

Educational Institution Problem

By: Norman Dwi Febrio



Business Understanding

Background Story

Jaya Jaya Institute is one of the higher educational institutes established in 2000. Until nowadays, they have created many graduates with good reputations. Yet, many students cannot complete their studies well, alias dropout.

This high dropout number indeed become a big problem for the institute. Therefore, Jaya Jaya Institute wants to detect students who are probably will dropout ASAP so that the institute can give special guidance.

Business Understanding

Business Problems

Problem 1: High Dropout Students Number

The number of dropout student needs to be reduced to keep the institute's reputation and attraction. The higher dropout number reflects the **dissatisfaction** of students regarding education quality or other supports. This problem could cause a **reduction in new registrants**, a decrement of prospective students and their parents, and also make the institute reputation's **bad** in a public point of view and other stakeholders.

Problem 2: No Device nor System Which Could Monitor The Student's Condition

Without an effective monitoring system, institutions find it **difficult** to identify potential dropout students in a timely manner and provide necessary interventions. This makes dropout prevention efforts **less efficient and undirected**. without a monitoring system, institutions cannot proactively address the problems of students experiencing academic or personal difficulties that could lead to dropouts



Data Understanding

A dataset created from a higher education institution (acquired from several disjoint databases) related to students enrolled in different undergraduate degrees, such as agronomy, design, education, nursing, journalism, management, social service, and technologies. The dataset includes information known at the time of student enrollment (academic path, demographics, and social-economic factors) and the students' academic performance at the end of the first and second semesters. The data is used to build classification models to predict students' dropout and academic success.

Feature Name	Description
Marital_status	Marital status of Student
Application_mode	Kind of application used by student to apply on the institute
Application_order	Application order (0 - 9)
Course	Kind of course choose by student



Data Understanding

Feature Name	Description
Daytime_evening_attendance	Time attendance
Previous_qualification	Previous qualification of student
Previous_qualification_grade	Grade of the previous qualification
Nacionality	Nationality of student
Mothers_qualification	Student's mother qualification
Fathers_qualification	Student's father qualification
Mothers_occupation	Student's mother occupation



Data Understanding

Feature Name	Description
Fathers_occupation	Student's father occupation
Admission_grade	Admission grade (0 - 200)
Displaced	Is the student a displaced person?
Educational_special_needs	Does the student with educational special need?
Debtor	Is the student a debtor?
Tuition_fees_up_to_date	Is the student's tuition fees up to date?
Gender	Gender of the student



Data Understanding

Feature Name	Description
Scholarship_holder	Is the student a scholarship holder?
Age at enrollment	Previous qualification of student
Previous_qualification_grade	Grade of the previous qualification
Nacionality	Nationality of student
Mothers_qualification	Student's mother qualification
Fathers_qualification	Student's father qualification
Mothers_occupation	Student's mother occupation



Data Understanding

Feature Name	Description
International	Is the student an international student?
Curricular_units_1_st_sem_credited	Number of curricular units credited in the 1st semester
Curricular_units_1_st_sem_enrolled	Number of curricular units enrolled in the 1st semester
Curricular_units_1_st_sem_evaluations	Number of evaluations to curricular units in the 1st semester
Curricular_units_1_st_sem_approved	Number of curricular units approved in the 1st semester
Curricular_units_1_st_sem_grade	Grade average in the 1st semester (between 0 and 20)



Data Understanding

Feature Name	Description
<code>Curricular_units_1st_sem_without_evaluations</code>	Number of curricular units without evaluations in the 1st semester
<code>Curricular_units_2nd_sem_credited</code>	Number of curricular units credited in the 2nd semester
<code>Curricular_units_2nd_sem_enrolled</code>	Number of curricular units enrolled in the 2nd semester
<code>Curricular_units_2nd_sem_evaluations</code>	Number of evaluations to curricular units in the 2nd semester
<code>Curricular_units_2nd_sem_approved</code>	Number of curricular units approved in the 2nd semester
<code>Curricular_units_2nd_sem_grade</code>	Grade average in the 2nd semester (between 0 and 20)



Data Understanding

Feature Name	Description
<code>Curricular_units_2nd_sem_without_evaluations</code>	Number of curricular units without evaluations in the 2nd semester
<code>Unemployment_rate</code>	Unemployment rate (%)
<code>Inflation_rate</code>	Inflation rate (%)
<code>GDP</code>	GDP
<code>Status</code>	Graduate, Dropout, or Enrolled



