

# SC1015 MINI PROJECT

Predicting how well  
students can cope in  
university.

FDAA Group 8

Yeo Ruizhi Carwyn  
Kent Karsten Pangestu  
James Ang Wen Jie





# Table of Contents

1

**Problem**  
Dataset  
Problem Definition

2

**Data**  
Preparation  
& Cleaning

3

**Exploratory Data  
Analysis**  
Visualisation

4

**Machine Learning  
Models**

5

**Conclusion**  
Insights  
Reflection



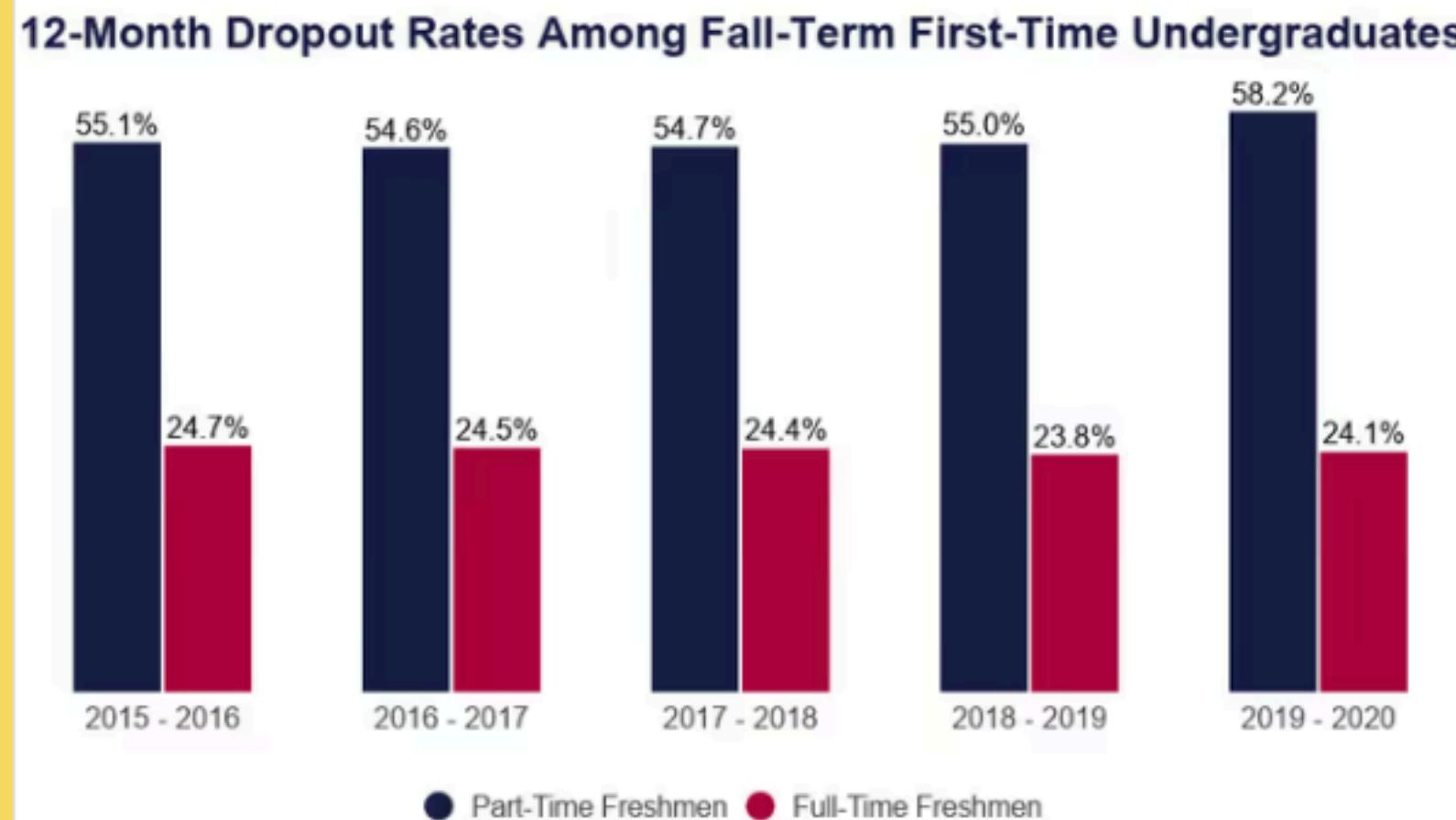


# Problem





# Problem



## College Dropout Rate [2023]: by Year + Demographics

Find data for American college students' dropout and graduation rates, including demographic averages and historical trends.

 Education Data Initiative

## Students cannot cope in University

College dropout rates indicate that up to 32.9% of undergraduates do not complete their degree program.

- First-time undergraduate first-year students have a 12-month dropout rate of 24.1%.
- Among first-time bachelor's degree seekers, 25.7% ultimately drop out; among all undergraduate students, up to 40% drop out.



# Problem Statement



Dropout Prediction: Can we classify students into those likely to drop out vs. those who will continue their education?

Academic Success Prediction: Can we predict the GPA, graduation likelihood, or other academic success indicators at the end of their course or semester?





# Dataset



**UCI Machine Learning Repository**  
Discover datasets around the world!  
[ics.uci.edu](http://ics.uci.edu)

“ Predict Students' Dropout and Academic Success”,  
Based on a paper by By Mónica V. Martins,  
Daniel Tolledo, Jorge Machado, Luís M. T.  
Baptista, and Valentim Realinho. 2021.



# DATA CLEANING

## VARIABLES REMOVED

'MARITAL STATUS',  
'APPLICATION MODE',  
'CURRICULAR UNITS 1ST SEM (CREDITED)',  
'CURRICULAR UNITS 1ST SEM (ENROLLED)',  
    'CURRICULAR UNITS 1ST SEM  
        (EVALUATIONS)',  
    'CURRICULAR UNITS 1ST SEM (WITHOUT  
        EVALUATIONS)',  
'CURRICULAR UNITS 2ND SEM (CREDITED)',  
'CURRICULAR UNITS 2ND SEM (ENROLLED)',  
    'CURRICULAR UNITS 2ND SEM  
        (EVALUATIONS)',  
    'CURRICULAR UNITS 2ND SEM (WITHOUT  
        EVALUATIONS)',  
    'UNEMPLOYMENT RATE',  
    'INFLATION RATE',  
    'GDP',  
    'DISPLACED',  
    'EDUCATIONAL SPECIAL NEEDS',  
    'DAYTIME/EVENING ATTENDANCE '

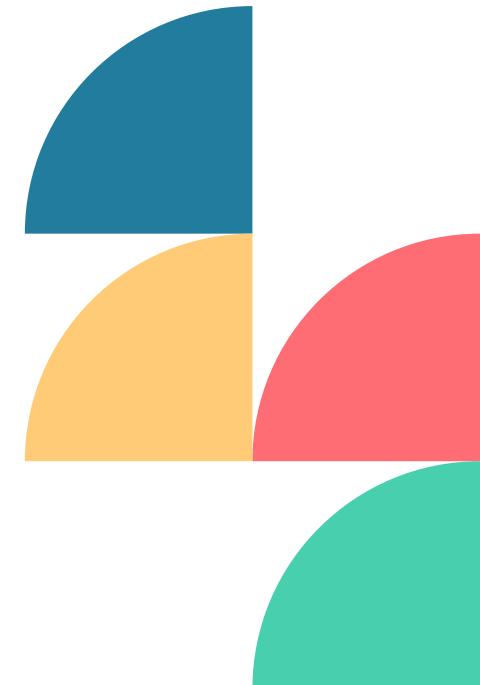
## VARIABLES REMOVED

```
1 # Variables we removed because we think they are irrelevant to our problem:  
2 Columns_to_remove = ['Marital status', 'Application mode',  
3 'Curricular units 1st sem (credited)', 'Curricular units 1st sem (enrolled)', 'Curricul  
4 'Curricular units 1st sem (without evaluations)', 'Curricular units 2nd sem (credited)',  
5 'Curricular units 2nd sem (evaluations)', 'Curricular units 2nd sem (without evaluations  
6 'GDP', 'Displaced', 'Educational special needs', 'Daytime/evening attendance      ']  
7  
8 data.drop(Columns_to_remove, axis=1, inplace=True)  
9  
10 # Streamline the 'Nationality' variable to only local students (i.e. Nationality == 1)  
11 data = data.loc[data['Nacionality'] == 1].copy()  
12 data.drop(columns=['Nacionality', 'International'], axis=1, inplace=True)
```

## VARIABLES RECATEGORISED

'TARGET'  
'COURSE'  
"MOTHER'S OCCUPATION"  
"FATHER'S OCCUPATION"  
"MOTHER'S QUALIFICATION"  
"FATHER'S QUALIFICATION"  
"PREVIOUS QUALIFICATION"

```
50 # We categorised the level of educations from 0 to 5, where:
51 def clean_Education(code):
52     if code in [34, 35, 36]:
53         return 0 #'Uneducated'
54     elif code in [37, 38]:
55         return 1 #'Primary Education'
56     elif code in [1, 10, 11, 12, 14, 18, 19, 26, 27, 29, 30, 39]:
57         return 2 #'Secondary/High School'
58     elif code in [22]:
59         return 3 #'Tertiary Education'
60     elif code in [2, 40]:
61         return 4 #'Bachelor's'
62     elif code in [3, 4, 5, 41, 42, 43, 44]:
63         return 5 #'Above Bachelor's'
64     else:
65         return -1 #'Others'
66
67 def clean_previous_qualification(code):
68     if code in [1, 9, 10, 12, 14, 15, 19, 38]:
69         return 2 #'Secondary/High School'
70     elif code in [39]:
71         return 3 #'Tertiary Education'
72     elif code in [2, 40]:
73         return 4 #'Bachelor's'
74     elif code in [3, 4, 5, 6, 42, 43]:
75         return 5 #'Above Bachelor's'
76     else:
77         return -1 #'Others'
```



```
<class 'pandas.core.frame.DataFrame'>
```

```
Index: 4157 entries, 0 to 4422
```

```
Data columns (total 19 columns):
```

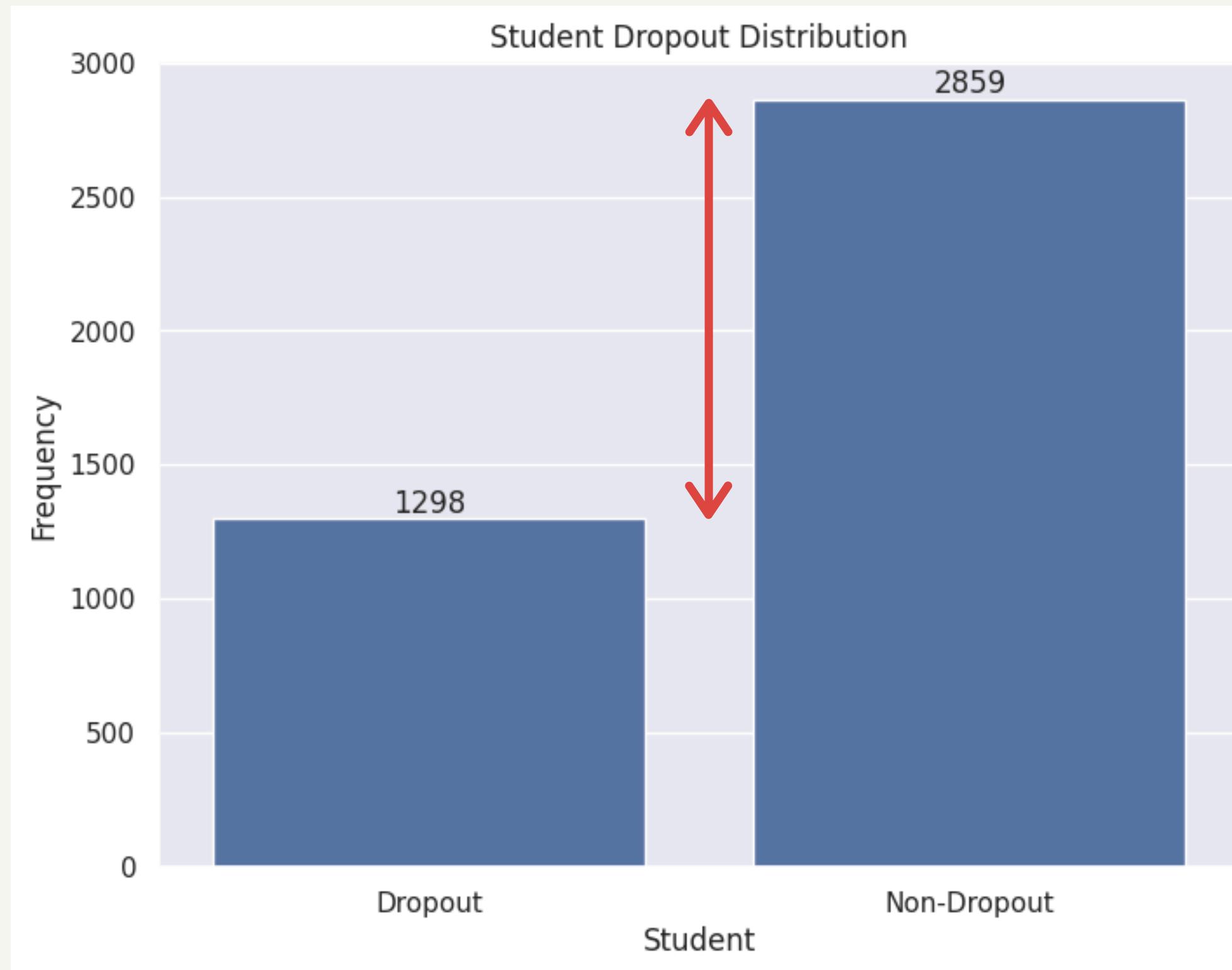
#	Column	Non-Null Count	Dtype
0	Application order	4157 non-null	int64
1	Course	4157 non-null	object
2	Previous qualification	4157 non-null	int64
3	Previous qualification (grade)	4157 non-null	float64
4	Mother's qualification	4157 non-null	int64
5	Father's qualification	4157 non-null	int64
6	Mother's occupation	4157 non-null	object
7	Father's occupation	4157 non-null	object
8	Admission grade	4157 non-null	float64
9	Debtors	4157 non-null	int64
10	Tuition fees up to date	4157 non-null	int64
11	Gender	4157 non-null	int64
12	Scholarship holder	4157 non-null	int64
13	Age at enrollment	4157 non-null	int64
14	Curricular units 1st sem (approved)	4157 non-null	int64
15	Curricular units 1st sem (grade)	4157 non-null	float64
16	Curricular units 2nd sem (approved)	4157 non-null	int64
17	Curricular units 2nd sem (grade)	4157 non-null	float64
18	Target	4157 non-null	object

