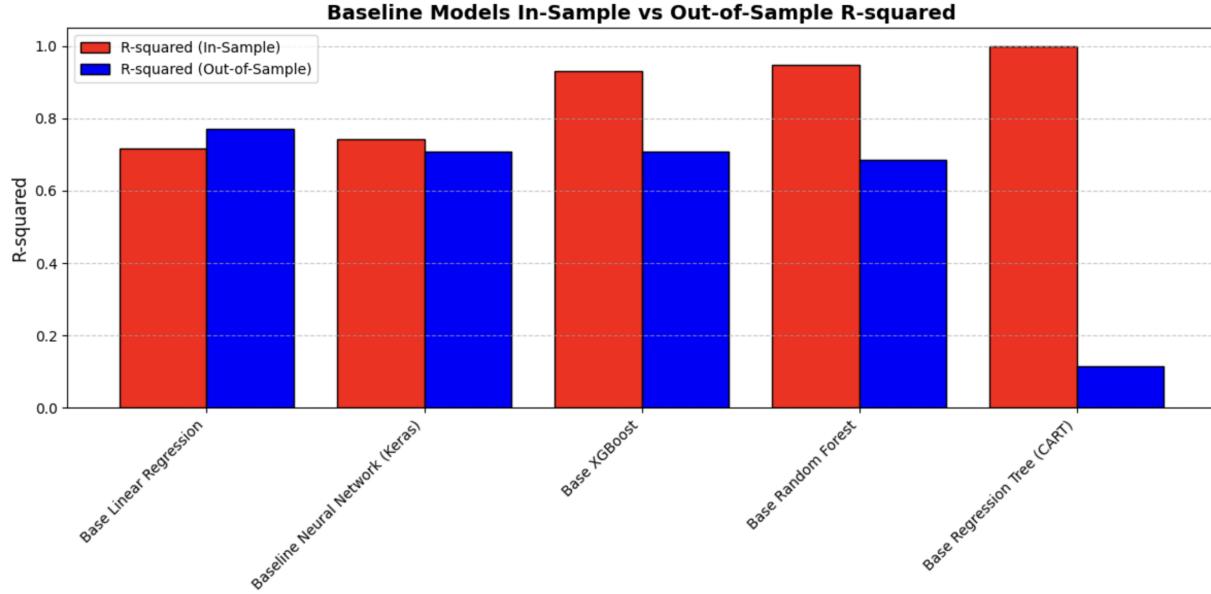


Figure 2: Baseline Regression Plots

Notice the XGBoost, Random Forest, and CART tree appear to be heavily overfitting to the training set, given their in-sample R-squared scores are much greater than the out-of-sample R-squared values. Consequently, we perform hyperparameter tuning across all our models and add the Lasso, Ridge, ElasticNet models. Furthermore, we used Bayesian Optimization to find optimal hyperparameters by balancing exploration of new values and refinement of promising ones, reducing the need for exhaustive searches. The table below shows the performance metrics of tuned models that performed better than their baselines.

Table I: Tuned Model Performance Metrics

Model	R^2 (In)	OOS R^2	MSE (In)	OOS MSE	MAE (In)	OOS MAE
ElasticNet	0.7169	0.7716	4.3529	3.23172	0.5011	0.4513
Ridge Regression	0.7171	0.7711	4.3506	3.2354	0.4970	0.4483
XGBoost	0.7739	0.7500	3.4767	3.5325	0.6280	0.6767
Lasso Regression	0.6827	0.7436	4.8793	3.6240	0.8317	0.7573
CART	0.6478	0.6353	5.1541	5.1541	0.9473	1.2378