Project Scope

**Exploratory Data Analysis
(EDA)**

Develop A Predictive Model

**Develop Prescriptive
Dashboards**

INSIGHTS

01

4424

Total Students

02

54.8%

Displaced Students

03

24.8%

Scholarship Holders

04

11.4%

Debtor Students

05

2.5%

International Students

06

1.2%

Students with Education Special Needs

Interesting Patterns

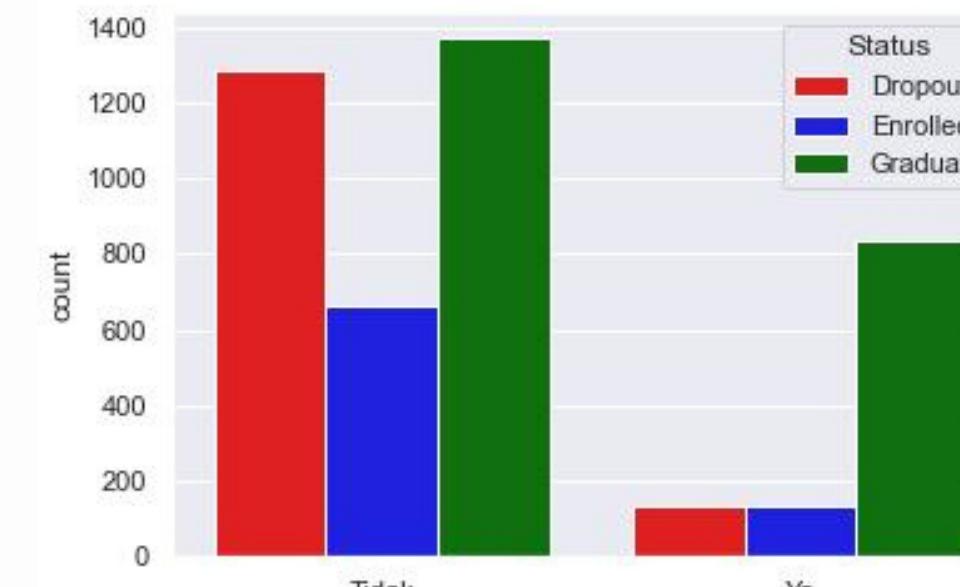
Student's Status Distribution by

Gender



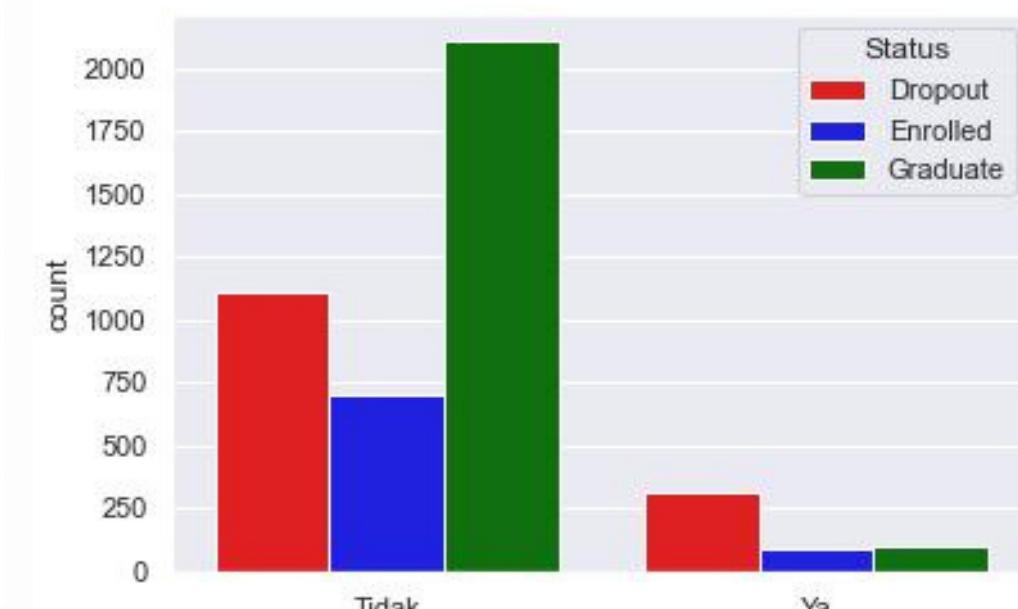
Men got more dropout student rather than women got