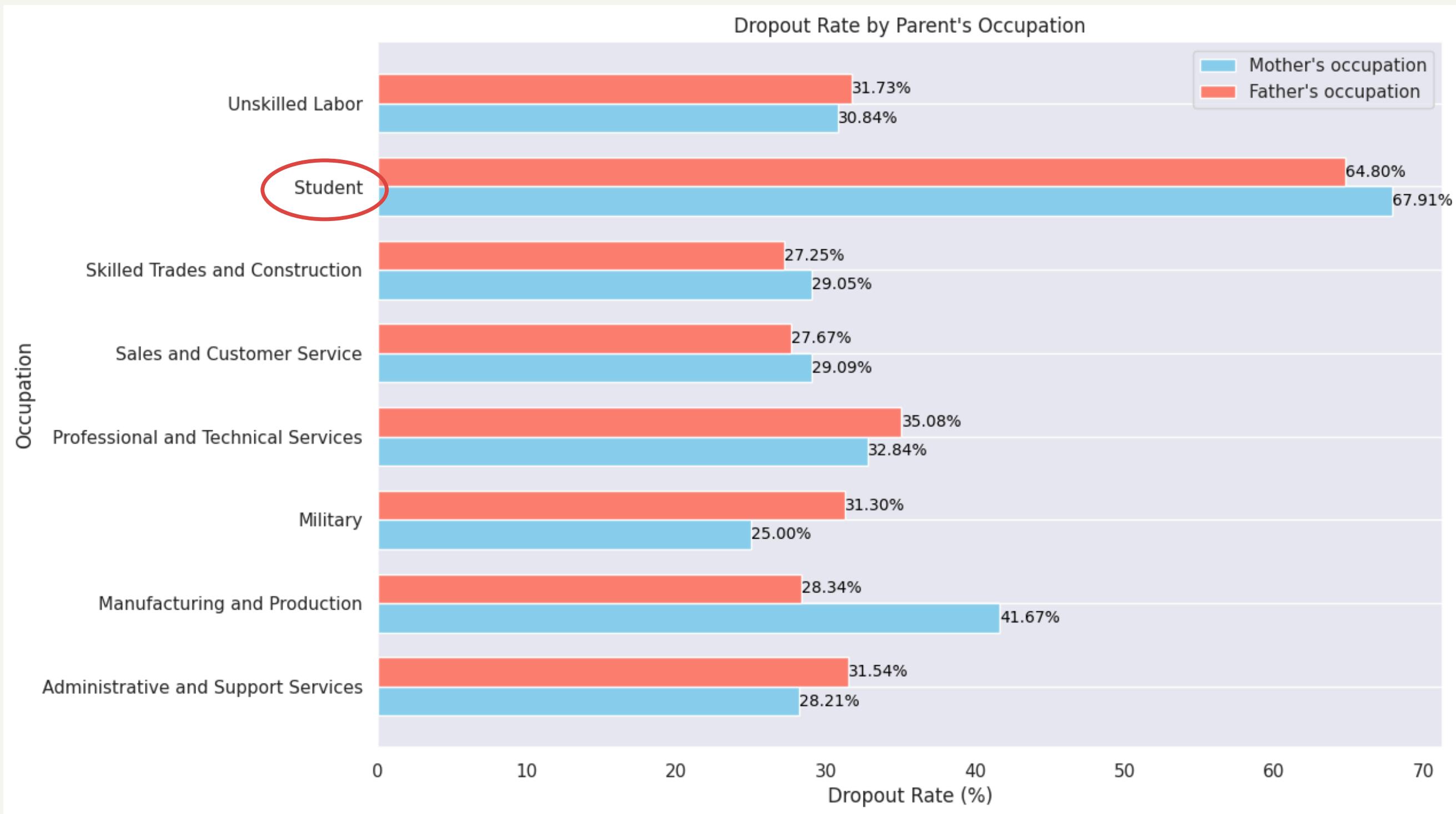
```
dtypes: float64(4), int64(11), object(4)
```

