

Student Dropout Prediction

Data Preprocessing Report



Description of Data Cleaning Steps:

1. Column Name Cleaning:

- Any tabs or whitespaces in column names were removed to standardize them.

2. Handling Missing Values:

- Missing values were handled by either dropping or filling them. The default method was dropping missing values, though the option to fill missing values with a specific value was available.

3. Outlier Detection and Treatment:

- The Interquartile Range (IQR) method was applied to detect outliers in numerical columns such as 'Admission grade,' 'Curricular units 1st sem (grade),' and 'Curricular units 2nd sem (grade).' These outliers were then capped at the 1st and 99th percentiles.

Summary of Data Quality Issues:

Missing Values: There were missing values in the dataset that were either dropped or filled depending on the selected method.

Outliers: Outliers in admission grades and curricular unit grades were identified and capped using the IQR method.

Justification for Chosen Data Transformation Methods:

Handling Missing Values: Dropping missing values was chosen to ensure data consistency for the analysis, though filling missing values can be justified when domain knowledge provides an appropriate fill strategy.

Outlier Capping: Capping outliers using the 1st and 99th percentiles ensured that extreme values did not overly influence the statistical models.

Scaling: Numerical features were scaled using the 'StandardScaler', which standardizes the data by removing the mean and scaling to unit variance. This is crucial when using machine learning algorithms sensitive to feature scaling.



Statistical Analysis Report

Descriptive Statistics for All Variables:

- Descriptive statistics for the numerical variables, such as the mean, standard deviation, minimum, and maximum values, were generated. The key statistics for the numerical columns are summarized:
 - Admission grade: Provided detailed summary statistics.
 - Curricular units (1st and 2nd sem) grades: Descriptive analysis showed the distribution of student performance across these metrics.

Correlation Matrix Heatmap:

- A correlation matrix heatmap was generated for all numerical variables. Strong correlations, if present, were highlighted to help assess multicollinearity or associations between variables.

Results and Interpretation of Hypothesis Tests:

T-Test for Admission Grade (Dropout vs Graduate): A t-test was conducted to compare the admission grades of students who dropped out and those who graduated. The test result showed a significant difference:

- T-statistic: -7.858
- P-value: (5.10×10^{-15})

Interpretation: Since the p-value is very small, we reject the null hypothesis, indicating that the admission grades of students who dropped out significantly differ from those who graduated.

Chi-square Test for Gender vs Dropout: A chi-square test was performed to examine the relationship between gender and dropout status. The result was:

- Chi-square statistic: 233.27
- P-value: (2.22×10^{-51})

Interpretation: The very small p-value suggests a significant association between gender and dropout status, meaning gender has an influence on whether a student is likely to drop out.



Data loaded successfully!

☒ Show raw data

	Marital status	Application mode	Application order	Course	Daytime/evening attendance	Previous qualification	Previous qualification (grade)	Nacionality
0	1	17	5	171	1	1	122	1
1	1	15	1	9,254	1	1	160	1
2	1	1	5	9,070	1	1	122	1
3	1	17	2	9,773	1	1	122	1
4	2	39	1	8,014	0	1	100	1

☒ Show initial data exploration

Dataset Shape:

Rows: 4424, Columns: 37

Data Types:

	0
Curricular units 2nd sem (credit)	int64
Curricular units 2nd sem (enrolled)	int64
Curricular units 2nd sem (evaluation)	int64
Curricular units 2nd sem (approved)	int64
Curricular units 2nd sem (grade)	float64
Curricular units 2nd sem (withheld)	int64
Unemployment rate	float64
Inflation rate	float64
GDP	float64
Target	object

Summary Statistics:

	Marital status	Application mode	Application order	Course	Daytime/evening attendance	Previous qualification	Previous qualification (grade)	Nacion:
count	4,424	4,424	4,424	4,424	4,424	4,424	4,424	4
mean	1.1786	18.6691	1.7278	8,856.6426	0.8908	4.5778	132.6133	1.8
std	0.6057	17.4847	1.3138	2,063.5664	0.3119	10.2166	13.1883	6.9
min	1	1	0	33	0	1	95	
25%	1	1	1	9,085	1	1	125	
50%	1	17	1	9,238	1	1	133.1	
75%	1	39	2	9,556	1	1	140	
max	6	57	9	9,991	1	43	190	

☒ Check for missing values

Missing Values per Column:

	0
empty	

☒ Show outliers

Outliers in Admission grade:

	0
21	776
22	1,002
23	1,037
24	1,073
25	1,126
26	1,177
27	1,215
28	1,254
29	1,283
30	1,503