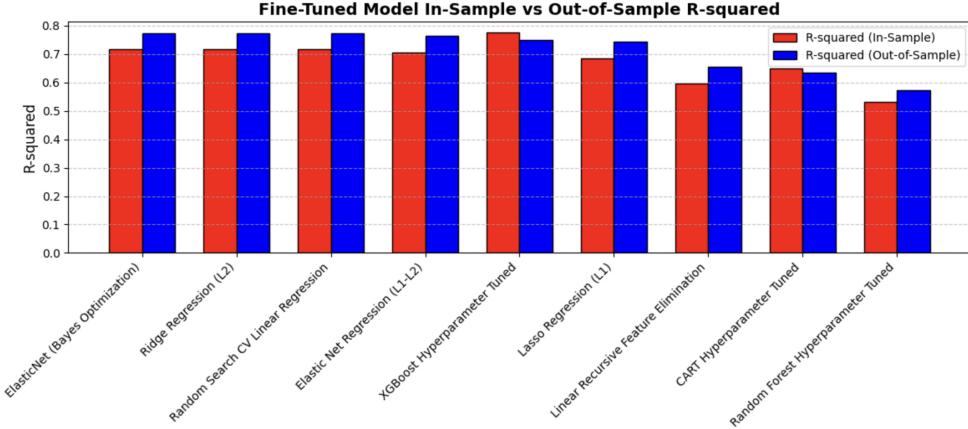


Figure 3: Hyperparameter Tuning Results for Different Models

Finally, here are our results from ensembling the top regression models:

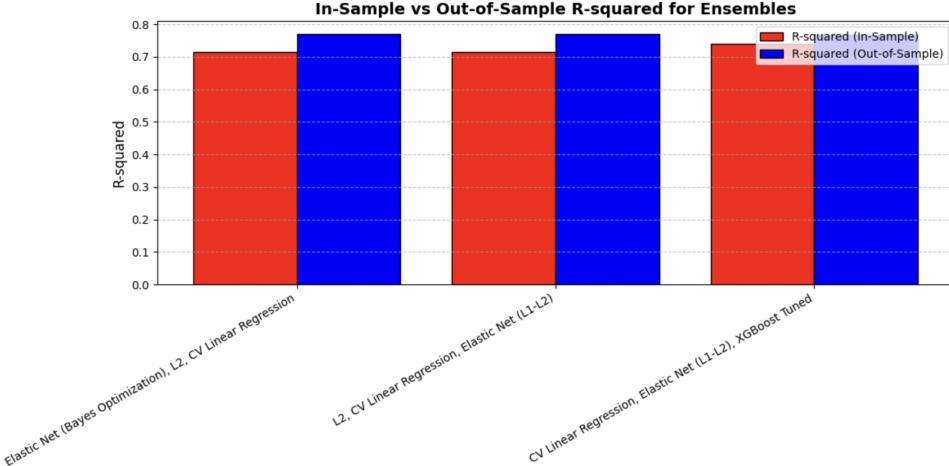


Figure 4: Ensemble of Top Regression Models

3.0.1 Lasso Regression

We used `RandomizedSearchCV` to tune Lasso’s α hyperparameter with 5-fold cross-validation. The best α was 0.1, resulting in a mean squared error (MSE) of 3.624 and an R^2 score of 0.743. Bayesian Optimization was not used, as Lasso showed the largest performance drop compared to Ridge and ElasticNet, likely due to its stricter sparsity constraints. Our best lasso model has 14 non-zero features, making it a bit more interpretable than models including all the features. From a business perspective, we can directly interpret the statistical significance of features and achieves comparable performance to the other models tested.

3.0.2 Ridge Regression

We used Bayesian Optimization to tune the α (regularization strength) hyperparameter for Ridge regression. The best α found was approximately 3.576, achieving a cross-validation score of -4.401 , a test MSE of 3.234, and an R^2 score of 0.771. Ridge is less interpretable than lasso given it does not zero out any of the coefficients, but we can still intuitively explain the relationship between the features and the final exam score outcome.

3.0.3 Elastic Net Regression

We used Bayesian Optimization to tune the α (regularization strength) and l1_ratio (balance between Lasso and Ridge) hyperparameters for ElasticNet. The optimization identified the best parameters as $\alpha \approx 0.0039$ and l1_ratio ≈ 0.5136 , achieving a cross-validation score of -4.399 , a test MSE of 3.232, and an R^2 score of 0.772. Elastic Net utilizes a combination of L1 and L2 regularization, so it can zero some terms while retaining others. This model had strong performance with out of sample R-squared of around 76%, close to our top three models, and would make a practical choice in a business setting. Given this is our highest performing model, we show the most important features selected:

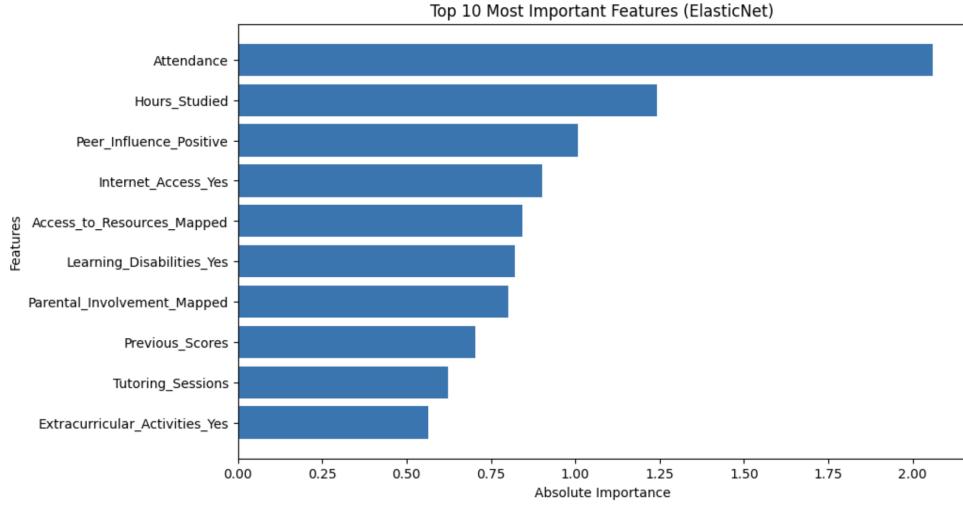


Figure 5: Most Important Features ElasticNet

Logically, attendance rates, number of hours studied, and positive peer influence make sense in determining a student's exam score. It appears that attendance is the most notable feature by quite a large margin, and we expect students who attend class more often are likely understanding the course material better than those who do not. ElasticNet balances interpretability of coefficients with the strongest performance.

3.0.4 CART

We used Bayesian Optimization to tune the hyperparameters for the Decision Tree (CART) model. The optimization process identified the best parameters as:

- `max_depth: 21`
- `min_samples_split: 2`
- `min_samples_leaf: 16`

These parameters achieved a cross-validation score of -6.988 , a test Mean Squared Error (MSE) of 5.810, and an R^2 score of 0.589. CART is highly interpretable, given we can map precisely how to arrive at a particular exam score prediction. However, it had the second lowest performance, with an out of sample R-squared of 63.54

3.0.5 Random Forest

The Random Forest Regressor performed close to the XGBoost model without hyperparameter tuning, but these models were still worse than our linear models overall. After hyperparameter tuning, Random Forest had the lowest out of sample R-squared score below 60%, which is surprising. Since our dataset has primarily linear relationships with the target, it seems the tuned CART tree outperformed random forest. We experimented with the number of trees, maximum number of features, and max depth of the tree, yet it still had the worst performance. Additionally, this model is not interpretable as it aggregates decision trees.

3.0.6 XGBoost

The best Bayesian Optimization with Hyperopt XGBoost hyperparameters include:

- `max_depth`: 3,
- `gamma`: 1.47,
- `reg_alpha`: 40,
- `reg_lambda`: 0.96,
- `colsample_bytree`: 0.97,
- `min_child_weight`: 3.

These parameters achieved a test MSE of 3.474 and an R^2 score of 0.754. XGBoost is also not interpretable as it aggregates shallow decision trees, but it achieved moderate performance. The linear ElasticNet and ridge regression models still outperform XGBoost for our dataset.

3.0.7 Neural Networks

We used Bayesian Optimization with Hyperopt to tune the hyperparameters of a neural network. The best configuration utilized ReLU with SGD optimizer and a batch size of 16. This resulted in an out-of-sample Mean Squared Error (MSE) of 4.076 and an R^2 score of 0.712. Bayesian Optimization enabled efficient exploration of the hyperparameter space and refinement of the model. The neural network model is also a black-box which is not interpretable, and it did not perform that well on our small dataset.

3.1 Regression Recommendation

The hyperparameter-tuned ElasticNet Regression model achieved the best performance, with an out-of-sample R^2 of approximately 0.7716. This is a marginal improvement over the baseline linear regression model, which achieved an out-of-sample R^2 of 0.7710. Ultimately, we recommend using ElasticNet as we retain 22 of the 27 features, so we can somewhat more easily identify the notable features relative to base linear regression. Furthermore, the coefficients for ElasticNet still have strong statistical properties and we can explain their effect on the outcome exam score value. ElasticNet maintains explainability with the best out of sample performance.