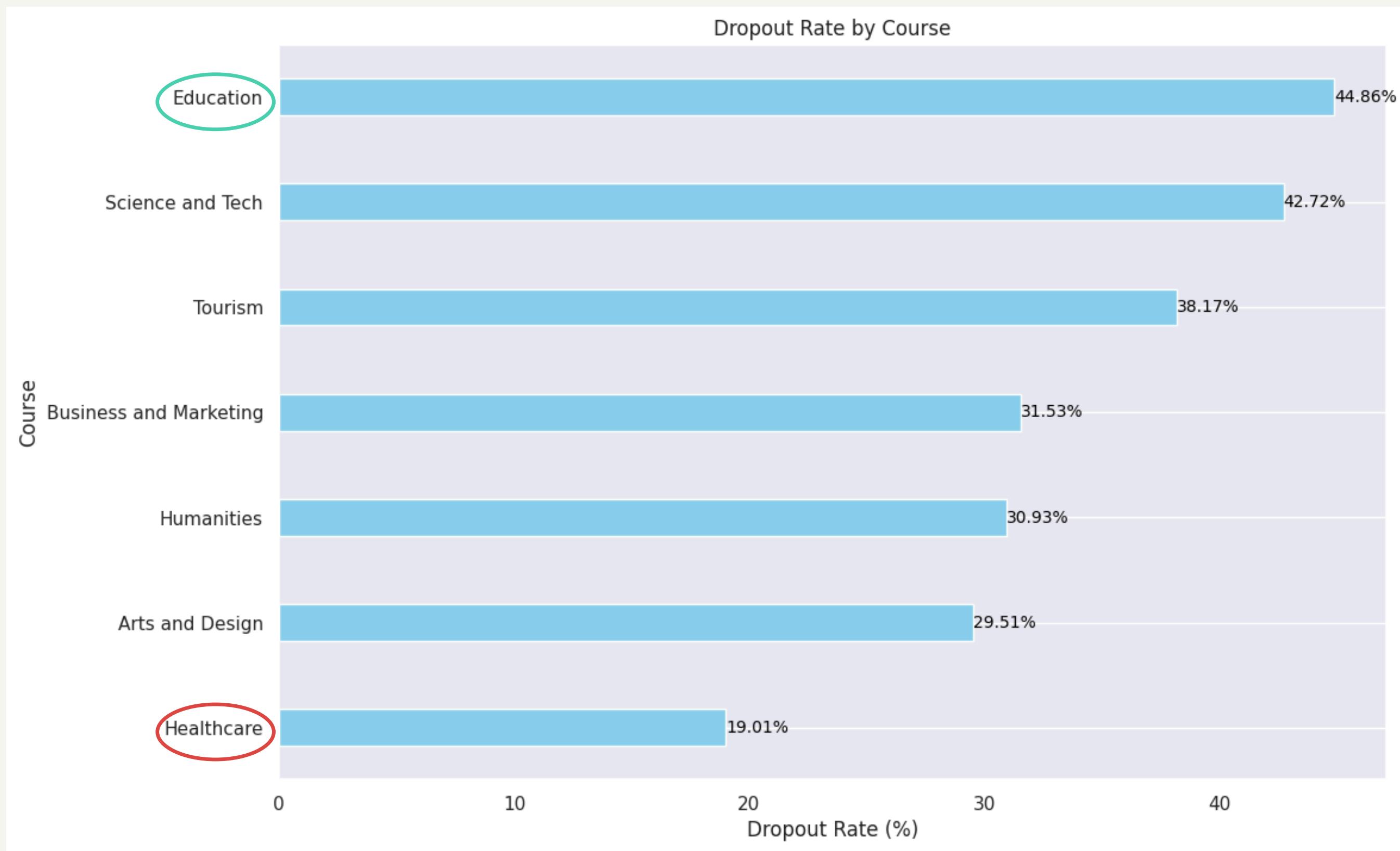
```
memory usage: 649.5+ KB
```

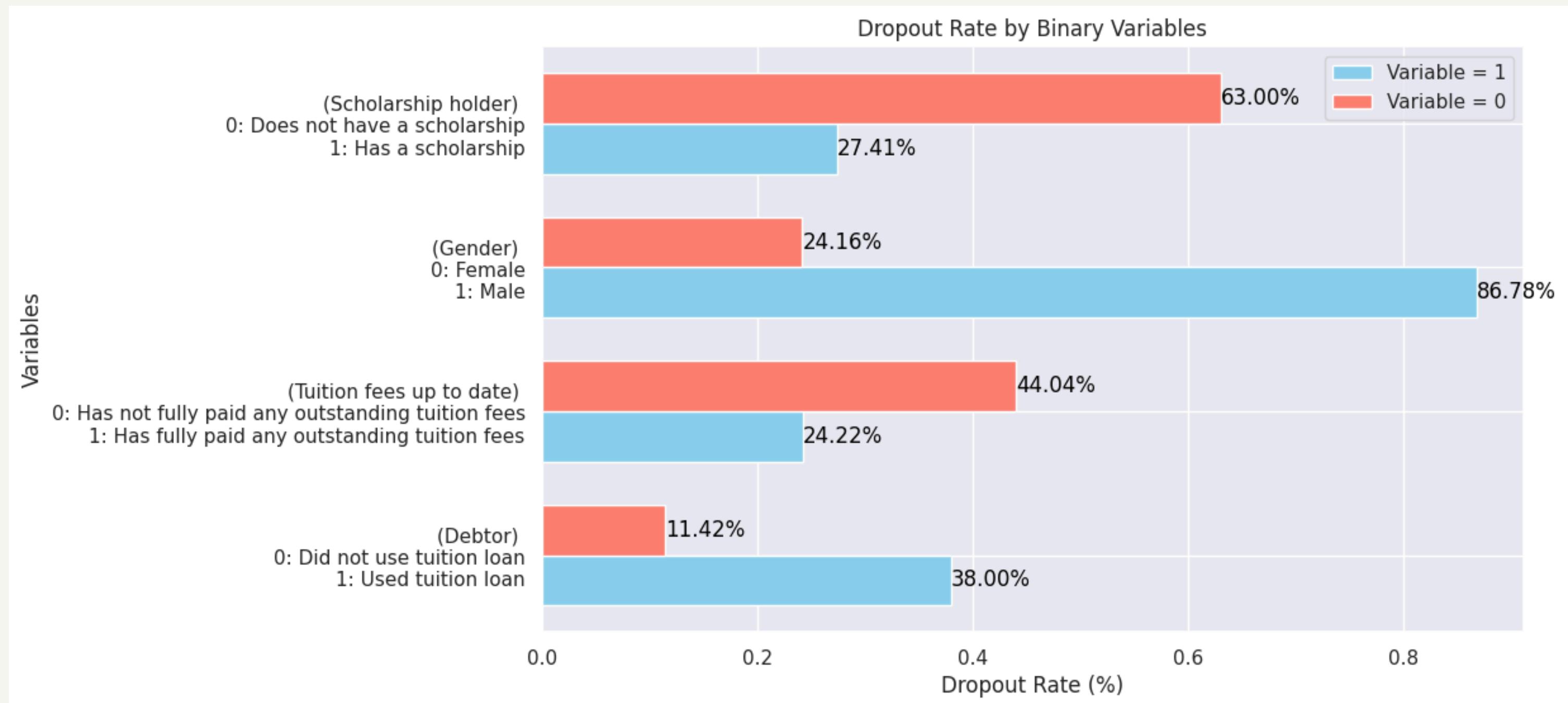
# **EXPLORATORY DATA ANALYSIS**

# Dropout vs Non Dropout







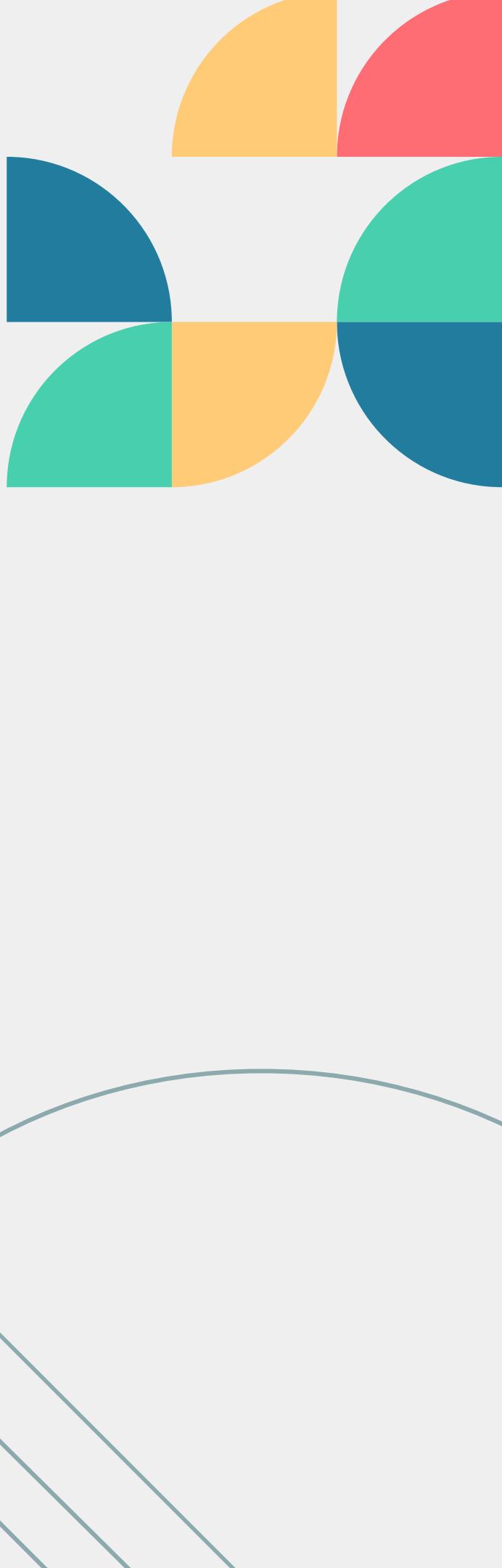
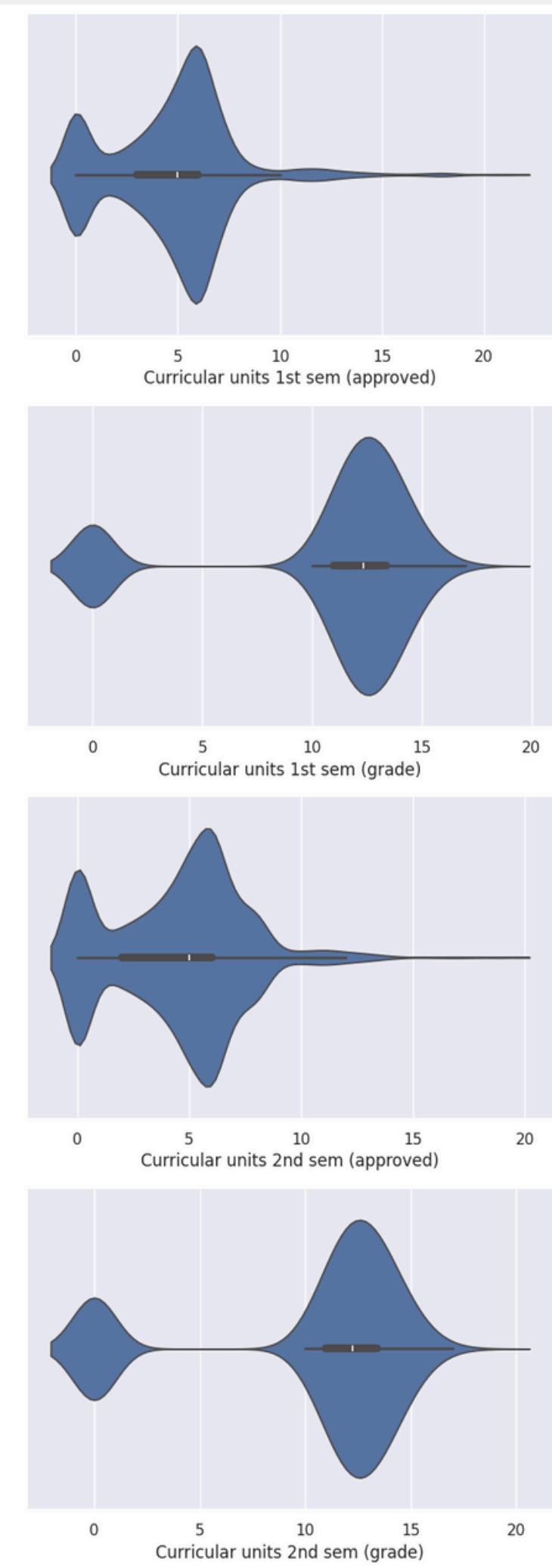
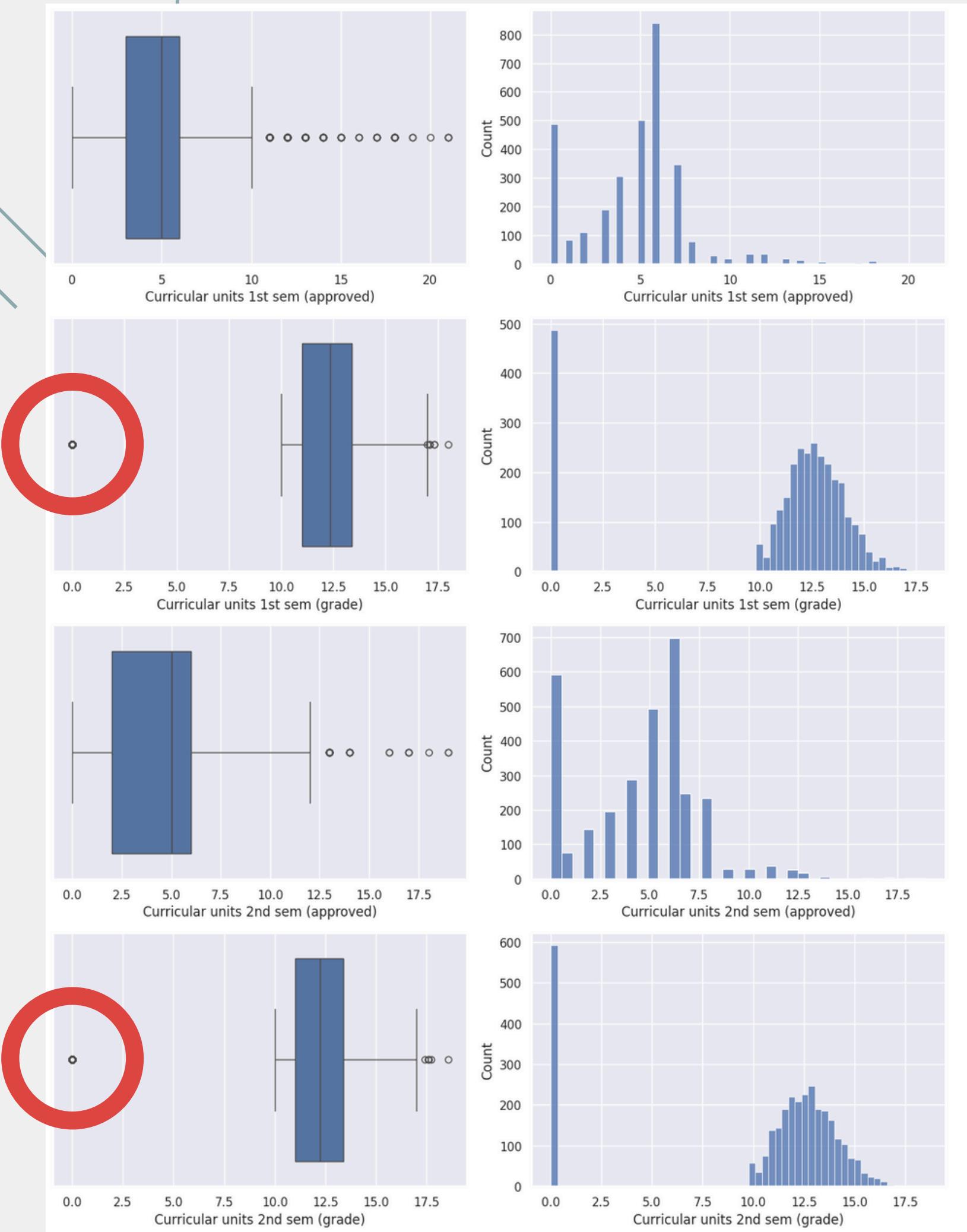


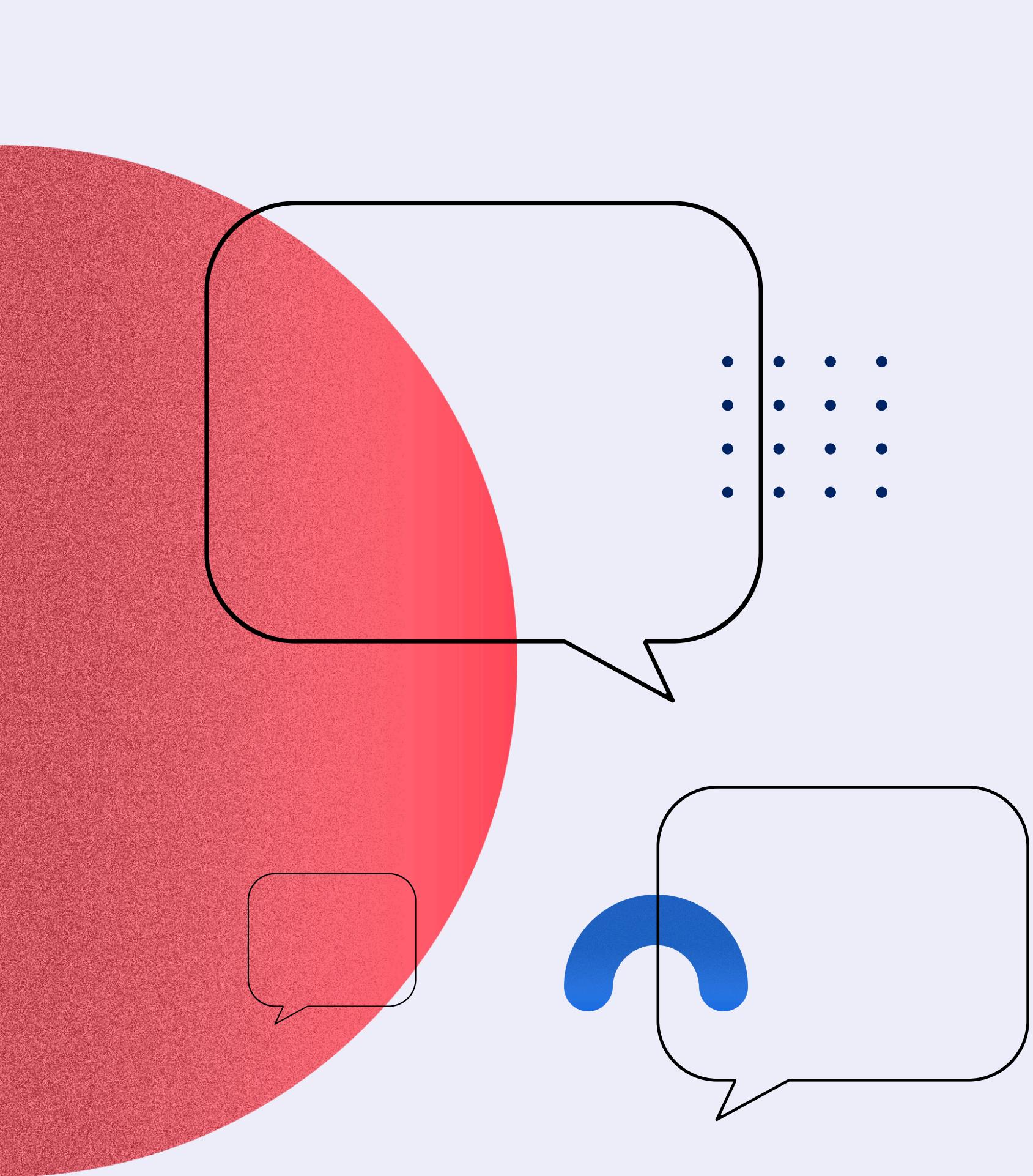
Correlation Matrix for Selected Variables

	Previous qualification (grade)	Admission grade	Age at enrollment	Curricular units 1st sem (approved)	Curricular units 1st sem (grade)	Curricular units 2nd sem (approved)	Curricular units 2nd sem (grade)	Target
Previous qualification (grade)	1	0.58	-0.11	0.047	0.056	0.049	0.05	-0.076
Admission grade	0.58	1	-0.0095	0.065	0.07	0.07	0.068	-0.086
Age at enrollment	-0.11	-0.0095	1	-0.037	-0.14	-0.095	-0.16	0.24
Curricular units 1st sem (approved)	0.047	0.065	-0.037	1	0.69	0.9	0.68	-0.47
Curricular units 1st sem (grade)	0.056	0.07	-0.14	0.69	1	0.67	0.84	-0.47
Curricular units 2nd sem (approved)	0.049	0.07	-0.095	0.9	0.67	1	0.76	-0.56
Curricular units 2nd sem (grade)	0.05	0.068	-0.16	0.68	0.84	0.76	1	-0.56
Target	-0.076	-0.086	0.24	-0.47	-0.47	-0.56	-0.56	1



Previous qualification (grade)  
Admission grade  
Age at enrollment  
Curricular units 1st sem (approved)  
Curricular units 1st sem (grade)  
Curricular units 2nd sem (approved)  
Curricular units 2nd sem (grade)  
Target