Scholarship Holder



Scholarship-holder students are more likely to graduate

Debtor

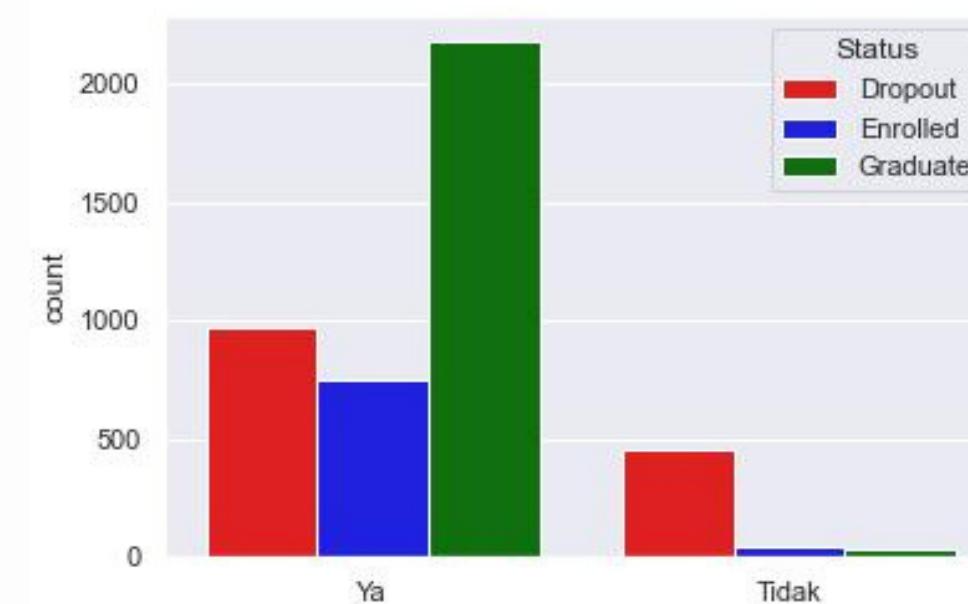


Debtor students tend to be dropped out students

Interesting Patterns

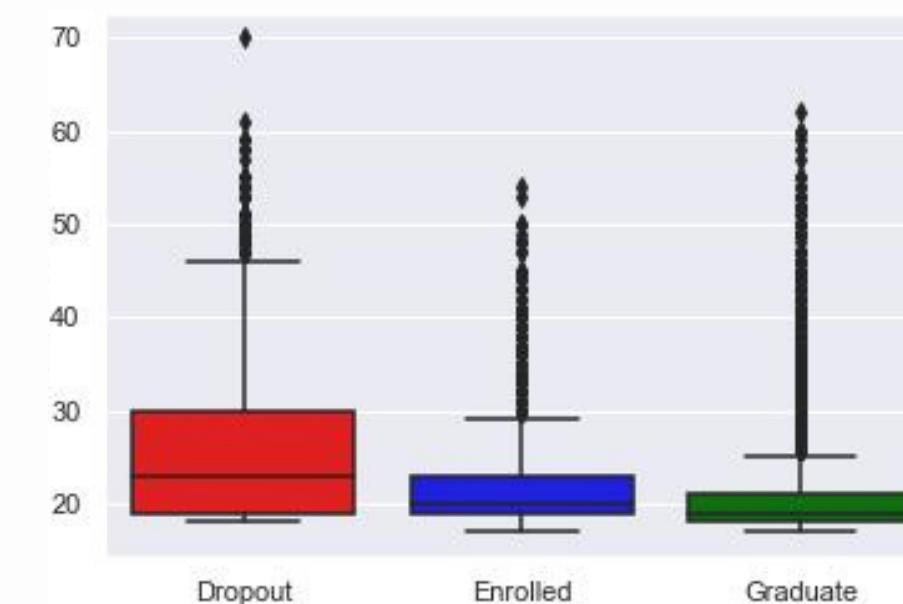
Student's Status Distribution by

**Tuition Fees
Up To Date**



Students whose tuition fees are not up to date are more likely to dropout

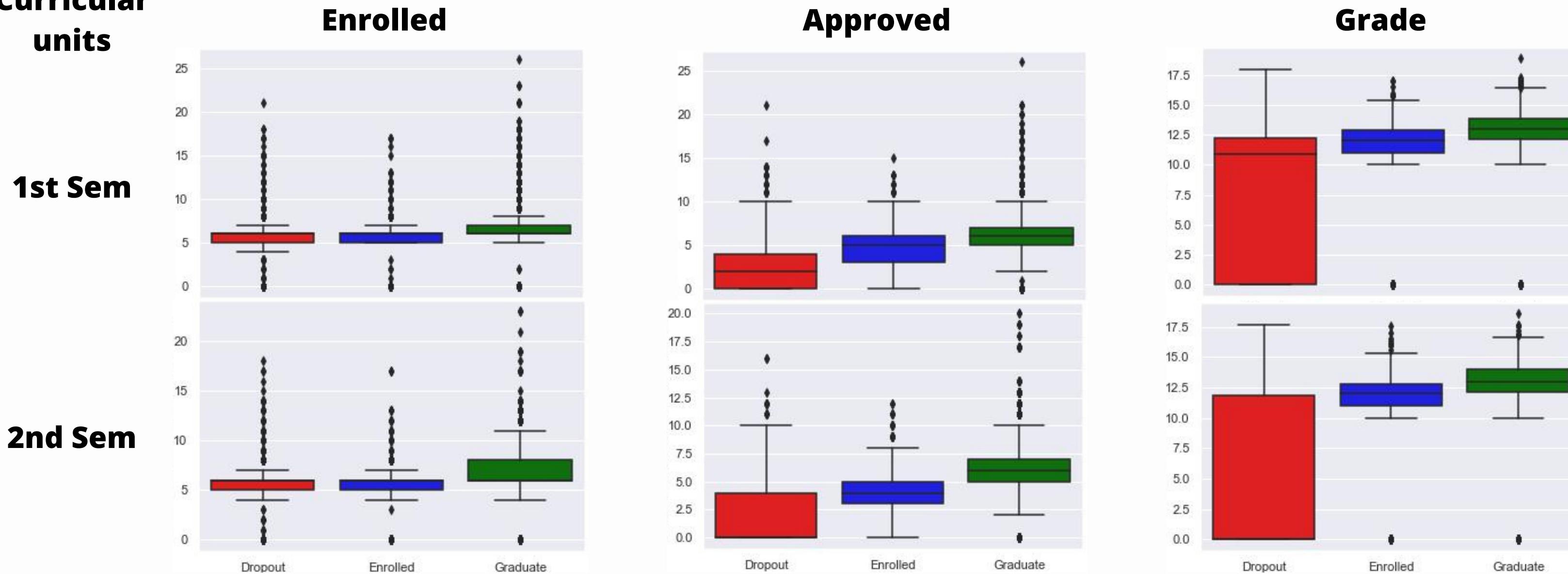
Age at Enrollment



The distribution of dropped out students have wider range

Interesting Patterns

**Curricular
units**

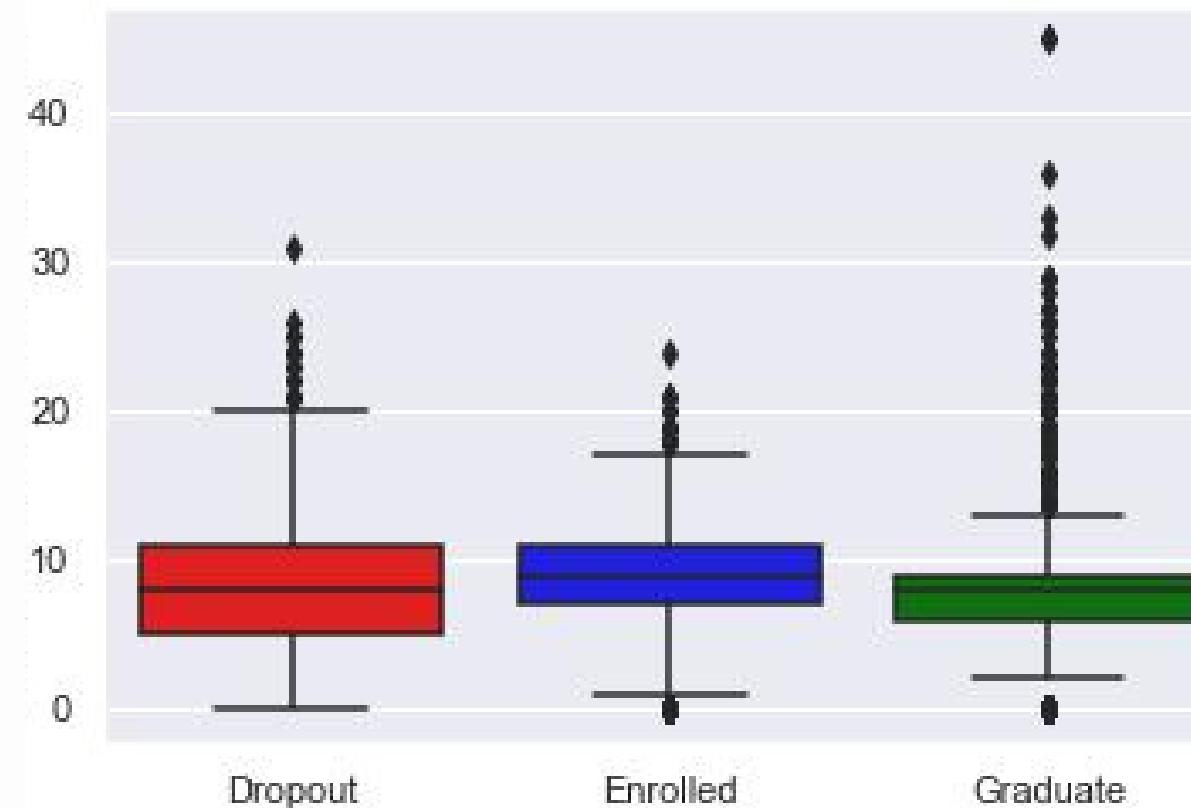


Similar patterns? Graduate students always have a higher median, either on enrolled, approved, or grade