Outliers in Curricular units 1st sem (grade):

	0
0	0
1	2
2	7
3	12
4	20
5	35
6	36
7	44
8	56
9	59

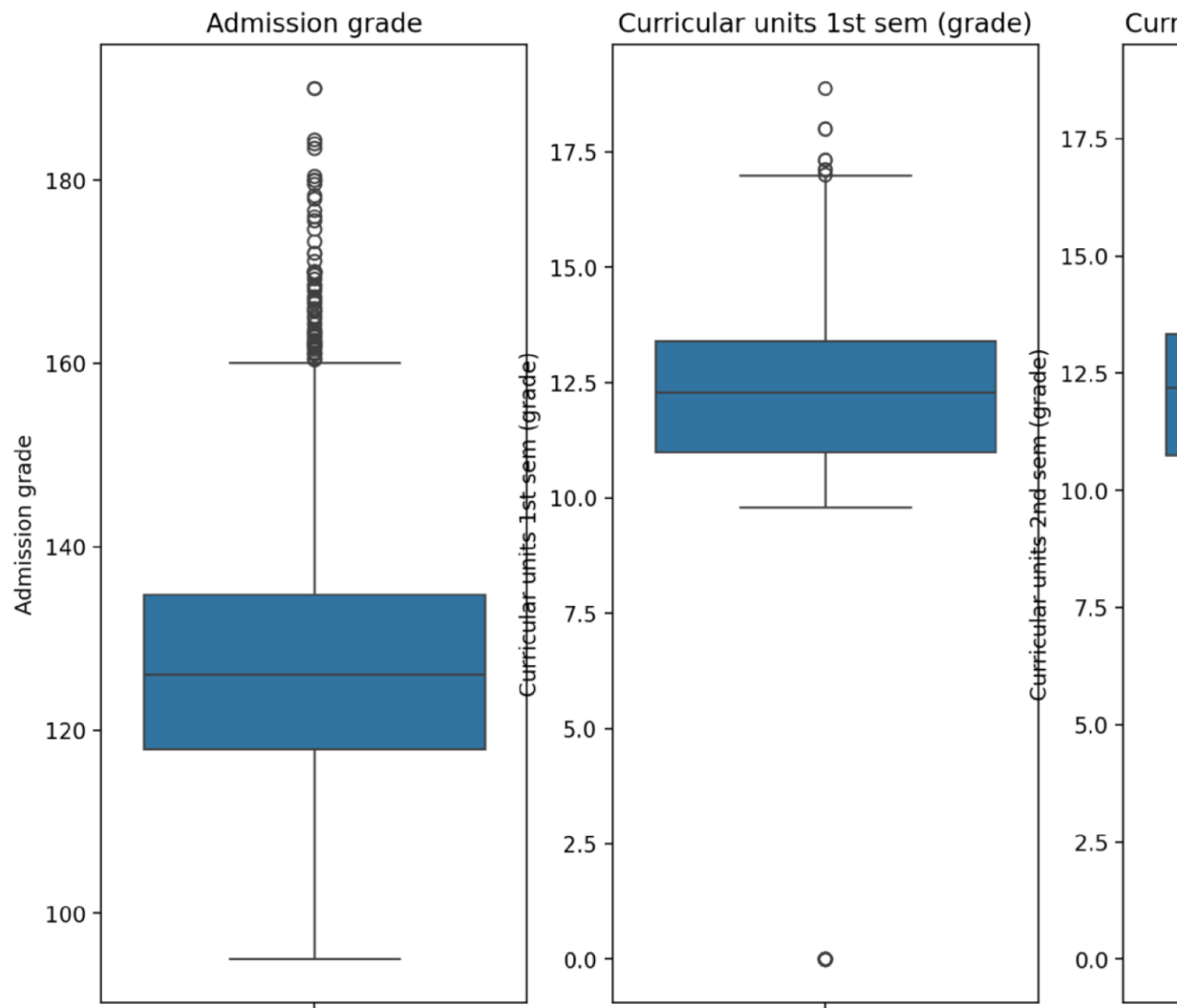


Outliers in Curricular units 2nd sem (grade):

	0
0	0
1	2
2	7
3	12
4	15
5	20
6	36
7	40
8	44
9	56

☒ Show box plots for numerical

columns Box plots for numerical columns:



☒ Show preprocessed data

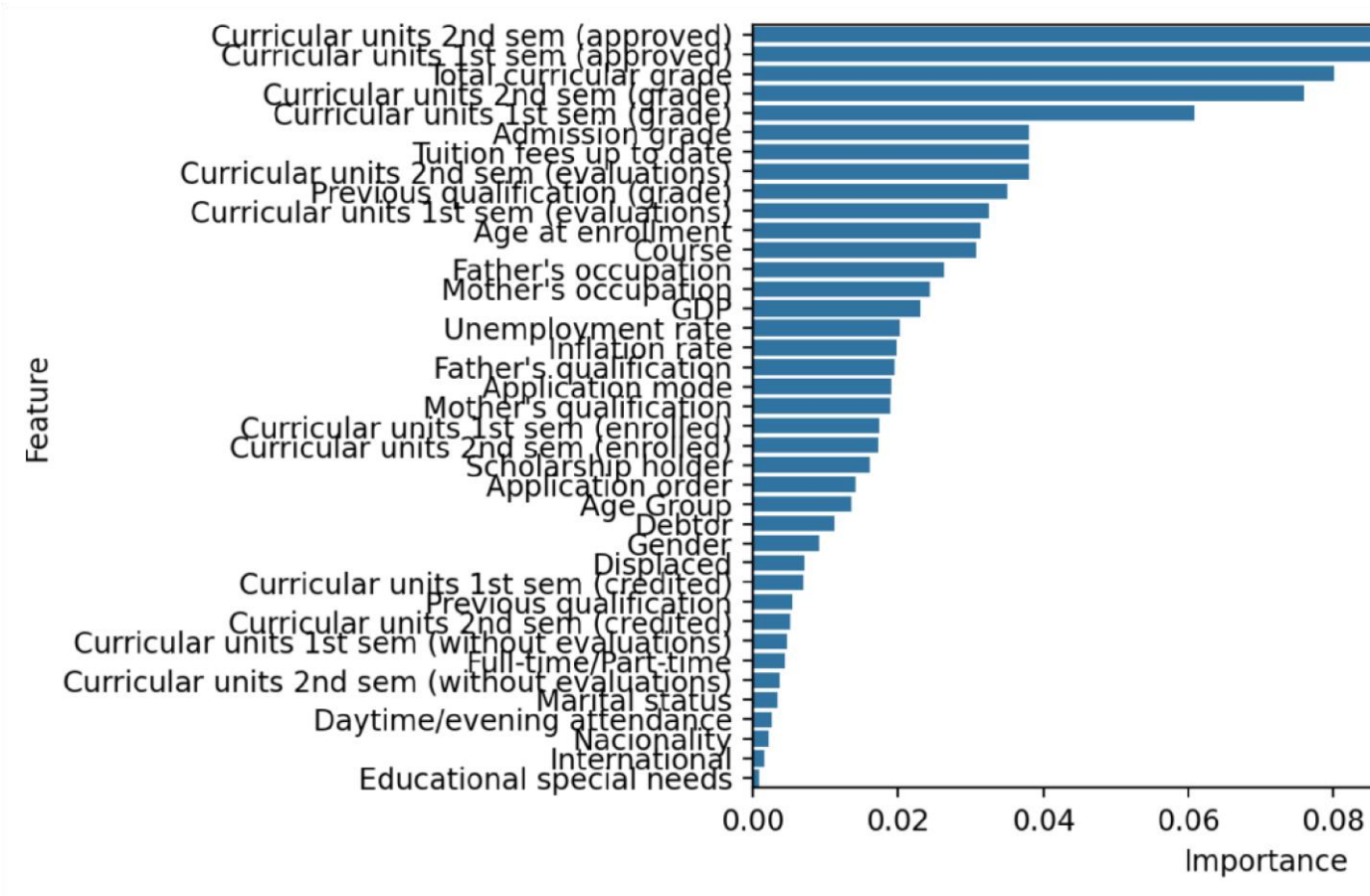
	Marital status	Application mode	Application order	Course	Daytime/evening attendance	Previous qualification	Previous qualification (grade)	Nacionality
0	1	17	5	171	1	1	122	1
1	1	15	1	9,254	1	1	160	1
2	1	1	5	9,070	1	1	122	1
3	1	17	2	9,773	1	1	122	1
4	2	39	1	8,014	0	1	100	1

☒ Show feature importance

Feature Importance using RandomForest:



	Feature	Importance
30	Curricular units 2nd sem (approved)	0.1346
24	Curricular units 1st sem (approved)	0.0863
37	Total curricular grade	0.0801
31	Curricular units 2nd sem (grade)	0.076
25	Curricular units 1st sem (grade)	0.0608
12	Admission grade	0.0381
16	Tuition fees up to date	0.038
29	Curricular units 2nd sem (evaluations)	0.038
6	Previous qualification (grade)	0.0351
23	Curricular units 1st sem (evaluations)	0.0325