# MACHINE LEARNING

# LINEAR REGRESSION

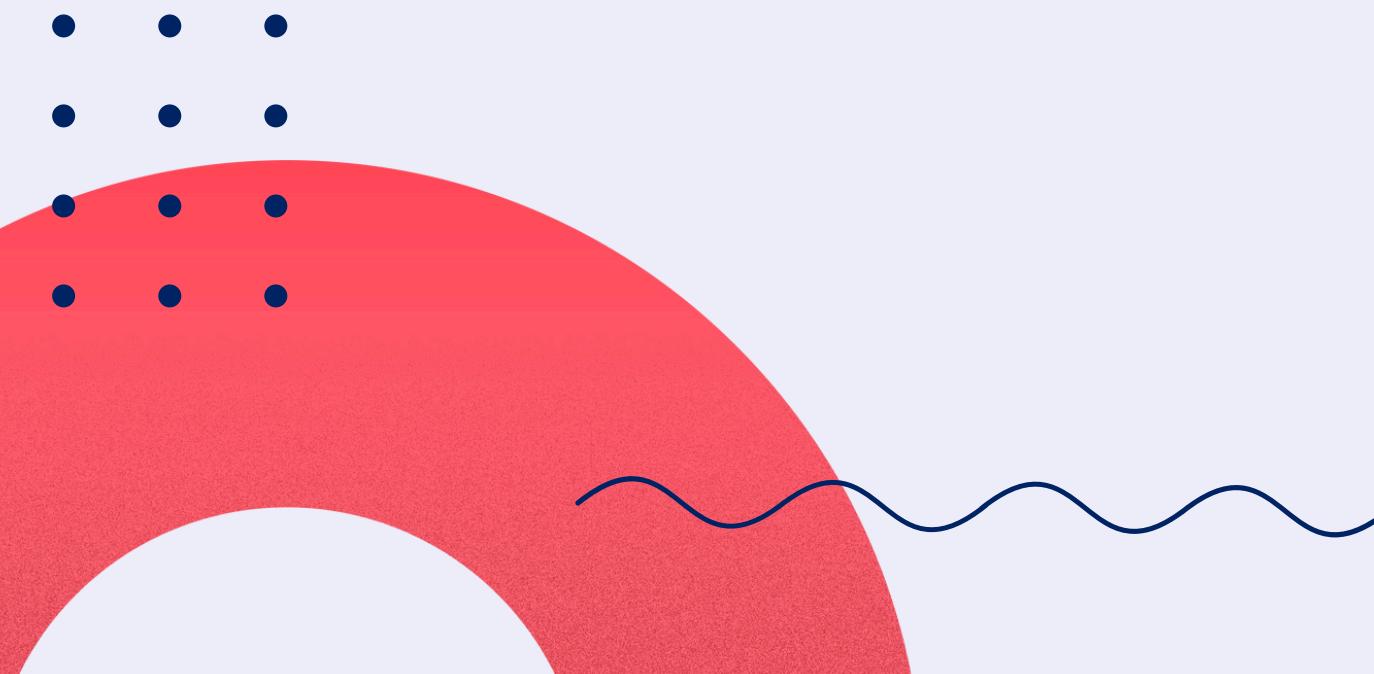
The Response Variables will be the Curricular units  
1st sem (grade) and Curricular units 2nd sem (grade)

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (\hat{Y}_i - Y_i)^2$$

Mean      Error      Squared

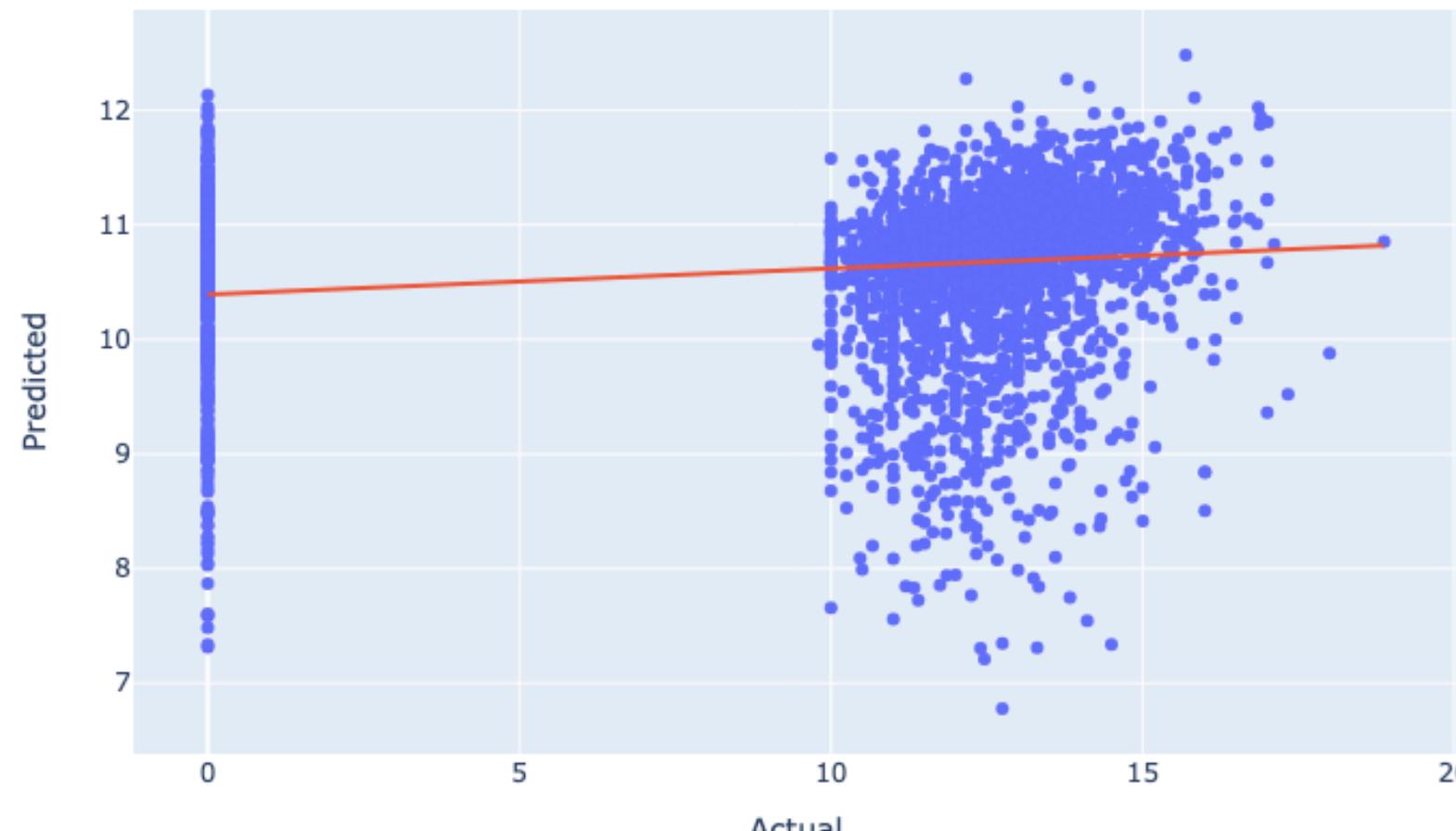
The Predictor Variables will be the variables that have floating point values:  
Previous qualification, Age at enrollment, Admission grade, Previous qualification (grade)

We will be using MSE as a metric for our regression models.

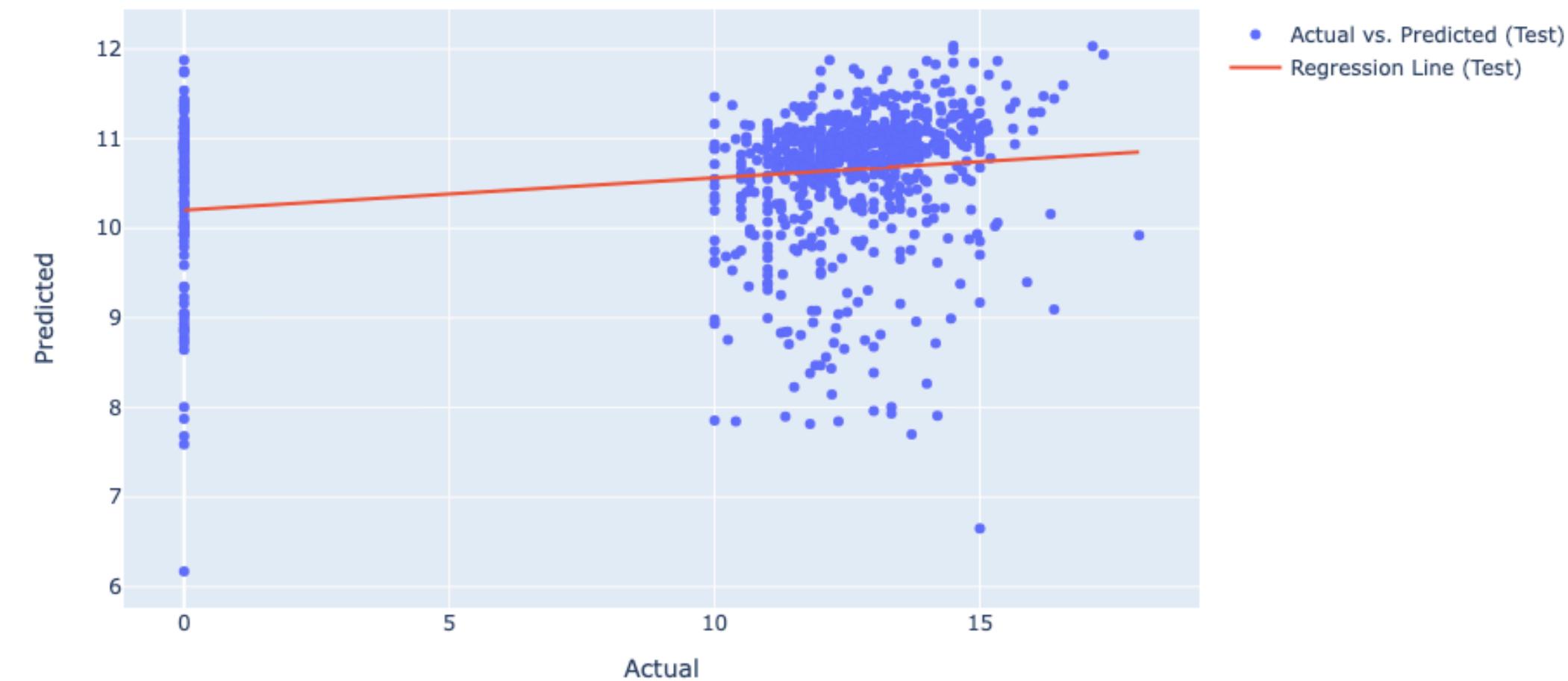


# LINEAR REGRESSION

Unclean Data: Actual vs Predicted Train Data



Unclean Data: Actual vs Predicted Test Data



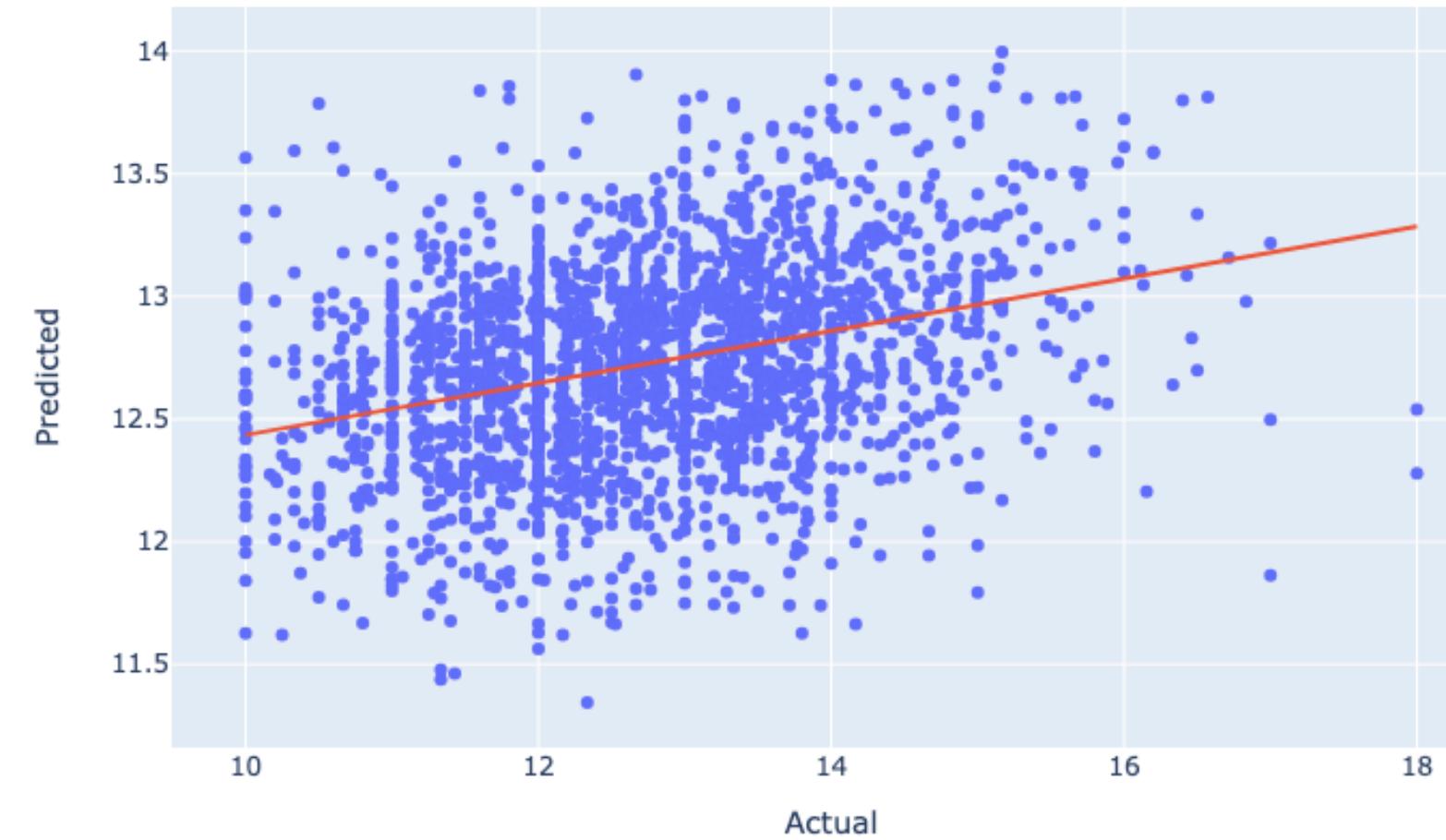
# CLEAN DATA

```
# Removing outliers
Q1 = X.quantile(0.25)
Q3 = X.quantile(0.75)
IQR = Q3 - Q1
X_clean = X[~((X < (Q1 - 1.5 * IQR)) | (X > (Q3 + 1.5 * IQR))).any(axis=1)]
y_clean = y.loc[X_clean.index]

# Remove all actual scores equal to 0 (as observed from the unclean data, it is obvious that they are outliers to be removed)
y_clean = y_clean[(y_clean["Curricular units 1st sem (grade)"] != 0) & (y_clean["Curricular units 2nd sem (grade)"] != 0)]
X_clean = X_clean.loc[y_clean.index]
```