4 Classification Task

We tested 8 different models, logistic regression, support vector machines, k-nearest neighbors, naive bayes, decision trees, and ensemble methods such as random forest, XGBoost, and neural networks. We binned our final exam score into 4 ranges of low, average, good, and excellent with roughly even proportions. The model pros and cons as well as interpretability considerations are similar to what we discussed above with the linear regression models. The newly tested and best performed models of logistic regression and support vector classifiers will be discussed further. Table II presents the results of fine-tuned classification models, sorted by out-of-sample performance.

Table II: Fine-Tuned Model Performance Metrics

Model	Accuracy (In-Sample)	Accuracy (OOS)	F1 (In-Sample)	F1 (OOS)
Support Vector Machine	0.962	0.963	0.962	0.963
Logistic Regression	0.952	0.955	0.952	0.955
Neural Network	0.964	0.937	0.964	0.938
XGBoost	0.996	0.870	0.996	0.871
Random Forest	0.997	0.759	0.997	0.756
Naive Bayes	0.679	0.695	0.683	0.700
K-Nearest Neighbors	1.000	0.685	1.000	0.688
CART	0.855	0.655	0.854	0.655

We also highlight the baseline and tuned model Receiver Operating Characteristic Curves:

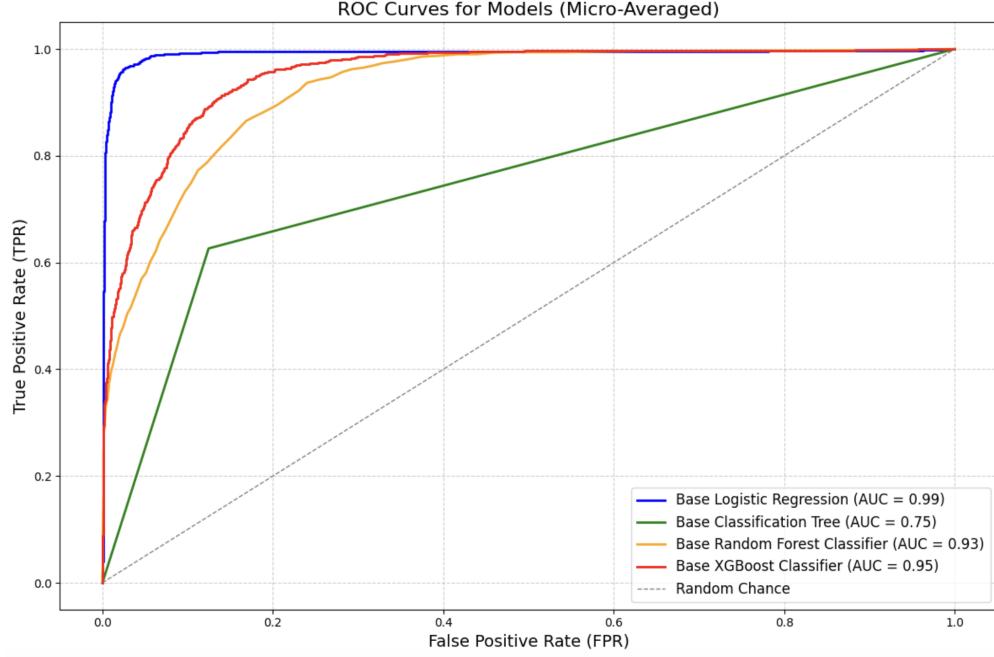


Figure 6: Baseline ROC Curves

The ROC curves reveal a clear improvement in model performance after fine-tuning. In the baseline models, Logistic Regression and Random Forest demonstrate strong discrimination capabilities with AUC scores of 0.93, while CART lags with an AUC of 0.75, highlighting its limitations in the baseline configuration.