☒ Show descriptive statistics

Numerical Descriptive Statistics:

	Marital status	Application mode	Application order	Course	Daytime/evening attendance	Previous qualification	Previous qualification (grade)	National average
count	4,424	4,424	4,424	4,424	4,424	4,424	4,424	4
mean	1.1786	18.6691	1.7278	8,856.6426	0.8908	4.5778	132.6133	1.8
std	0.6057	17.4847	1.3138	2,063.5664	0.3119	10.2166	13.1883	6.9
min	1	1	0	33	0	1	95	
25%	1	1	1	9,085	1	1	125	
50%	1	17	1	9,238	1	1	133.1	
75%	1	39	2	9,556	1	1	140	
max	6	57	9	9,991	1	43	190	

☒ Show correlation matrix heatmap

Marital status	-1.0020.0.09.20.06.02.0110.10.03.06.01.20.0306.09.01.0552.0306.05.06.03.0503.06.04.04.04
Application mode	0.21.00.20.01.30.42.04.01.02.08.05.09.01.30.0312.0.16.10.52.00.25.16.20.03.1205.24.13.10
Application order	-1.13.2900.06.10.10.06.02.06.05.04.03.1030.06.0706.0900.20.01.0.02.0904.06.01.0.08.01
Course	0.05.07.01.00.0400.03.03.05.03.00.0.09.02.0302.1002.04.03.1030.20.18.39.06.0940.28
Daytime/evening attendance	-1.27.0.16.0.00.0705.01.20.10.02.0201.25.03.01.04.0109.40.09.10.04.0502.06.06.11.00.00
Previous qualification	0.0640.0.00.01.00.10.03.01.0101.02.16.12.0110.0706.0716.0310.03.13.02.00.00.14.06.10
Previous qualification (grade)	-0.02.04.06.0805.1.00.06.06.03.01.0258.01.00.0406.0506.11.06.01.03.0705.06.00.02.03.01
Nacionality	-0.01.00.02.0302.0301.00.05.0904.02.02.01.0106.03.02.01.0170.00.01.01.00.00.00.01.02.01
Mother's qualification	0.19.12.0609.20.01.06.01.00.54.03.06.05.03.0202.03.0505.29.0405.05.06.01.0400.03.04.04
Father's qualification	0.13.08.0509.20.01.04.0951.00.05.06.05.0500.01.02.0710.19.0904.04.04.04.01.0204.02.00
Mother's occupation	0.03.06.0406.0200.01.03.03.01.00.90.04.05.0010.01.00.0206.04.01.00.02.0201.00.01.00
Father's occupation	0.03.04.0306.0202.0202.05.00.91.00.04.0501.10.00.01.0203.02.00.00.00.02.01.00.01.01.01
Admission grade	-0.01.01.10.0.01.13.53.02.05.05.04.01.00.00.02.0205.01.02.0302.04.03.0707.07.01.04.03.01
Displaced	-1.23.0.36.0926.10.01.01.03.06.05.05.01.00.00.0910.0.00.30.01.10.06.0305.06.02.03.03.01
Educational special needs	-0.03.0306.0206.01.00.01.0200.00.00.0201.00.00.00.00.00.00.00.00.00.00.00.00.00.00.00
Debtor	0.03.12.07.0301.10.0405.02.0110.10.02.0901.00.40.06.0710.03.06.01.04.10.10.00.06.03.00
Tuition fees up to date	-0.09.0.06.02.04.0706.03.03.02.01.00.05.10.00.4100.1014.10.04.00.06.03.24.26.0501.09.00
Gender	-0.0116.09.10.0106.05.02.06.0700.0101.10.0206.1.00.0.16.0302.10.01.10.10.10.10.10.10.10
Scholarship holder	-0.05.0.07.02.09.0706.0105.10.02.0202.07.02.0710.1700.10.03.03.00.0515.10.01.02.03.00
Age at enrollment	0.5250.20.09.40.16.10.0209.19.06.06.01.30.0410.10.19.1900.0120.14.14.05.10.06.20.09.00
International	-0.0300.03.0306.0301.70.04.0904.02.02.01.0006.04.03.03.01.00.01.00.01.01.01.03.00.01.00
Curricular units 1st sem (credited)	0.0620.10.10.10.10.01.0005.04.00.00.00.10.0203.00.02.0920.01.00.70.54.63.12.13.94.64.40
Curricular units 1st sem (enrolled)	0.0516.0236.0406.03.0105.04.00.00.03.06.03.0106.10.00.14.00.77.00.68.70.38.13.75.94.60
Curricular units 1st sem (evaluations)	0.0626.0920.0516.07.0105.04.02.00.07.03.0304.06.02.0614.01.54.63.00.52.42.24.52.60.70
Curricular units 1st sem (approved)	-0.03.0304.18.02.02.06.00.01.00.02.0207.06.02.1120.0.16.0501.63.70.51.00.70.160.70.54
Curricular units 1st sem (grade)	-0.06.1206.39.06.00.06.00.03.01.02.0107.06.01.1020.0.10.10.12.38.40.71.00.0710.41.49
Curricular units 1st sem (without evaluations)	0.03.06.0303.05.00.00.01.00.0201.00.00.00.00.02.01.00.05.01.0506.03.12.13.24.01.01.12.10.14
Curricular units 2nd sem (credited)	0.0620.10.10.10.01.00.04.04.00.01.04.09.0203.01.02.0821.00.94.70.52.60.10.11.00.60.40
Curricular units 2nd sem (enrolled)	0.04.13.03.40.00.06.03.0204.02.01.01.03.04.03.0300.0.03.09.01.60.90.60.70.40.11.63.00.60
Curricular units 2nd sem (evaluations)	0.0210.0520.01.10.06.0302.01.00.00.06.03.0102.06.04.0206.00.43.60.70.54.49.14.43.61.00
Curricular units 2nd sem (approved)	-0.04.0707.20.06.01.06.02.0100.03.0303.06.00.10.20.0.20.10.0149.0.44.90.60.0152.70.41
Curricular units 2nd sem (grade)	-0.07.1206.35.05.00.06.01.03.01.02.0207.07.01.10.30.0.10.10.0113.36.36.69.80.0513.40.41
Curricular units 2nd sem (without evaluations)	0.02.06.0206.00.00.02.01.00.01.01.01.03.01.06.0706.0506.01.06.07.18.05.0758.07.07.14
Unemployment rate	-0.0209.1001.06.10.06.00.10.03.09.1000.0.05.02.01.02.06.06.01.03.04.06.05.00.0501.06.00
Inflation rate	0.00.02.0102.02.0502.01.06.06.02.06.02.01.00.02.00.00.0303.01.02.04.01.01.03.0501.02.00
GDP	-0.03.0206.0202.06.0506.03.0712.18.0206.01.06.00.01.04.0504.03.03.1002.09.10.02.01.00
Target	-0.09.20.09.03.08.0510.01.04.00.01.0012.10.00.20.40.0.30.20.00.05.16.04.53.49.0705.18.00
Total curricular grade	-0.01.0.06.38.06.00.06.00.03.01.02.0103.07.01.10.20.0.10.10.00.13.39.40.72.90.0713.42.49
Full-time/Part-time	0.0509.01.80.0706.09.0304.03.02.00.10.03.01.00.0901.00.0206.52.40.32.45.04.06.59.40