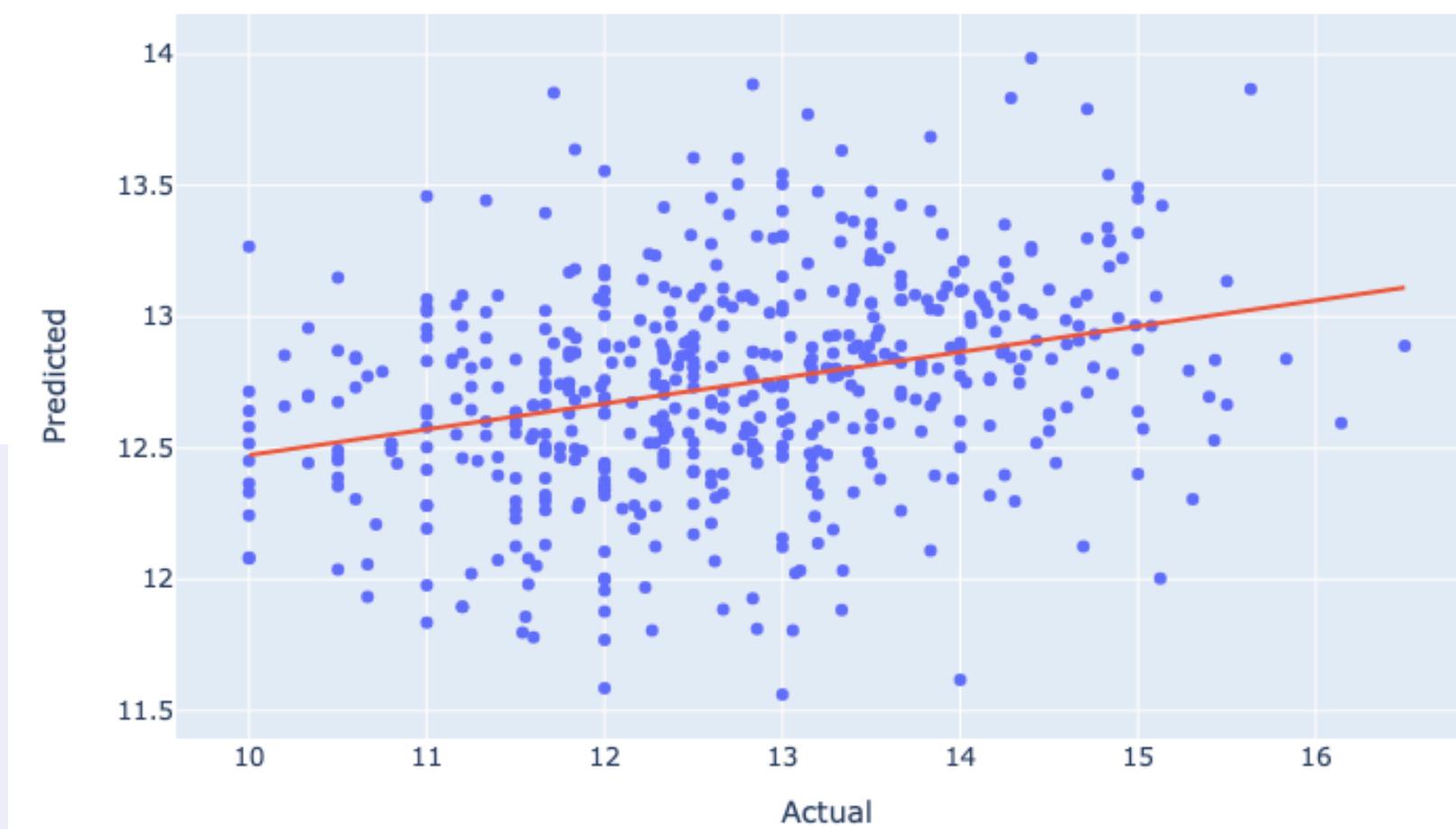
# CLEAN DATA

Clean Data: Actual vs Predicted Train Data



● Actual vs. Predicted (Train)  
— Regression Line (Train)

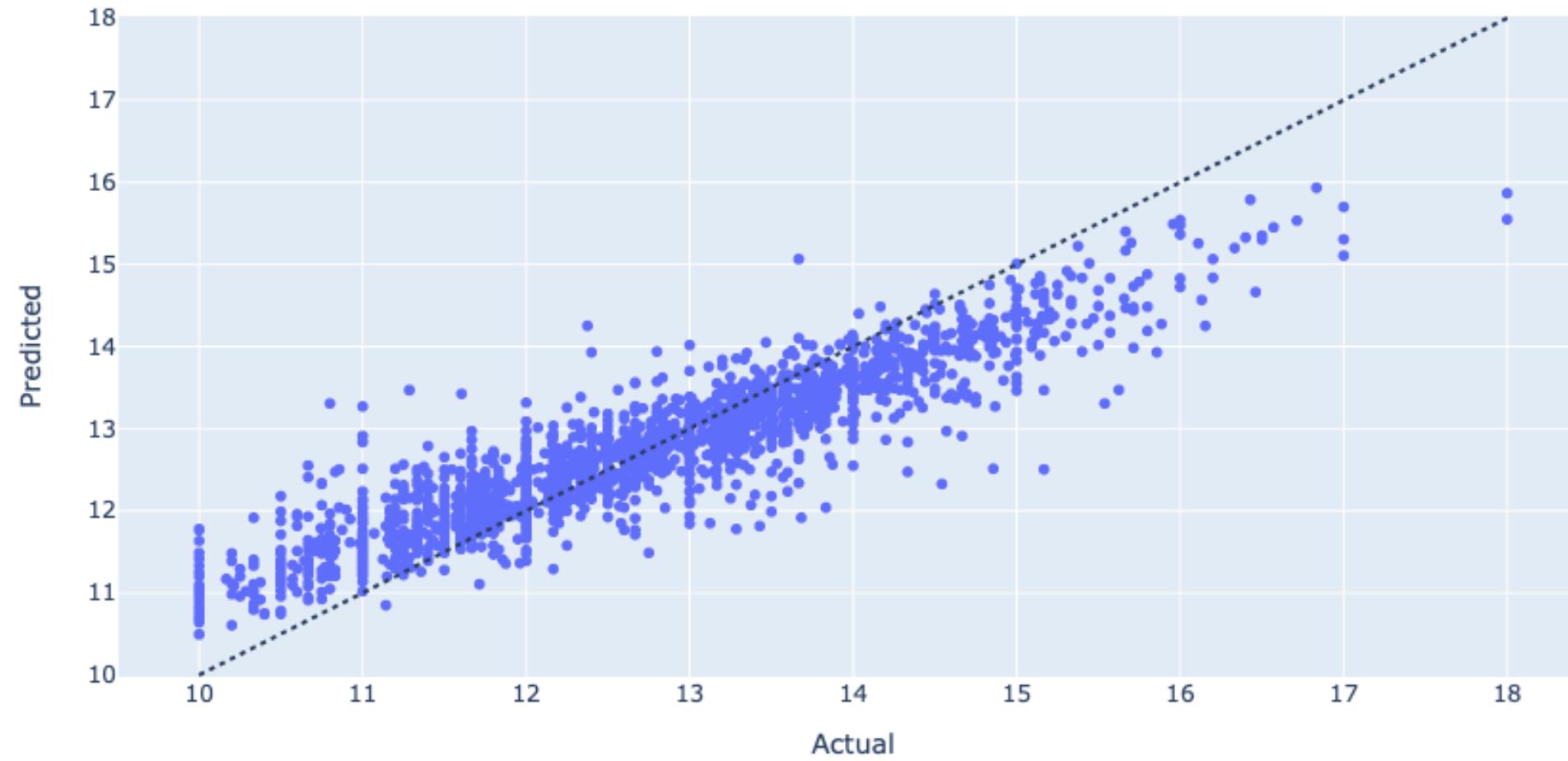
Clean Data: Actual vs Predicted Test Data



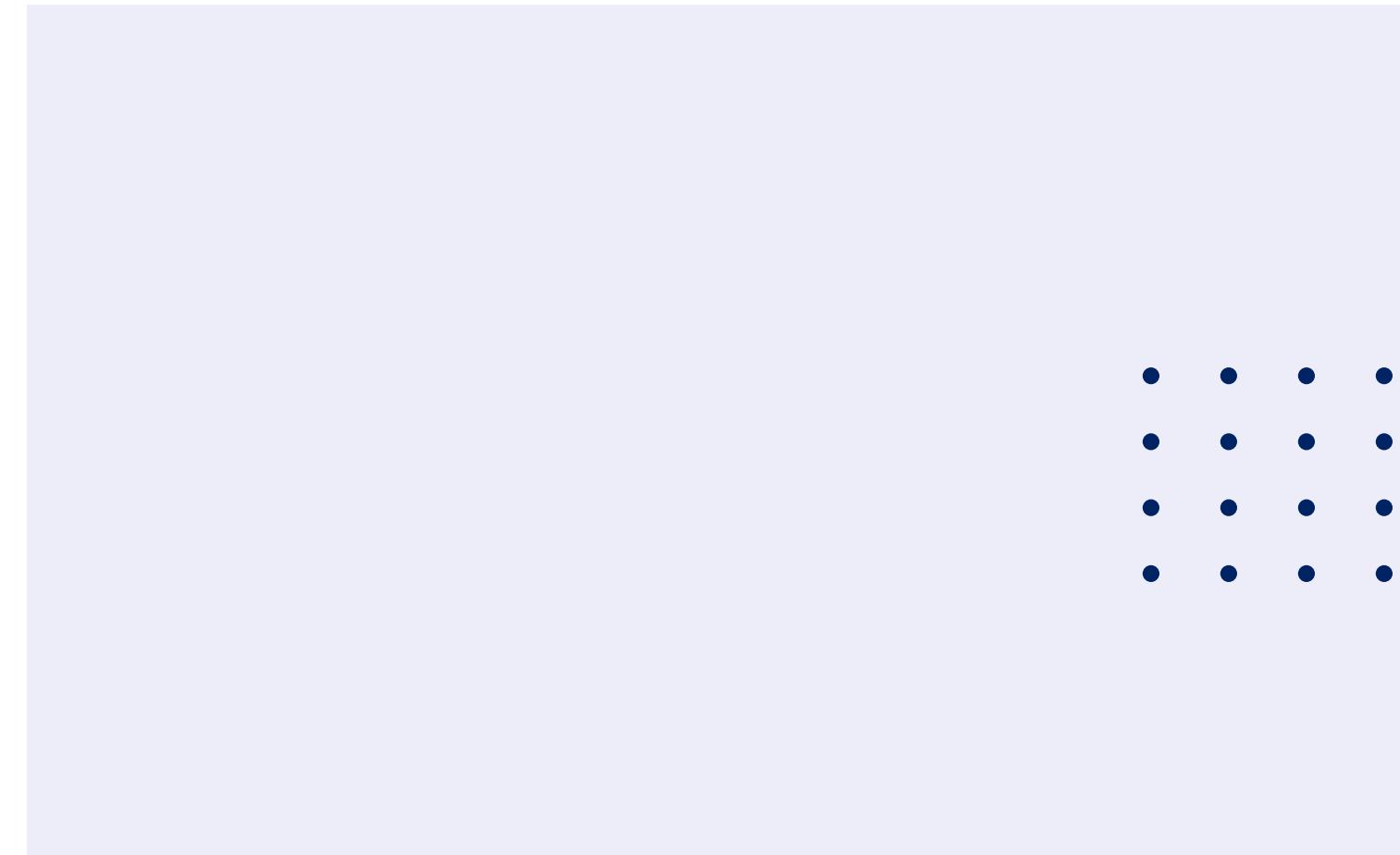
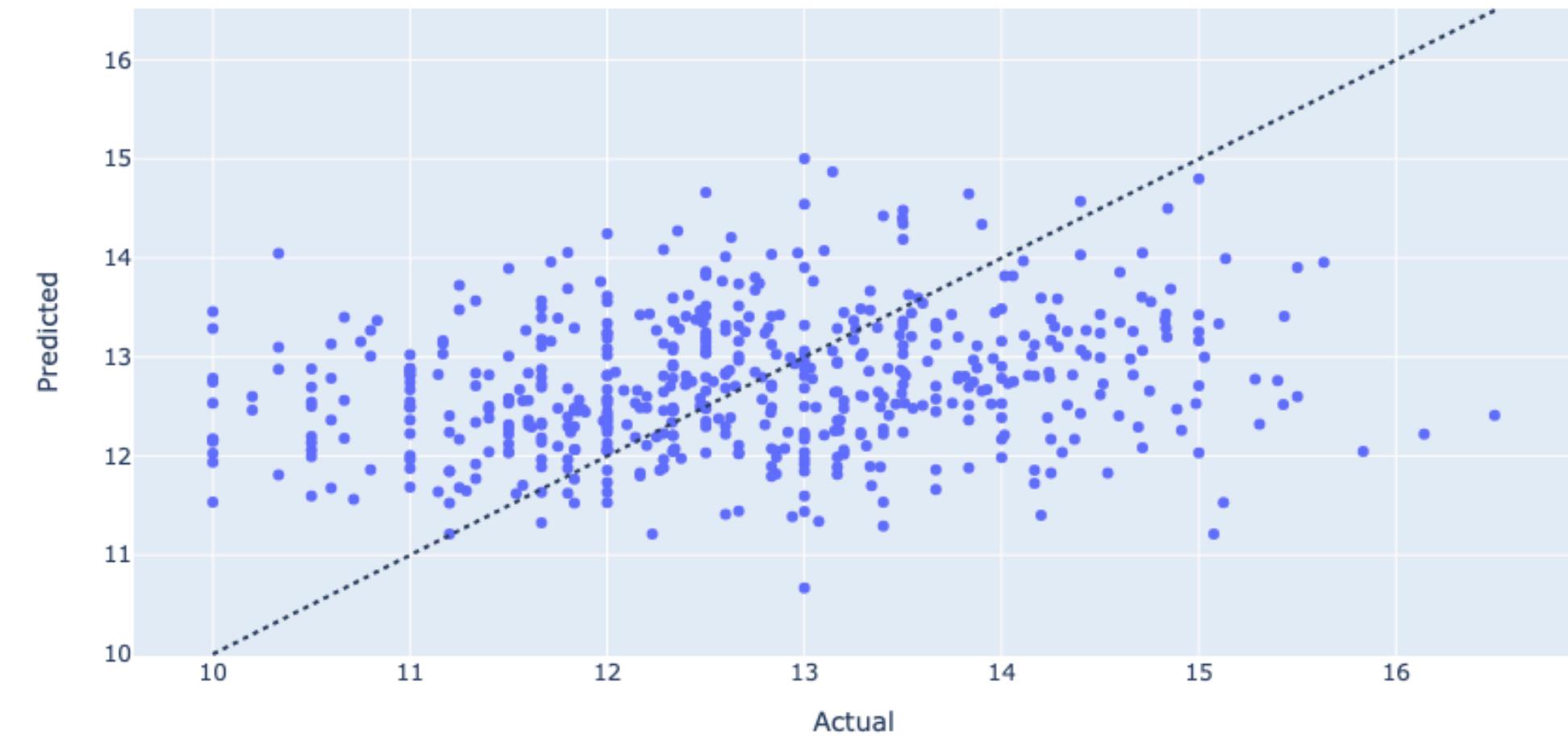
● Actual vs. Predicted (Test)  
— Regression Line (Test)

