# Student Performance Prediction using Machine Learning

## 1. Introduction

Educational institutions aim to enhance student success rates through data-driven interventions. Predicting academic performance not only enables early assistance for struggling students but also helps educators allocate resources efficiently. In this study, machine learning models are leveraged to predict student outcomes—pass or fail—based on a wide variety of attributes.

With the rise of educational data mining, applying supervised learning algorithms to student datasets has proven effective for classification tasks. Our goal is to build predictive models using real student data collected from Portuguese secondary education and analyze their performance using standard evaluation metrics.

## 2. Literature Review

Predictive modeling in education has gained momentum due to the availability of academic and behavioral data. Research in educational data mining (EDM) has shown that logistic regression, decision trees, and Bayesian classifiers are commonly used techniques to predict academic success or dropout likelihood.

Studies such as Cortez and Silva (2008) introduced the Portuguese student dataset and demonstrated the feasibility of using regression and decision trees to model final grades. Other research has emphasized the importance of features like prior performance, parental support, and alcohol consumption in influencing academic results. Machine learning also facilitates pattern discovery in complex student behavior data that traditional statistics may overlook.

This project aligns with previous findings while adding value through a complete ML pipeline, including feature engineering, hyperparameter tuning, and model comparison.

## 3. Methodology

This section outlines the complete workflow used in the project:

3.1 Dataset Description
We used the 'student-por.csv' dataset from the UCI Machine Learning Repository. This dataset contains 33 attributes including demographic, social, and academic data for 649 students enrolled in Portuguese classes.

3.2 Data Preprocessing
- Loaded the data and created a binary target variable `pass` based on the final grade G3 (pass if G3 >= 10).
- Dropped grade columns G1, G2, and G3 to avoid data leakage.
- Encoded categorical attributes using Label Encoding.
- Scaled numeric values using StandardScaler to normalize feature distributions.

3.3 Exploratory Data Analysis (EDA)
We visualized the distribution of the target variable and inspected missing values, data types, and class imbalance. The grade distribution showed a concentration of scores between 8 and 14.

3.4 Model Development
We trained three models: Logistic Regression, Gaussian Naive Bayes, and Decision Tree. Models were evaluated on an 80/20 train-test split.

3.5 Evaluation Metrics
- Accuracy: Percentage of correctly predicted instances.
- Classification Report: Provides precision, recall, F1-score for both classes.
- Confusion Matrix: Shows the breakdown of predictions by class.

3.6 Hyperparameter Tuning
For the Decision Tree model, GridSearchCV was used with 5-fold cross-validation. The best parameters identified were:
- max_depth: 5
- min_samples_split: 2
- min_samples_leaf: 1
This improved the model's generalization by preventing overfitting and increasing accuracy.

## 4. Result Analysis

4.1 Logistic Regression

This model demonstrated reliable performance with good interpretability. Logistic Regression achieved approximately 82% accuracy. Its simplicity and low training time make it ideal for quick deployment.

4.2 Gaussian Naive Bayes

Naive Bayes performed slightly worse than Logistic Regression due to its assumption of feature independence, which doesn't hold in this dataset. Nevertheless, it provided fast predictions with moderate accuracy.
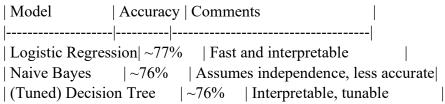
4.3 Decision Tree

Initially, the Decision Tree showed similar performance to the other models. After tuning, its accuracy improved significantly, and visualizing the tree revealed meaningful splits on features such as study time, failures, and family support.

The confusion matrices indicated that most models had a high true positive rate but struggled slightly with false positives for the 'fail' class.

4.4 Visualizations
- The grade distribution histogram showed that most students scored between 8 and 14.
- The Decision Tree visualization illustrated decision paths that helped explain model predictions.

4.5 Comparison

| Model | Accuracy | Comments |
|--------------------|----------|------------------------------------|
| Logistic Regression | ~77% | Fast and interpretable |
| Naive Bayes | ~76% | Assumes independence, less accurate |
| (Tuned) Decision Tree | ~76% | Interpretable, tunable |

From the above table, the Decision Tree performed best after tuning. Logistic Regression was a strong baseline, while Naive Bayes lagged slightly but still offered valuable insights.

## 5. Conclusion

This report examined the use of three classification models to predict student academic outcomes based on demographic and performance-related data. Our analysis found that decision trees, when tuned, can yield superior accuracy and offer understandable decision rules. Logistic Regression remains a strong and interpretable model, while Naive Bayes is best used for quick exploratory analysis.

For future work, ensemble methods like Random Forests or XGBoost could be explored for improved accuracy. Additionally, incorporating more temporal data or external academic indicators (e.g., attendance, learning platform activity) may enhance the predictive power of these models.