☒ Show t-test results for Admission Grade (Dropout vs Graduate)

T-test results: T-statistic: 0.834913824313304, P-value: 0.4038563525431528

☒ Show chi-square test results for Gender vs Dropout

Chi-square test results for Gender vs Dropout: Chi-square stat: 233.26643249623856, P-value: 2.2224795668092454e-51

Enter details for prediction:

Admission Grade

Model Prediction



0.00

Curricular Units 1st Sem Grade

0.00

Curricular Units 2nd Sem Grade

0.00

Age Group

<18

Predict

Insights, Recommendations, and Conclusions

Insights

- The dataset contains 4,424 rows and 37 columns.
- Outliers were detected in key numerical variables such as Admission grade and semester grades.

The T-test showed no significant difference between dropouts and graduates in Admission grades.

Gender shows a significant impact on dropout rates according to the Chi-square test.

Recommendations

- Investigate missing values and handle them appropriately to improve model performance.
- Consider outlier treatment techniques like transformations to avoid skewing model predictions.

Enhance feature engineering and explore additional models (e.g., XGBoost) for better accuracy.

Address dataset imbalance (if present) to improve dropout predictions, especially related to gender.



Conclusions

- Admission grades alone may not be a strong predictor of dropouts.
- Gender has a notable influence on dropout rates, indicating a need for gender-targeted interventions.

The current model is a good baseline, but improvements can be made with additional data processing and model tuning.