# RANDOM FOREST REGRESSOR

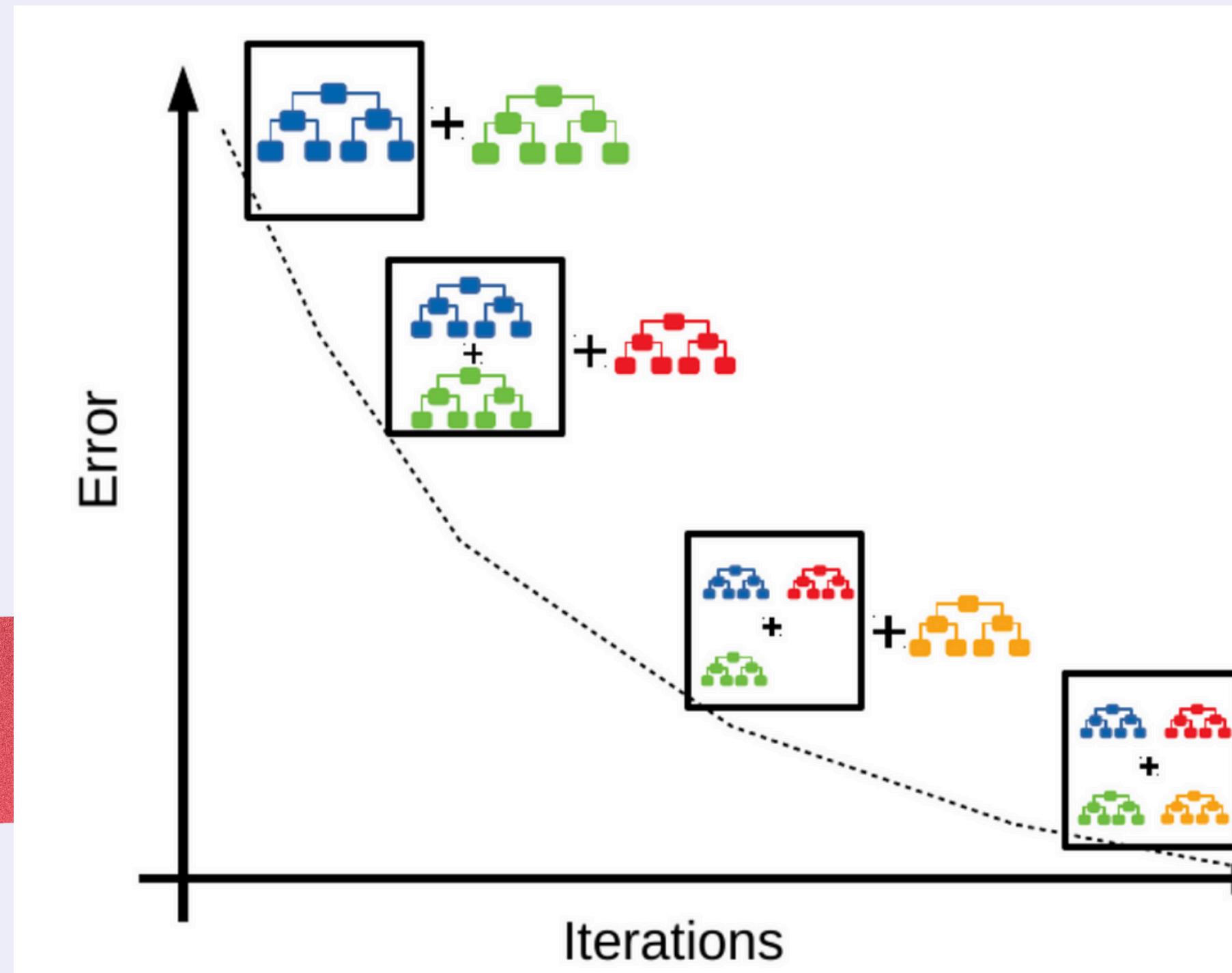
Actual vs Predicted - Training Data (Cleaned)



Actual vs Predicted - Testing Data (Cleaned)



# GRADIENT BOOSTING



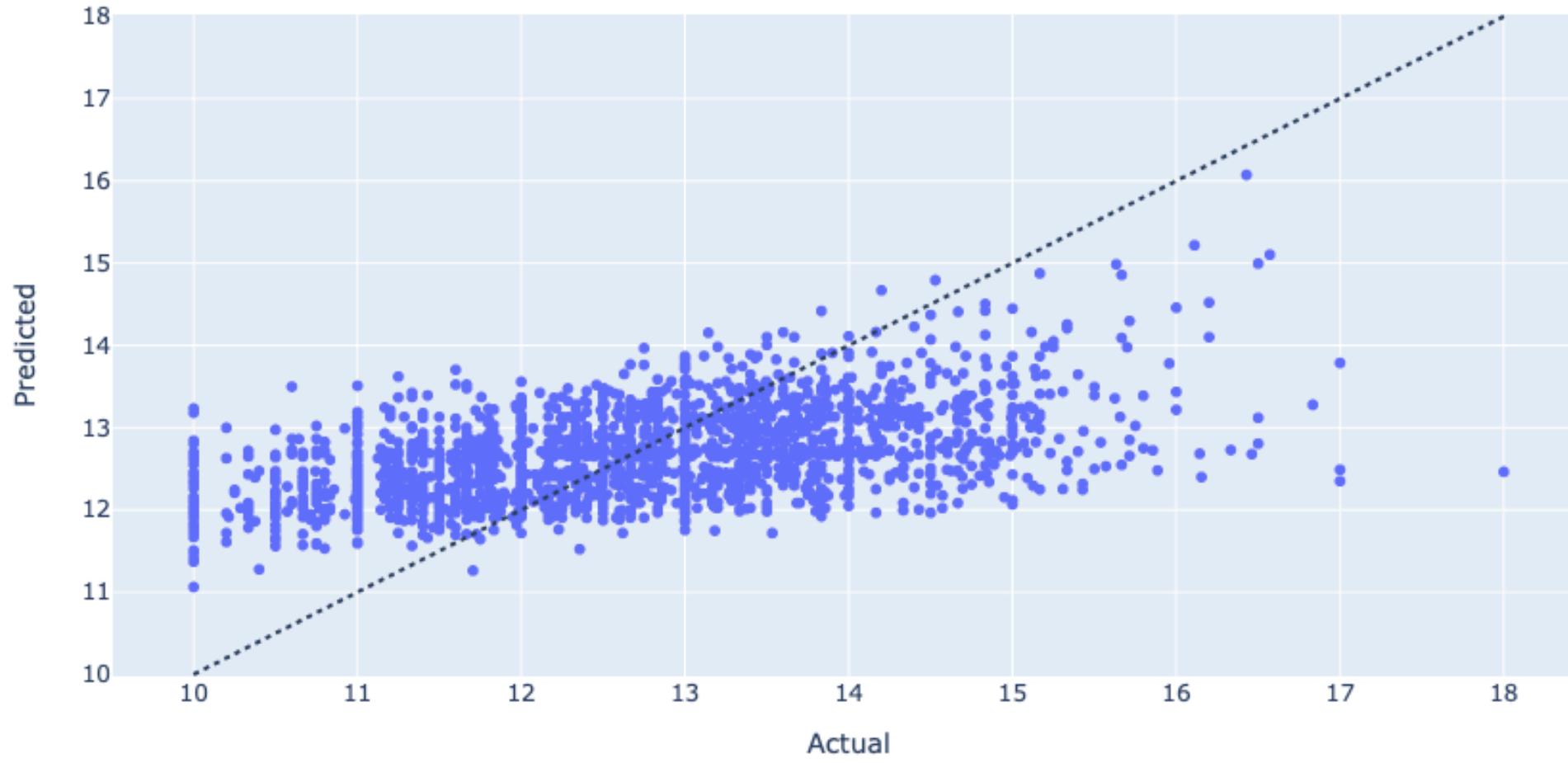
Building a series of weak learners (usually decision trees) sequentially, with each new learner correcting the errors made by the previous ones.

It fits a sequence of models to the data, where each new model focuses on the examples that were poorly predicted by the previous models

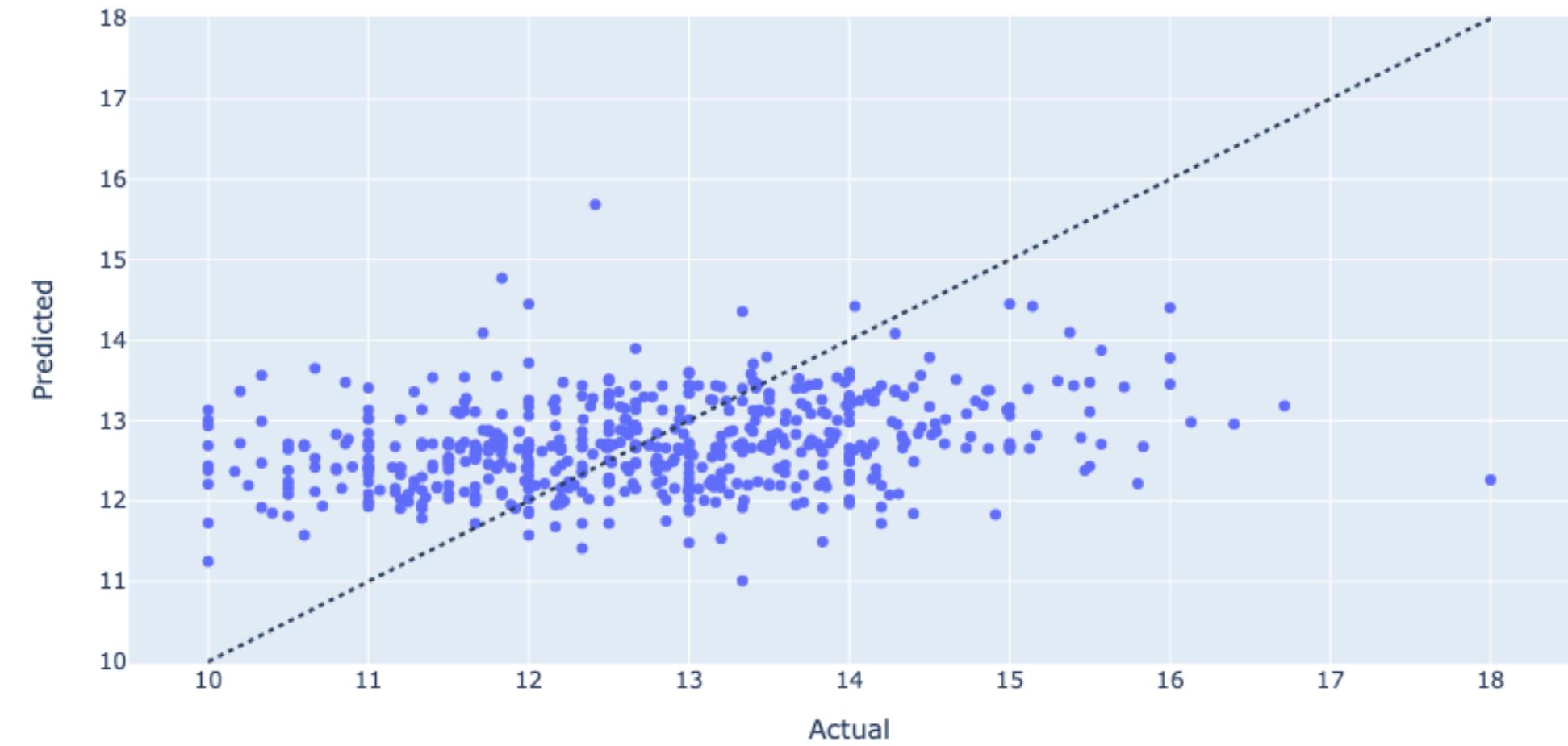
Gradient boosting typically produces highly accurate predictions compared to other machine learning algorithms.