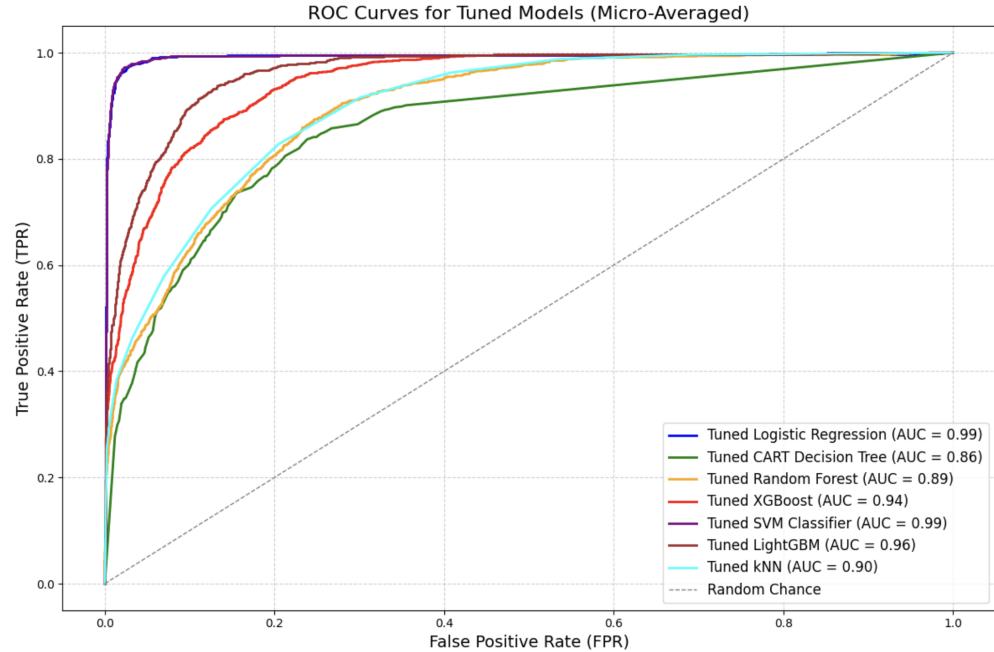


Figure 7: Tuned ROC Curves

After fine-tuning, nearly all models exhibit substantial performance gains, with Logistic Regression and SVM achieving near-optimal AUC scores of 0.99. The performance gaps between models are narrow, indicating that tuning mitigates variability in performance across algorithms. However, CART remains an under-performer even after tuning, suggesting inherent limitations in its suitability for this classification task. Similarly from table II, tuned models display improved consistency and effectiveness, with Logistic Regression, SVM, and Random Forest emerging as the most reliable classifiers.

4.1 Logistic Regression

Best parameters identified for Logistic Regression:

- **C:** 10
- **Penalty:** L1
- **Solver:** saga
- **Max Iterations:** 500

The logistic regression model appears to perform best initially without hyperparameter tuning relative to the other models. This model is reasonably interpretable as we can assess the direct impact of the features on the log odds of the outcome variable. We subsequently fine-tune the complexity parameter which essentially regularizes the coefficients, and we experiment with balancing the class weights. The fine-tuned model achieved an out-of-sample accuracy of 95.5% and an F1 score of 95.5%. This is a strong model choice for the schools to implement, as it balances some interpretability with high performance.

4.2 Support Vector Classifiers

Best parameters identified for SVM:

- **C:** 100
- **Kernel:** Linear
- **Gamma:** Scale

Support vector classifiers performed on par with the tuned logistic regression model, indicating both offer robust performance. It achieves an out-of-sample accuracy of 96.3% and an F1 score of 96.3%. SVCs can handle high-dimensional data well and create boundaries to maximize the margin between the support vectors and given classes. This model can also capture non-linear relationships depending on the kernel. SVCs are somewhat challenging to interpret as they are mathematically rigorous.

4.3 Classification Recommendation

We recommend Logistic Regression as the primary model for schools to predict student exam score ranges. This model achieves an out-of-sample accuracy of 95.5% and offers a strong balance between performance and interpretability. Its ability to provide probabilities for outcome categories allows for granular assessments of student performance, making it highly practical for identifying at-risk students. While SVM performed similarly, it lacks the same level of interpretability. Models like XGBoost and Random Forest seemed to overfit while CART underperformed significantly. Schools can use logistic regression predictions to guide interventions such as tutoring, counseling, or study groups, ensuring resources are effectively directed to improve outcomes and foster equitable learning opportunities.

5 Future Direction

With our recommendations, schools can utilize ElasticNet regression model to predict exam scores and the Logistic Regression model to classify the performance range. These methods allow schools to determine the primary features affecting academic performance while maintaining strong performance potential. Consequently, they can identify at-risk students and provide additional support to improve academic outcomes. In the future, we would acquire more student data, further fine-tune the model parameters, and enhance machine learning with an optimization task to allocate resources to particular students. Our models provide a strong foundation to help school districts decide the best approach for their students to succeed.