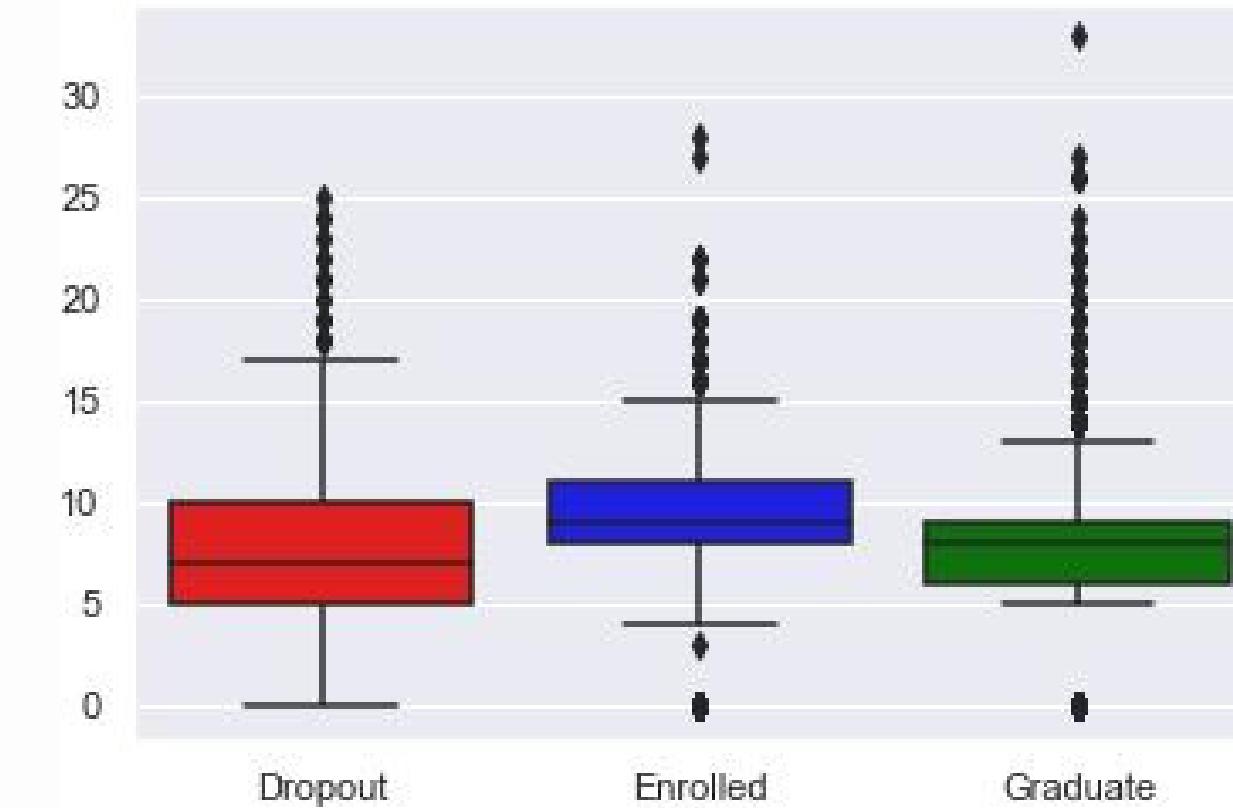
Interesting Patterns

Curricular units

1st Semester Evaluations

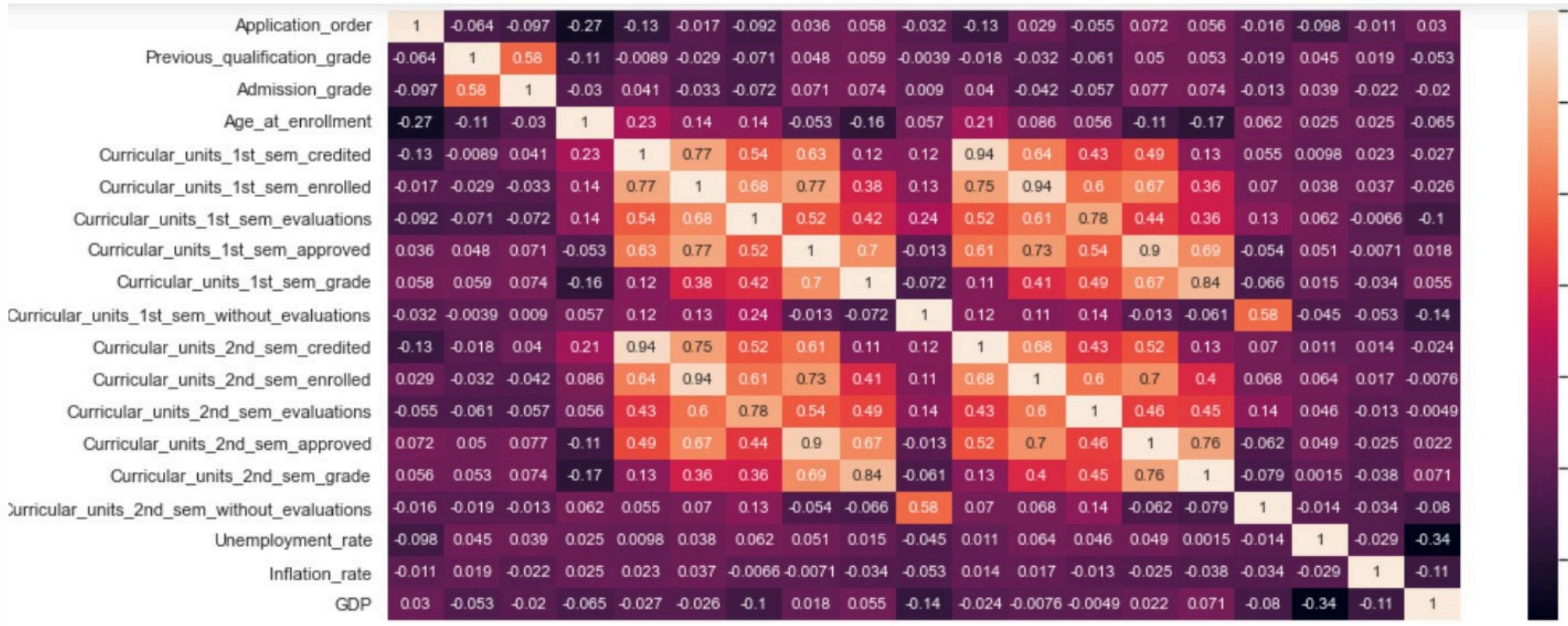


2nd Semester Evaluations



Dropout students always have **wider spread** and graduate students always have **narrower spread**.
Means that graduate students have **more consistent** in their evaluations rather than dropout students