# GRADIENT BOOSTING

Actual vs Predicted - Training Data (Cleaned)



Actual vs Predicted - Testing Data (Cleaned)



# Best Model

## LINEAR REGRESSION

Linear Regression performed the best with a lower average MSE when performing K-Fold Cross Validation.



### Linear Relationships

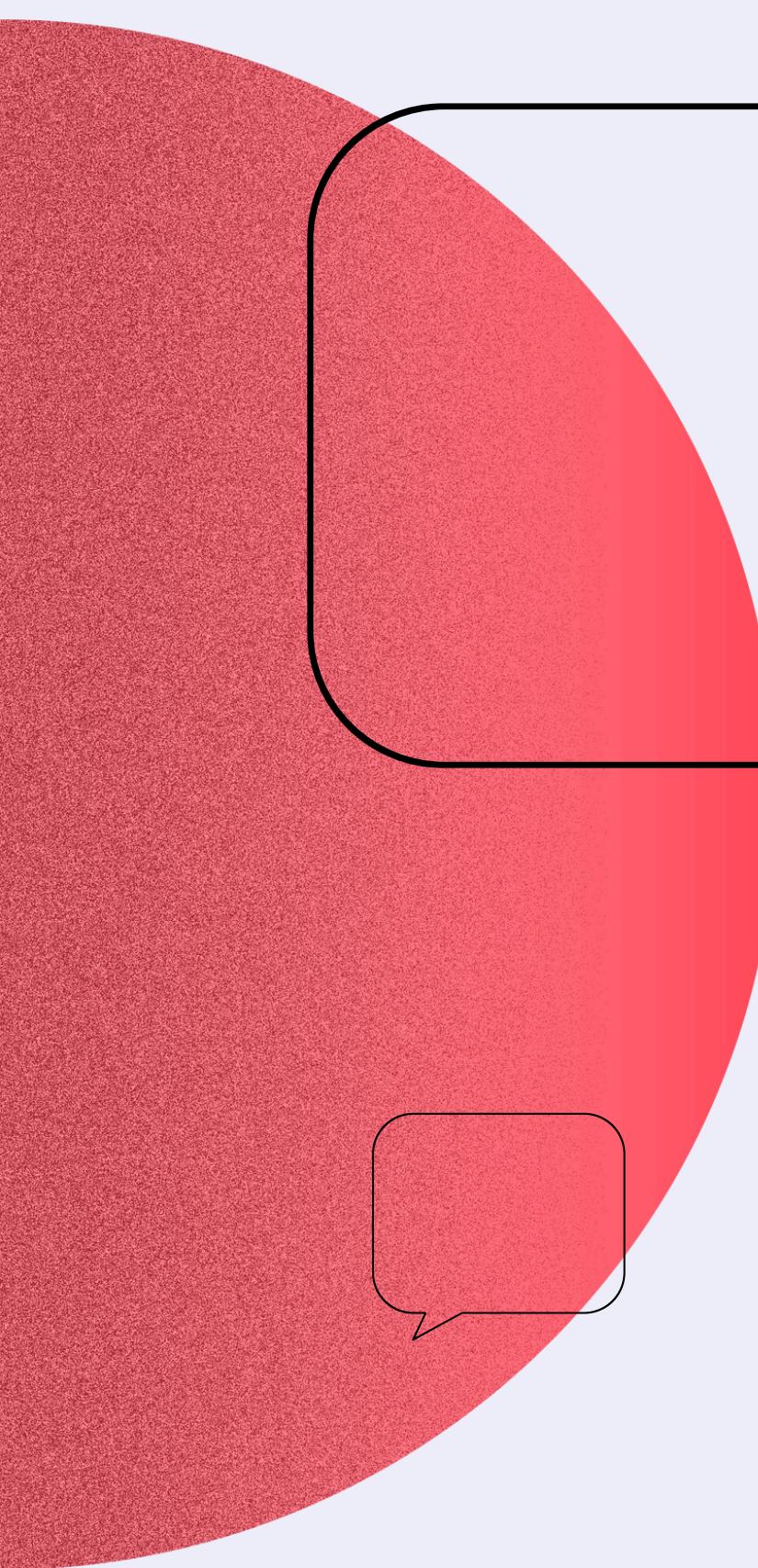
If the relationship between the features and the target variable is approximately linear, a linear regression model may provide a better fit

### Avoidance of Overfitting

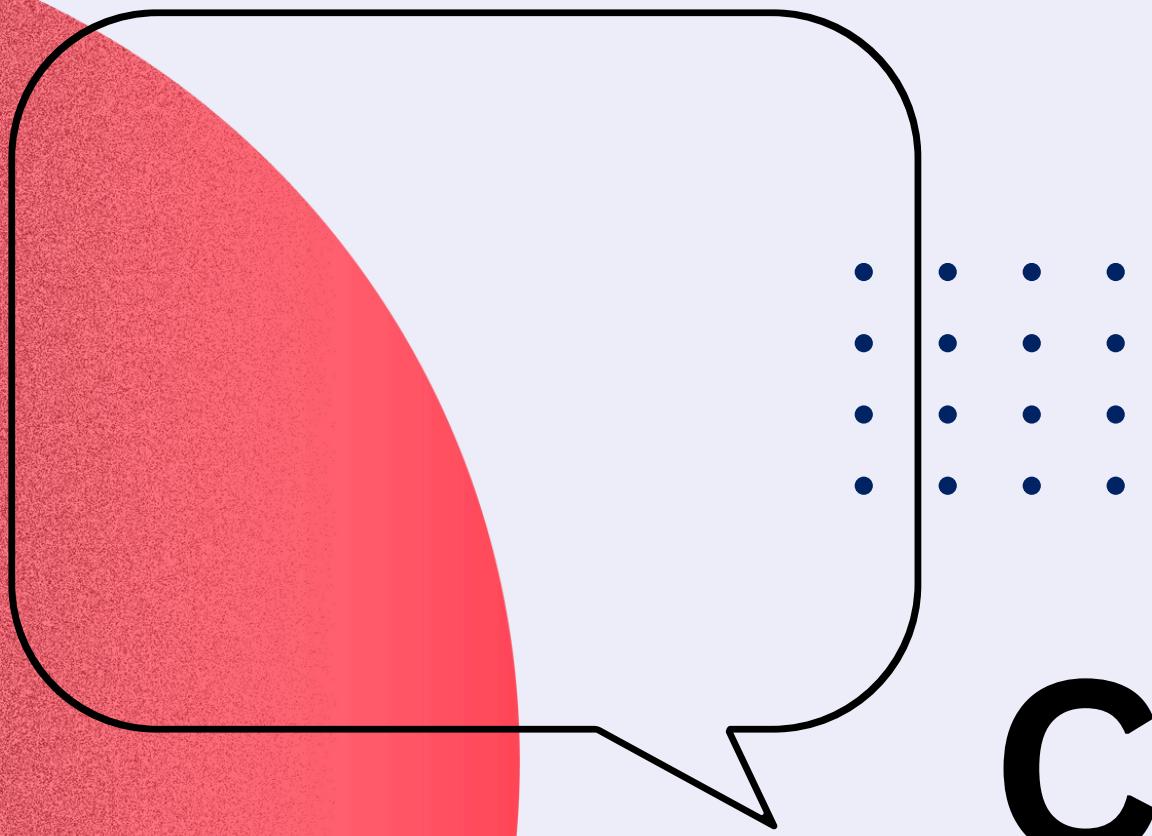
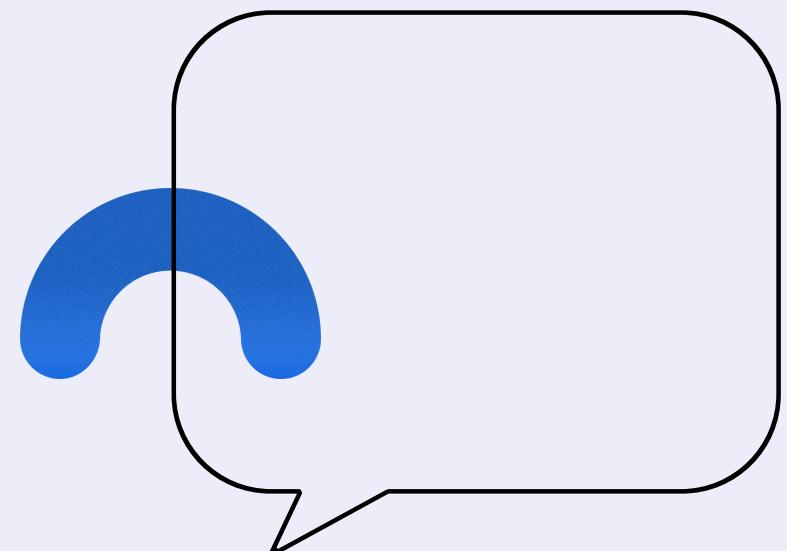
If there is limited training data or a risk of overfitting, a simpler linear regression model might generalize better.