6 Appendix

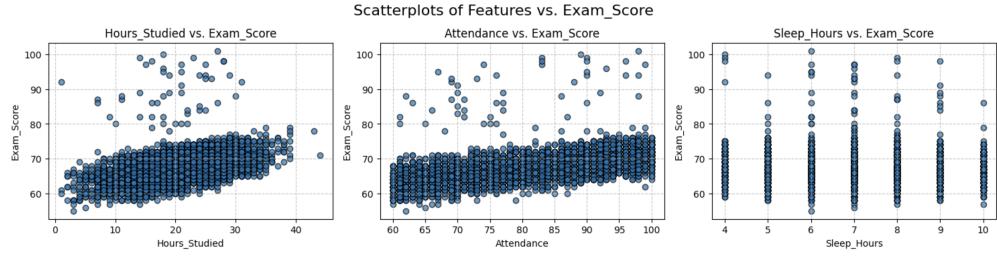


Figure 8: Relationships with Outcome

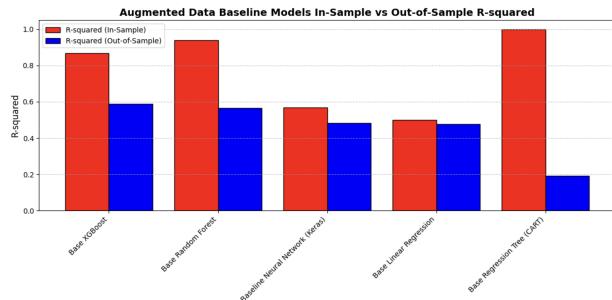


Figure 9: Augmented Data Baseline Results

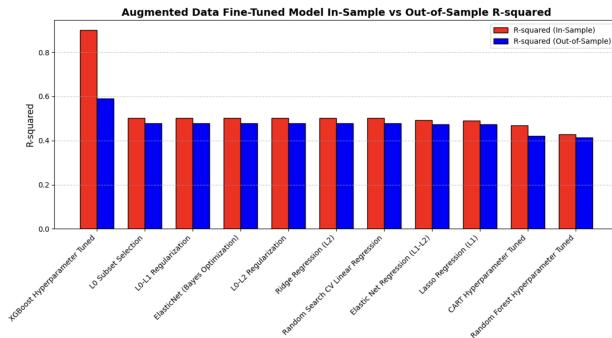


Figure 10: Augmented Data Tuned Results

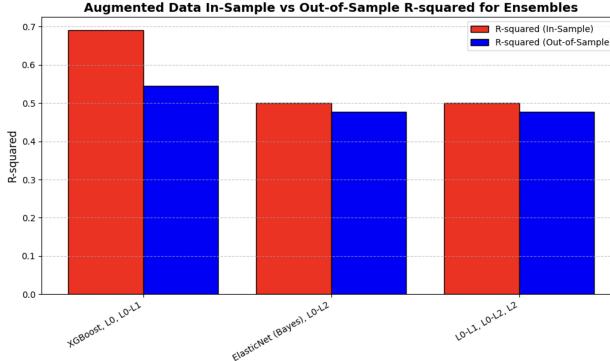


Figure 11: Augmented Data Ensemble Results

Model	R-squared (In-Sample)	R-squared (Out-of-Sample)	MSE (In-Sample)	MSE (Out-of-Sample)	MAE (In-Sample)	MAE (Out-of-Sample)
Base Linear Regression	0.7714611992305818	0.7710816276470018	4.350563974868887	3.2357778362359515	0.49681450833166685	0.4481122945667927
Baseline Neural Network (Keras)	0.7419852102394104	0.7097357511520386	3.964561506774853	4.102906119173954	0.8510619121569558	0.8730191379942367
Base XGBoost	0.9269869293928955	0.7096007595962256	1.0814553429931202	4.11094913140925	0.4391705169244056	0.783404723995029
Base Random Forest	0.948176664280157	0.6856264752109408	0.797092771902153	4.443601511205929	0.4130259001324496	1.087972768532525
Base Regression Tree (CART)	1.0	0.6189482392068542	8.0	12.5278789712596	0.0	1.8913767019667172