Numeric Feature's Correlations



Some features are **highly correlated**. Those are features in the 1st semester correlated with its 2nd semester. This problem could cause **multicollinearity** and **reduce the model's accuracy**. Therefore, the **Principle Component Analysis (PCA)** technique will be used to preprocess these features.

Model Development

Data Preprocessing

Feature Selection:
All features in The Interesting Patterns

Numeric Features:

- PCA
- PowerTransformer

Categorical Features:

- One-hot Encoding

Data Splitting

X_train
X_test
y_train
y_test

Model Training

Use Randomized SearchCV to get the best params and train the model with this best params

- Decision-trees Classifier
- Random Forest Classifier
- Gradient Boosting

Evaluations

Choose a model with the highest Recall score on the “Dropout” status

Evaluation Metrics

Precision

Precision is the ratio of correctly predicted positive observations to the total predicted positives. High precision indicates a low false positive rate. It answers the question, "Of all the instances predicted as positive, how many were actually positive?"

Recall

Recall is the ratio of correctly predicted positive observations to all the observations in the actual class. High recall indicates a low false negative rate. It answers the question, "Of all the instances that were actually positive, how many were correctly predicted as positive?"

F1-Score

The F1-score is the harmonic mean of precision and recall. The F1-score is a balance between precision and recall. It is useful when the class distribution is imbalanced. It answers the question, "How well does the model balance precision and recall?"

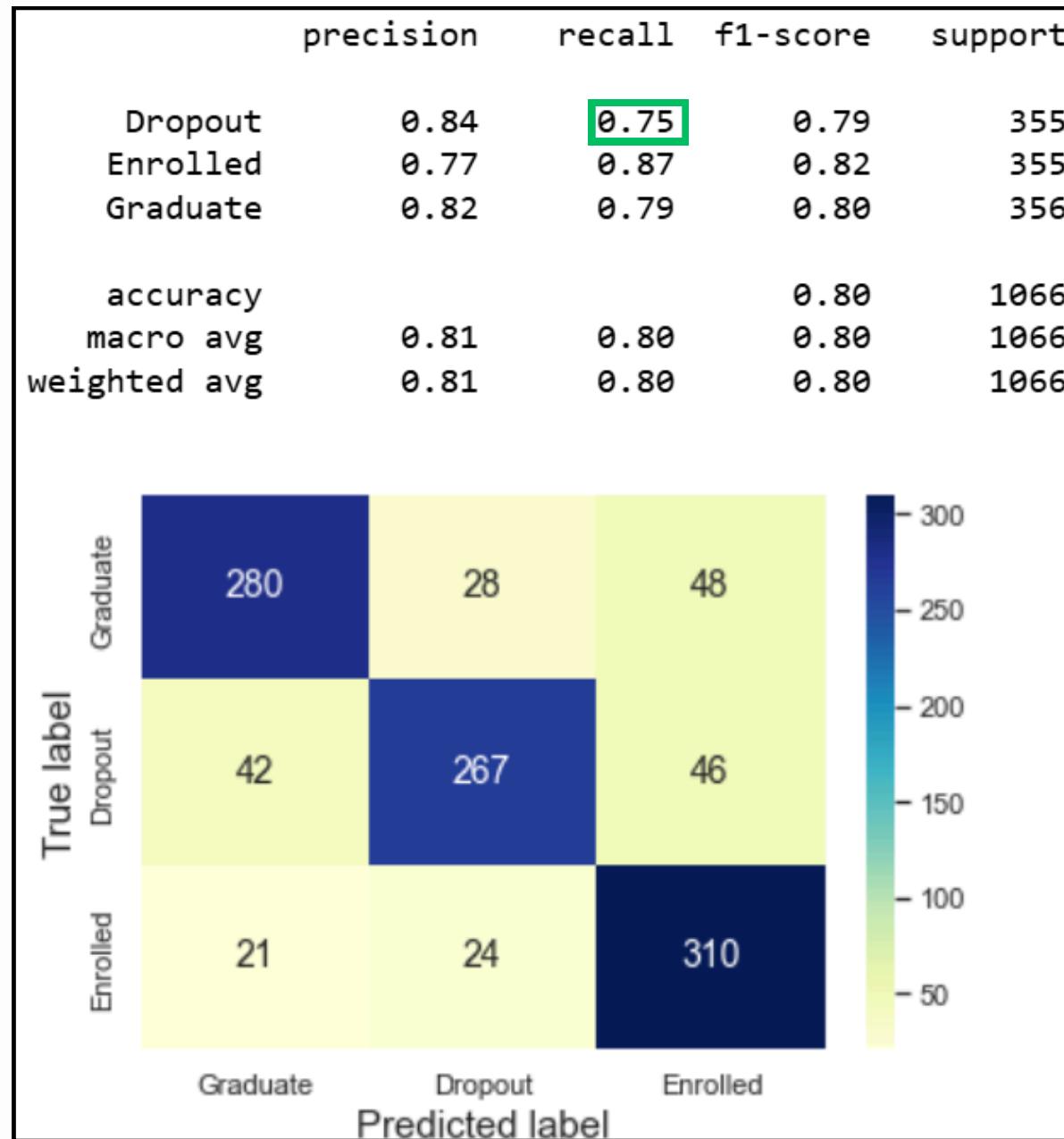
Accuracy

Accuracy is the ratio of correctly predicted observations to the total observations. It is the most intuitive performance measure and represents the overall effectiveness of the model. It answers the question, "How often is the model correct?"