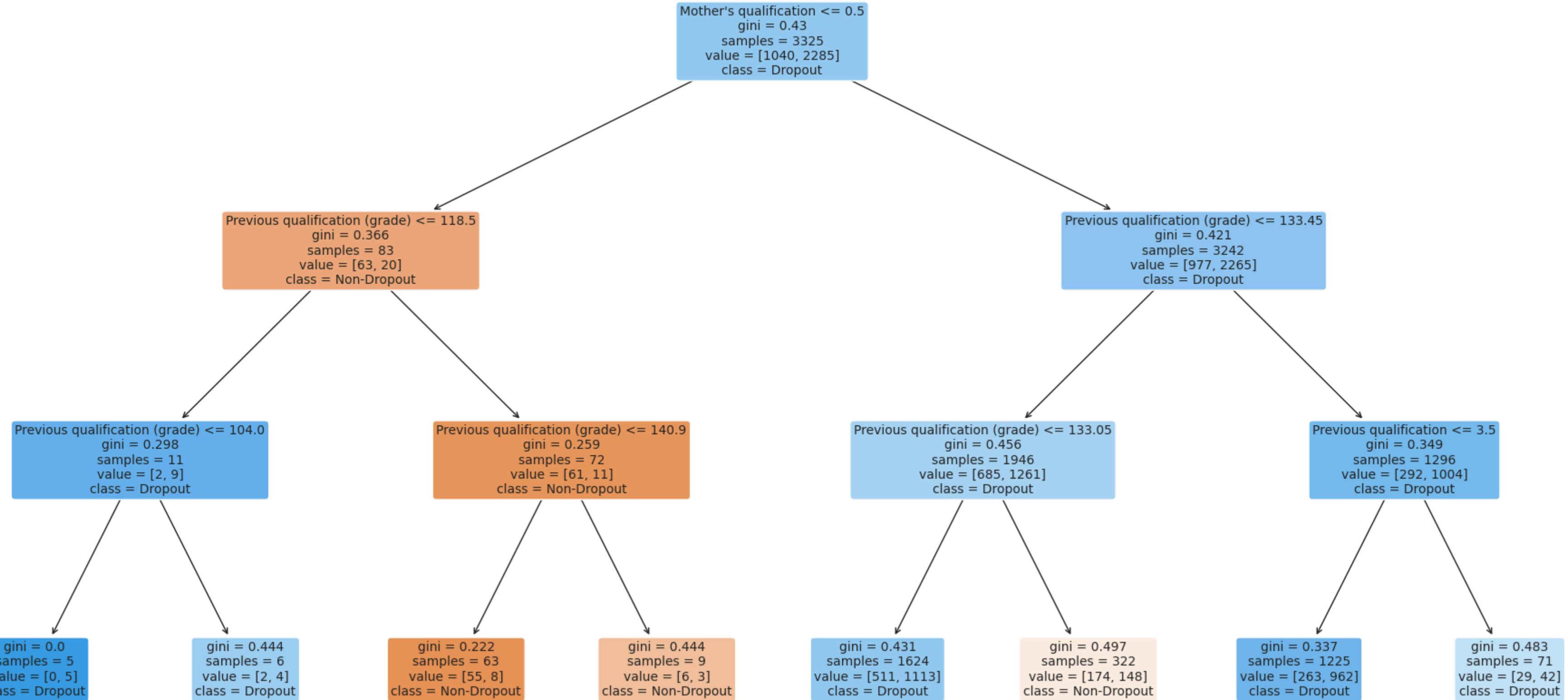
### Low-Dimensional Data

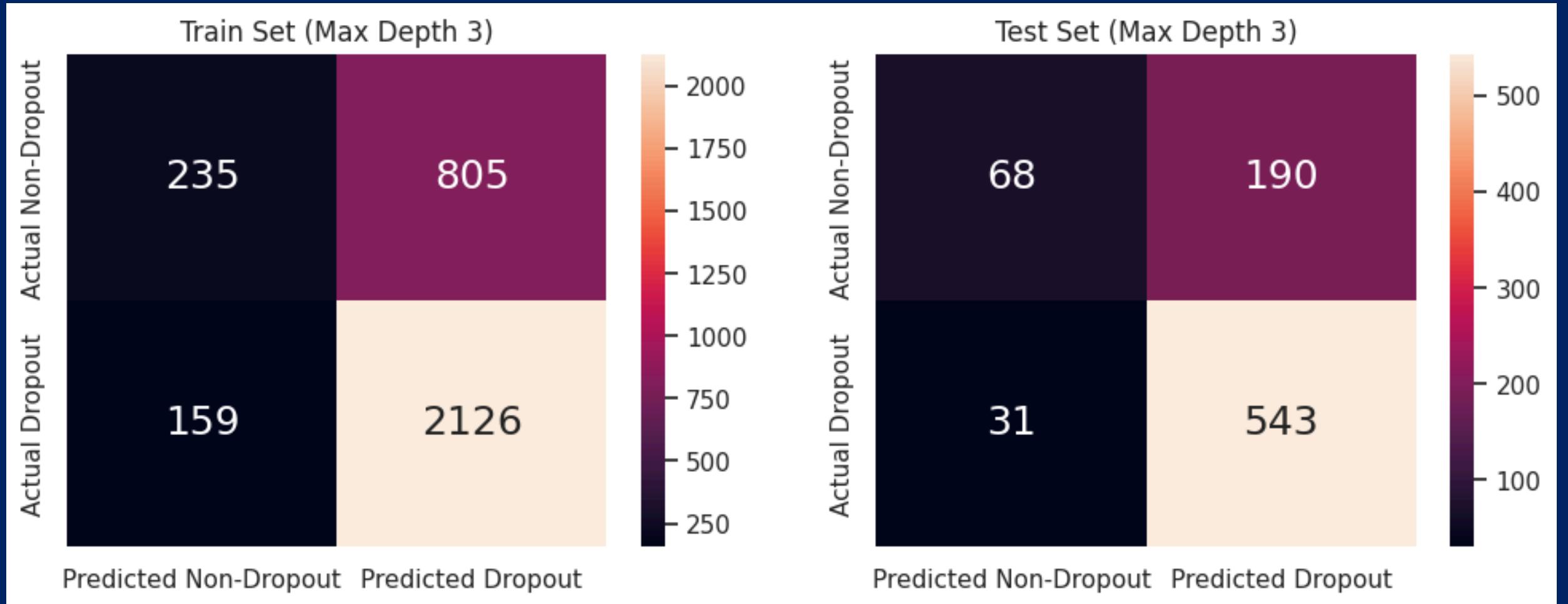
Tree-based models may be prone to overfitting in such cases.



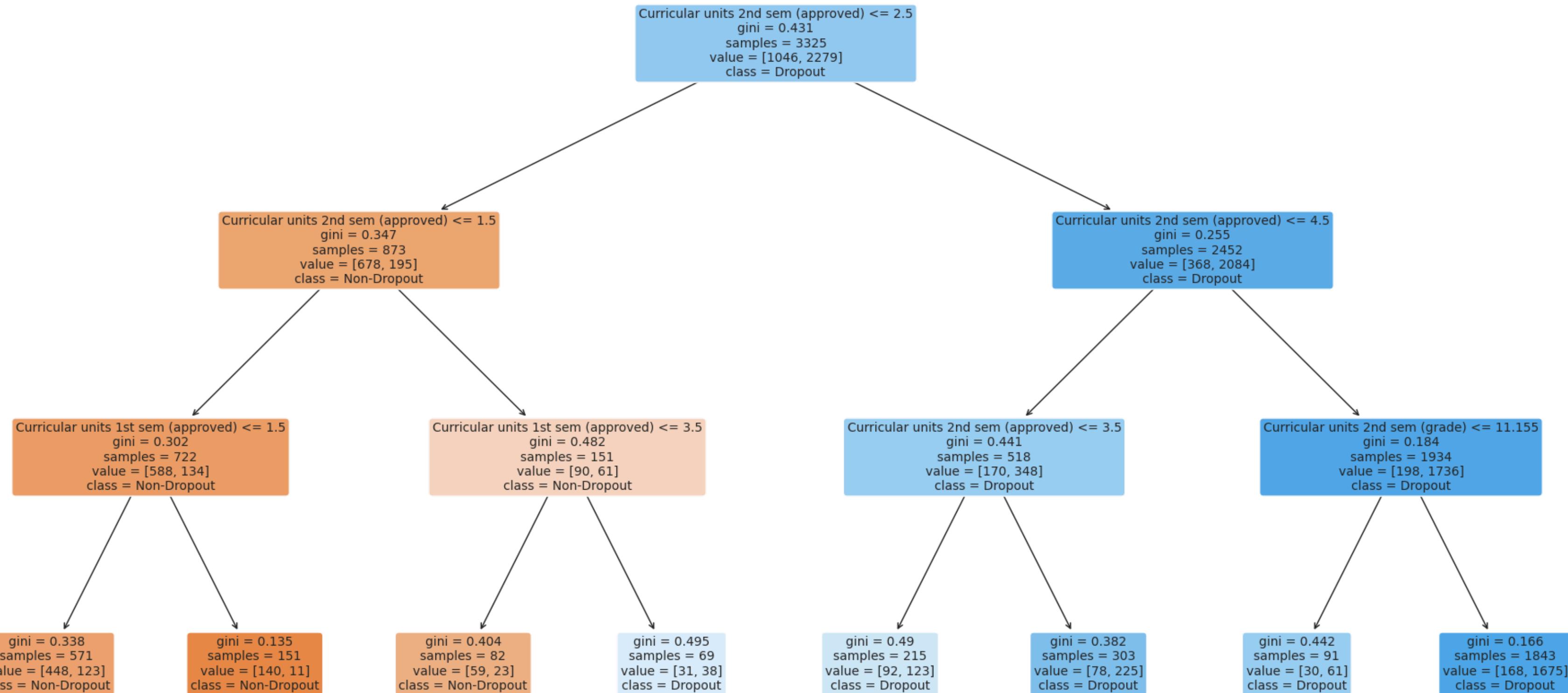
# Classification Model

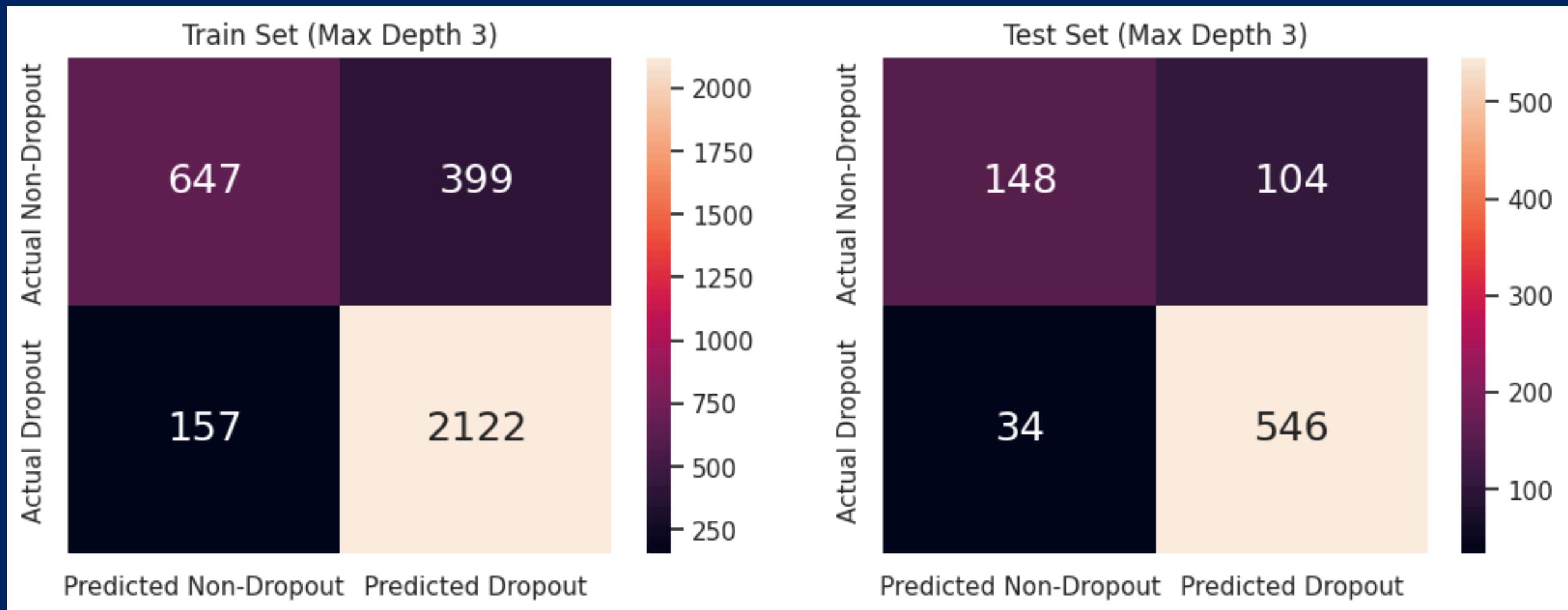




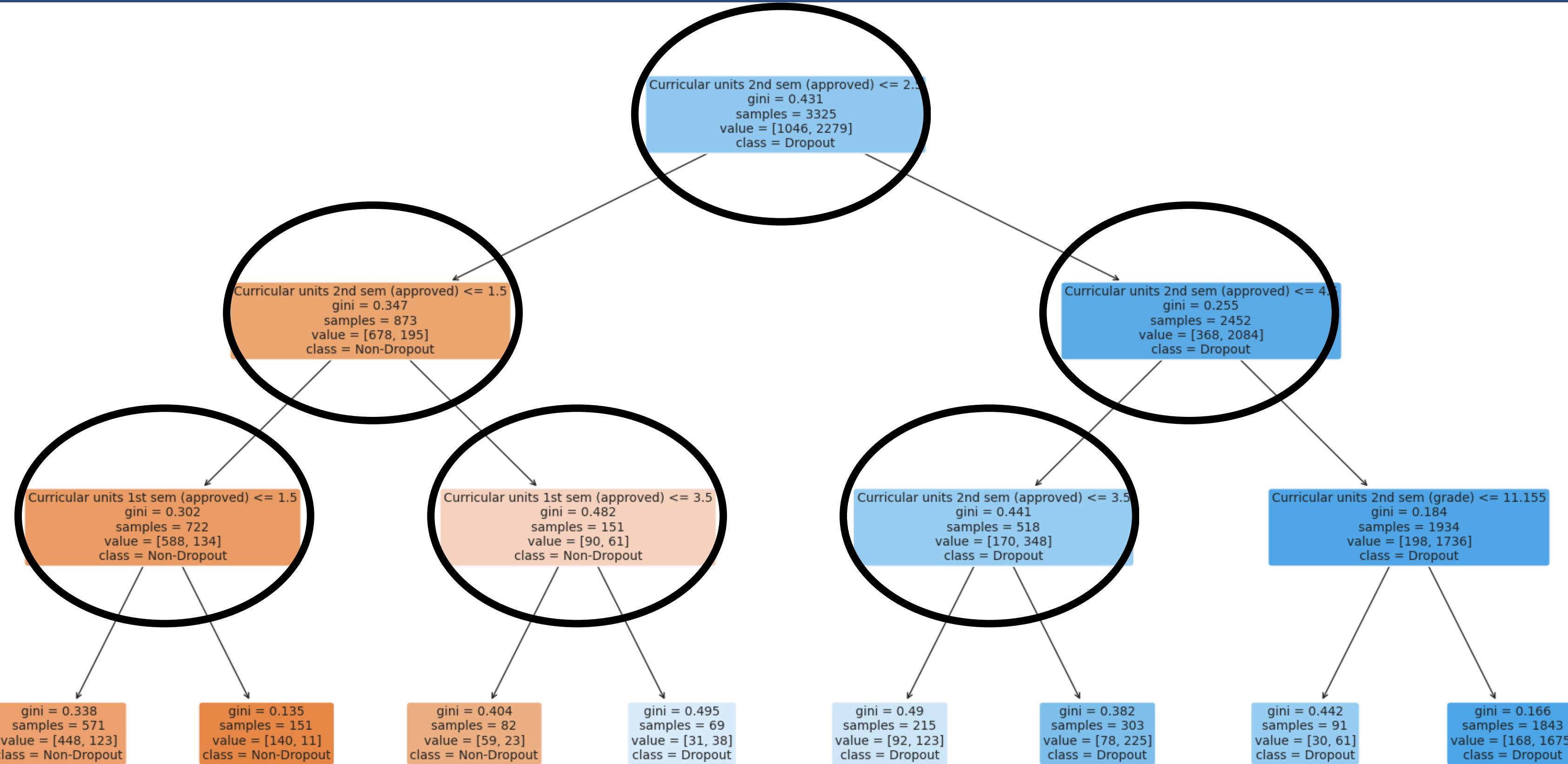


Goodness of Fit of Model	Train Dataset
Classification Accuracy	: 0.7100751879699249
F1 Score	: 0.3277545327754533
Goodness of Fit of Model	Test Dataset
Classification Accuracy	: 0.734375
F1 Score	: 0.380952380952381





Goodness of Fit of Model	Train Dataset
Classification Accuracy	: 0.8327819548872181
F1 Score	: 0.6994594594594594
Goodness of Fit of Model	Test Dataset
Classification Accuracy	: 0.8341346153846154
F1 Score	: 0.6820276497695853





# PROJECT OUTCOME

- Predictions using Regression Models
- Predictions using Classification Models

# CONCLUSION

Helping aid for  
Universities/Teachers to  
identify students to  
administer help to reduce  
dropout rates





**THANK YOU**