Figure 12: Baseline Regression Table Results

Model	R-squared (In-Sample)	R-squared (Out-of-Sample)	MSE (In-Sample)	MSE (Out-of-Sample)	MAE (In-Sample)	MAE (Out-of-Sample)
ElasticNet (Bayes Optimization)	0.7149929173952996	0.7731032364464269	4.3579214349030685	5.21597222854872	0.5011021612954711	0.451199773230239
Ridge Regression (L2)	0.71714542730125	0.7711021860233069	4.350756047172496	3.235497189454942	0.49707823583200284	0.44833916817216225
Random Search CV Linear Regression	0.7171461892305818	0.7710816276470018	4.350563974868887	3.2357778362359515	0.49681450833166685	0.4481122945667927
Linear Net Regression (L1-L2)	0.7046374024441397	0.762642586787766	4.54296116059586	3.3604081165908347	0.648682102145992	0.5821731042907829
XGBoost Hyperparameter Tuned	0.7739548683366504	0.7509057710388318	3.476792341207756	3.5325571034713762	0.628069831061247	0.67673429870782
Lasso Regression (L1)	0.6827688564764478	0.74361441802337031	4.87932045520541	3.624032132293864	0.8317469539325287	0.757354334499784
Linear Recursive Feature Elimination	0.9366151521389918	0.6540597471903234	6.704441045970849	4.890017760464271	1.022186150291856	1.2021168021885
CART Hyperparameter Tuned	0.6478714015013996	0.6336638730162024	5.41598647114475	5.154122516036305	0.9473892970671713	1.2178972556723223
Random Forest Hyperparameter Tuned	0.5306664436493116	0.5711553077444302	7.210798950561008	6.061750804717776	1.5891870167607653	1.53884947974958

Figure 13: Tuned Regression Table Results

Model	Accuracy (In-Sample)	Accuracy (Out-of-Sample)	ROC AUC (In-Sample)	ROC AUC (Out-of-Sample)	Precision (In-Sample)	Precision (Out-of-Sample)	Recall (In-Sample)	Recall (Out-of-Sample)
Base Logistic Regression	0.9468306527909177	0.9447069635409077	0.9597218298018233	0.960532302009123	0.9470914468913717	0.9452963536133738	0.9468306527909177	0.9447069635409077
Base Neural Network	0.9379375991296121	0.9144977307110439	0.8947215434574749	0.8941028119402913	0.9378844010036752	0.9124346511309790	0.9379375991296121	0.9144977307110439
Base XGBoost Classifier	1.0	0.737821534093344	1.0	0.948103612918376	1.0	0.775797474393951	1.0	0.737821534093344
Base Random Forest Classifier	1.0	0.732903328290468	1.0	0.916884733459462	1.0	0.73147936165363243	1.0	0.732903328290468
Base Classification Tree	1.0	0.626327318910741	1.0	0.7524369929418073	1.0	0.6298609043639778	1.0	0.626327318910741

Figure 14: Baseline Classification Table Results

Model	Accuracy (In-Sample)	Accuracy (Out-of-Sample)	ROC AUC (In-Sample)	ROC AUC (Out-of-Sample)	Precision (In-Sample)	Precision (Out-of-Sample)	Recall (In-Sample)	Recall (Out-of-Sample)
Tuned SVM Classifier	0.9578051087984482	0.9546142208774584	0.9937984552979398	0.991291358305191	0.9583437396438681	0.95199845312727	0.9578051087984482	0.9546142208774584
Tuned Logistic Regression	0.9538315988647115	0.9531013915733737	0.9934750417255698	0.99148322434121	0.9540528115938838	0.9534628156333605	0.9538315988647115	0.9531013615733737
Tuned LightGBM	0.999621570424977	0.8071104872791982	0.9999979324693857	0.9572947294906731	0.99962013867899	0.8117178478421	0.9996215704824977	0.807110487291982
Tuned XGBoost	0.997350933774835	0.770545855791225	0.9999714712670825	0.9371976112417354	0.9973560617783077	0.7724179916129876	0.997350933774835	0.770545855791225
Tuned KNN	0.7623462630905147	0.6724459606505081	0.9348022332044392	0.88514579301490998	0.7592381261346607	0.6677864477850473	0.7623462630905147	0.6724459606505081
Tuned Random Forest	0.7165562913907285	0.6488638124054463	0.9133038348636291	0.8880515346563476	0.760603832530216	0.6648703131495229	0.7165562913907285	0.6488638124054463
Tuned CART Decision Tree	0.7952968310312204	0.643721633880484	0.907986126999123	0.8519324737936453	0.7951934549263677	0.6439071162255183	0.7952968310312204	0.643721633880484

Figure 15: Tuned Classification Table Results