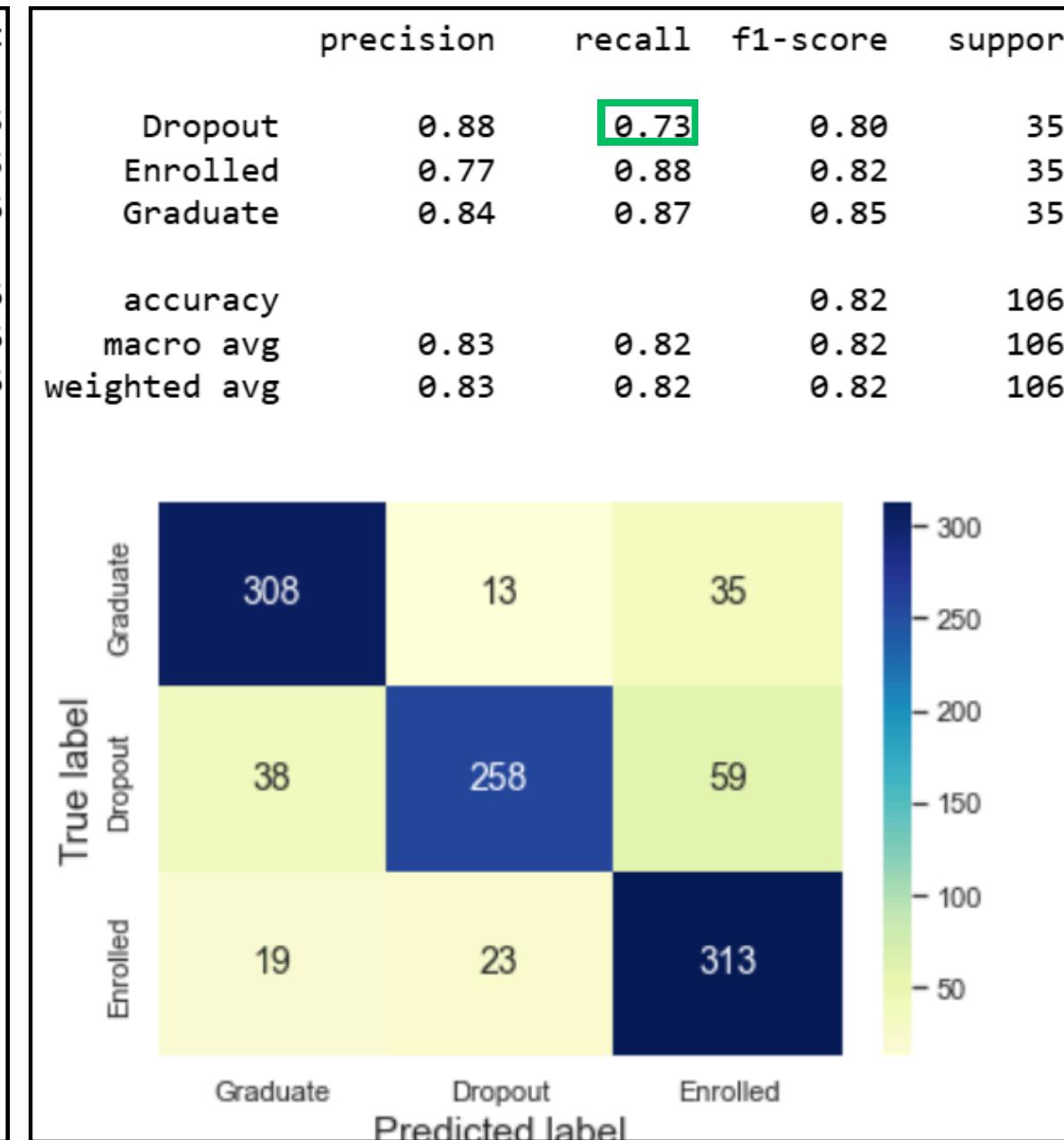
Among these evaluation metrics, recall will be used to evaluate the model's performance. Why recall? We want a model that could accurately predict the "Dropout" status so that the institute could give guidance and direction for students classified as dropout

Evaluations

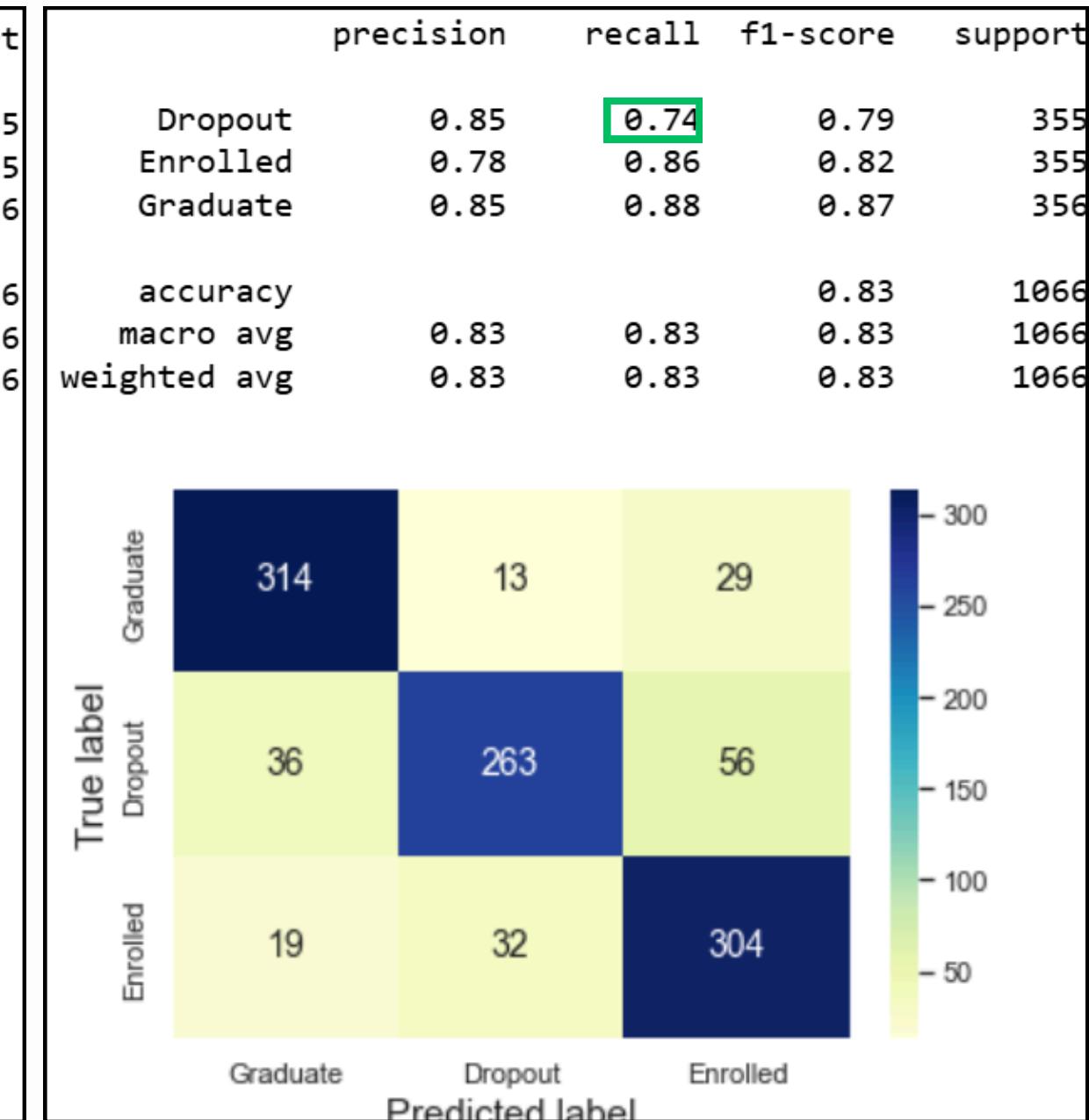
Decision-trees Classifier



Random Forest Classifier

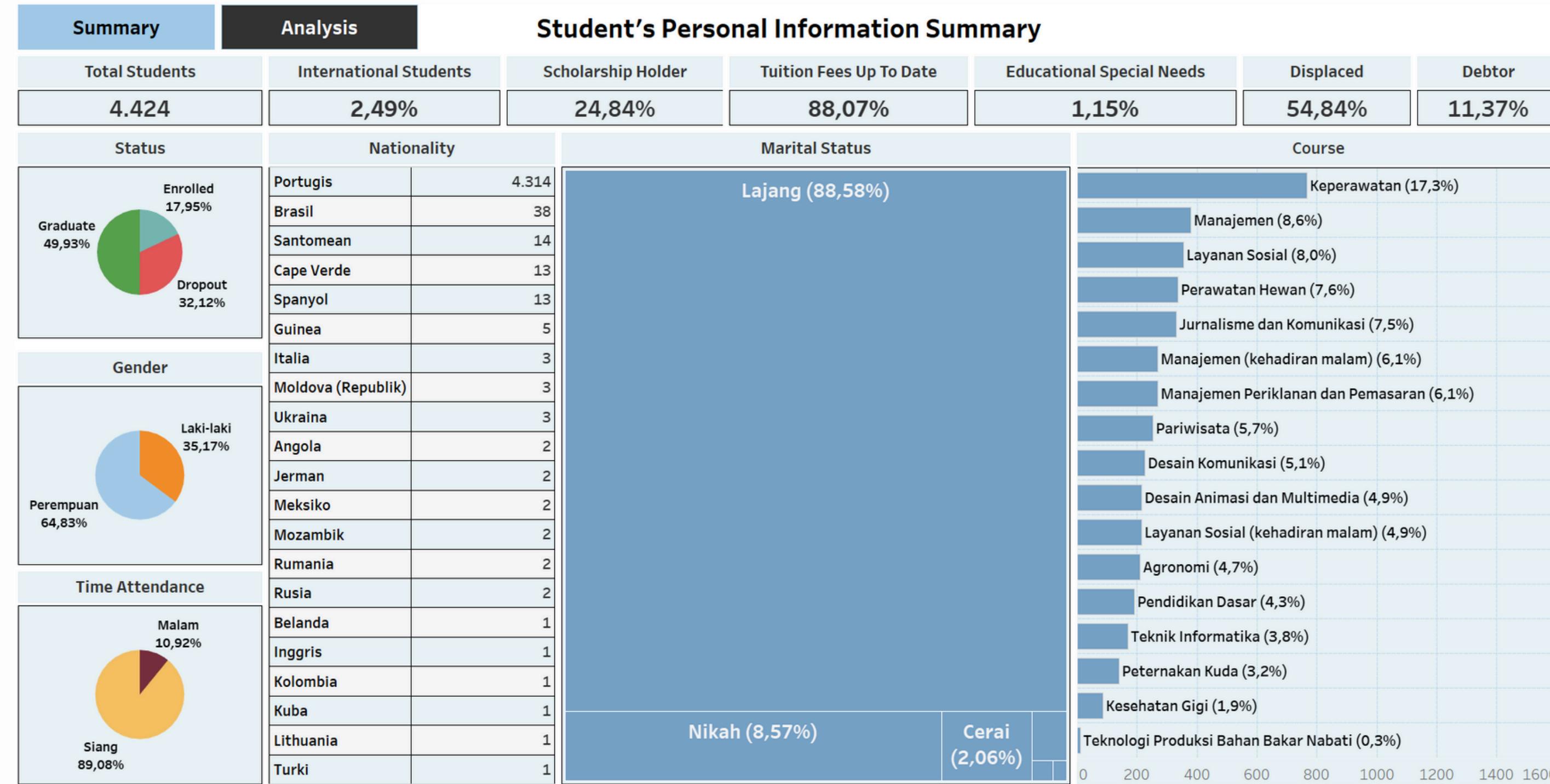


Gradient Boosting



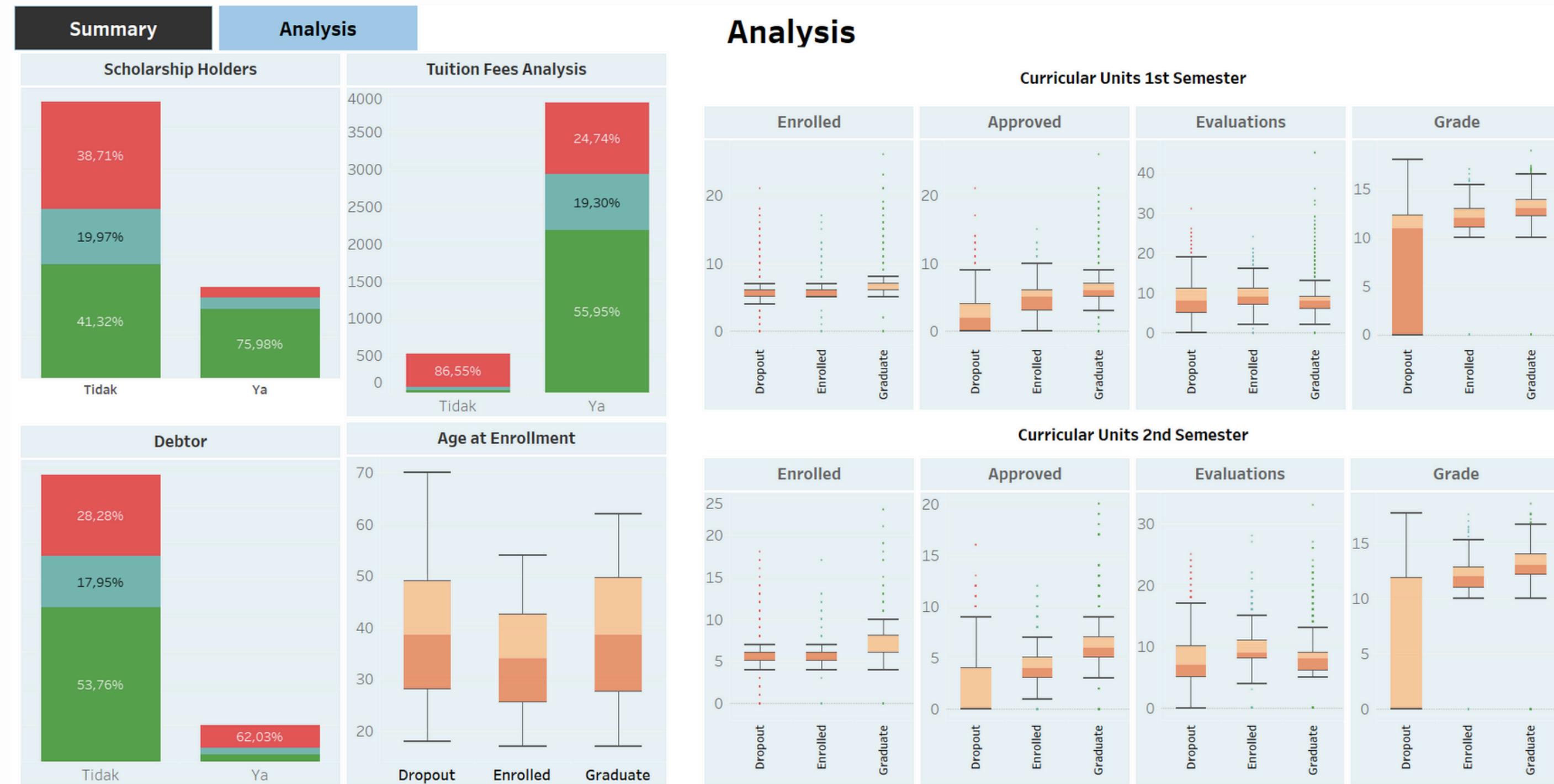
Based on the previous explanation, we need a model with the highest recall score on the “Dropout” status. From the pictures above, a model with the highest recall score on the “Dropout” status is **Decision-trees Classifier** model. This model will be deployed on Streamlit to predict the student’s status

Dashboard



*This is an interactive dashboard. The information will automatically filtered if you clicked one or more of the charts (pie charts, table, treemaps, bar chart).

Dashboard



Model Prototype

Share    

Jaya Jaya Institute

Personal Information

Gender: Laki-laki | Age at Enrollment: 22

Debtor: Tidak | Scholarship Holder: Tidak | Tuition Fees Up To Date: Tidak

Curricular Units 1st Semester Information

Enrolled (0 - 30): 20 | Evaluations (0 - 50): 35 | Approved (0 - 30): 20 | Grade (0 - 20): 17,00

Curricular Units 2nd Semester Information

Enrolled (0 - 30): 20 | Evaluations (0 - 50): 35 | Approved (0 - 30): 20 | Grade (0 - 20): 17,00

[Manage app](#)

Model Prototype

[Share](#)

Curricular Units 2nd Semester Information

Enrolled (0 -30)

20

-

+

Evaluations (0 - 50)

35

-

+

Approved (0 - 30)

20

-

+

Grade (0 - 20)

17,00

-

+

Overall Information

	Gender	Age_at_enrollment	Debtors	Scholarship_holder	Tuition_fees_up_to_date	Curricular_
0	Laki-laki	22	Tidak	Tidak	Tidak	

Predict

View the preprocessed data

	Transformed_Age_at_enrollment	Transformed_pca1	Transformed_pca2	Transformed_pca3	Trar
0	-1,709.6552	0.6853	4.5915	1.8819	

Unfortunately, you are Dropouted

Manage app

Conclusions

The number of dropped out students in this higher education reaches 1421 out of 4424 students (32%). Indeed that was a high number. After analysis, some factors affect the student's status, whether they will dropout or not. The strongest factor that affects the student will dropout is **tuition fees up to date** feature. Students who dropped out because the tuition fees un-updated reaches 86%. Next, there are some interesting patterns obtained by students in the 1st and 2nd semester. Thee graduated students tend to have **higher grade** rather than enrolled and dropped out students.

Recommendation Action Items

Keep tuition fees updated to the recent info

An up-to-date tuition fee will ensure that students do not experience **unexpected financial problems**, which may cause them to drop out. A possible action is to conduct **regular audits** every semester to ensure that all tuition fees data has been updated. The finance team can send **automated reminders** to students and parents about tuition fees updates.

Pay more attention to students who get low grades (below 13)

Actions that can be given to students with these criteria are such as providing **academic guidance**, encouraging students to **consult with supervisors** regarding problems or obstacles faced by students when studying, and providing **remedial classes** for students with low scores.

THANK YOU!



Norman Dwi Febrio

Prospective Data Scientist



[linkedIn.com/in/norman-dwi-febrio](https://www.linkedin.com/in/norman-dwi-febrio)



github.com/NormanFebrio