# eda

October 24, 2024

## 1 EDA Analysis

### 1.1 Background

The dataset sourced from the UCI Machine Learning Repo contains 4424 student entries with
36 features and one categorical target varaible. The goal is to predict student academic success
(dropout or graduate) from the 36 input features.

```python
[2]: !pip install ucimlrepo  --quiet
     from visualization_utils import *
     from file_utils import *
     import matplotlib.pyplot as plt
     from ucimlrepo import fetch_ucirepo
```

```python
[3]: # fetch dataset
     predict_students_dropout_and_academic_success = fetch_ucirepo(id=697)

     # data (as pandas dataframes)
     X = predict_students_dropout_and_academic_success.data.features
     y = predict_students_dropout_and_academic_success.data.targets

     # metadata
     metadata = predict_students_dropout_and_academic_success.metadata

     # variable information
     variable_info = predict_students_dropout_and_academic_success.variables
     df = X
     df['dropout'] = y
     print(f"{df.shape[0]} entries with {df.shape[1]} features")
```

```
4424 entries with 37 features
```

```python
[4]: predict_students_dropout_and_academic_success.variables['name']
     ds_vars = predict_students_dropout_and_academic_success.variables['name']
     ds_desc = predict_students_dropout_and_academic_success.variables['description']
     # need to create forward mappings for each one of the variables
     quantitative_vars = {'Application mode',
     'Application order',
     'Previous qualification (grade)',
```

```
'Admission grade',
'Age at enrollment',
"Curricular units 1st sem (credited)",
"Curricular units 1st sem (enrolled)",
"Curricular units 1st sem (evaluations)",
"Curricular units 1st sem (approved)",
"Curricular units 1st sem (grade)",
"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 (approved)",
"Curricular units 2nd sem (grade)",
"Curricular units 2nd sem (without evaluations)",
"Unemployment rate",
"Inflation rate",
"GDP",
"Target",
}

variable_map = generate_variable_map(quantitative_vars, ds_vars, ds_desc)
```
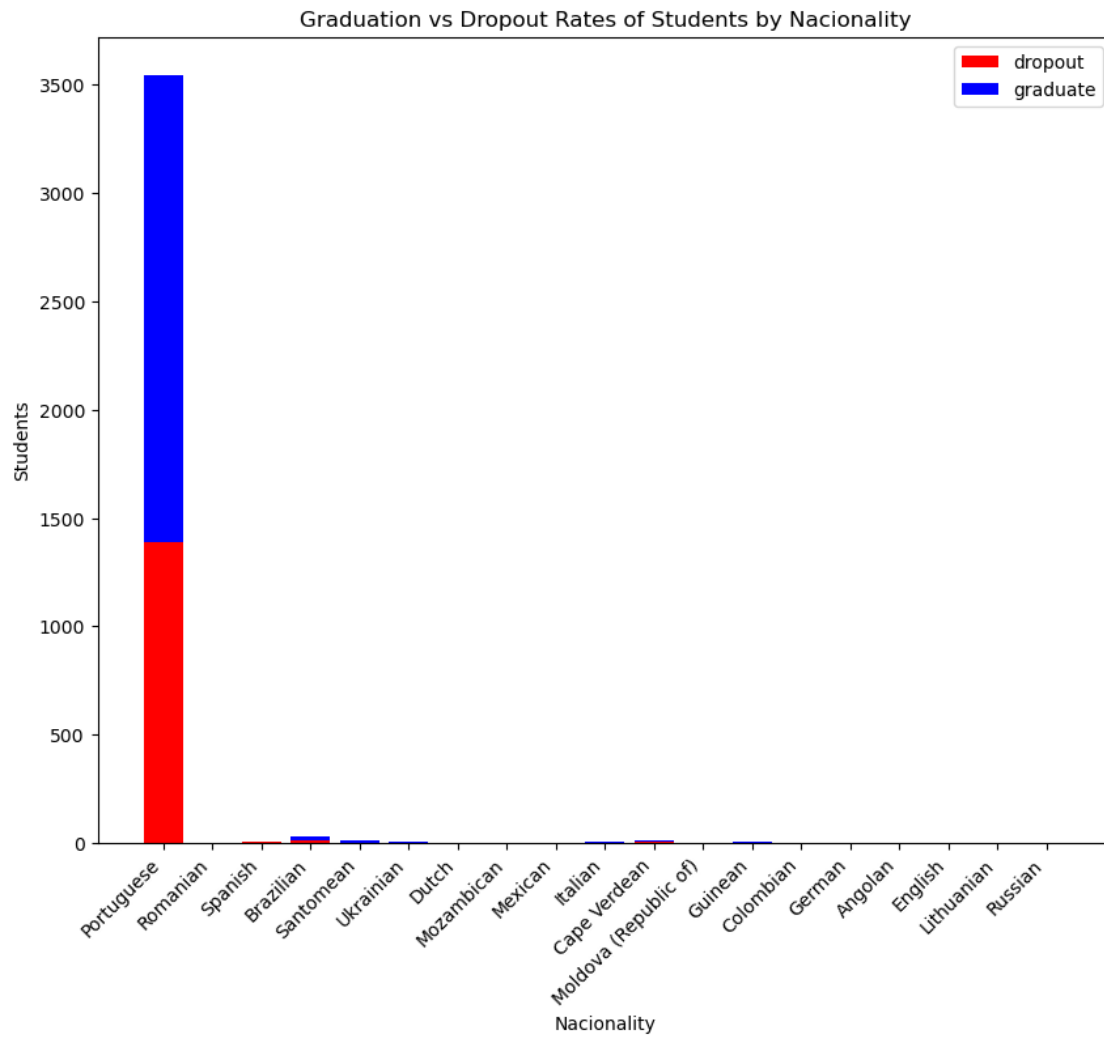
## 2 Analysis of Categorical Variables Affecting Student Performance
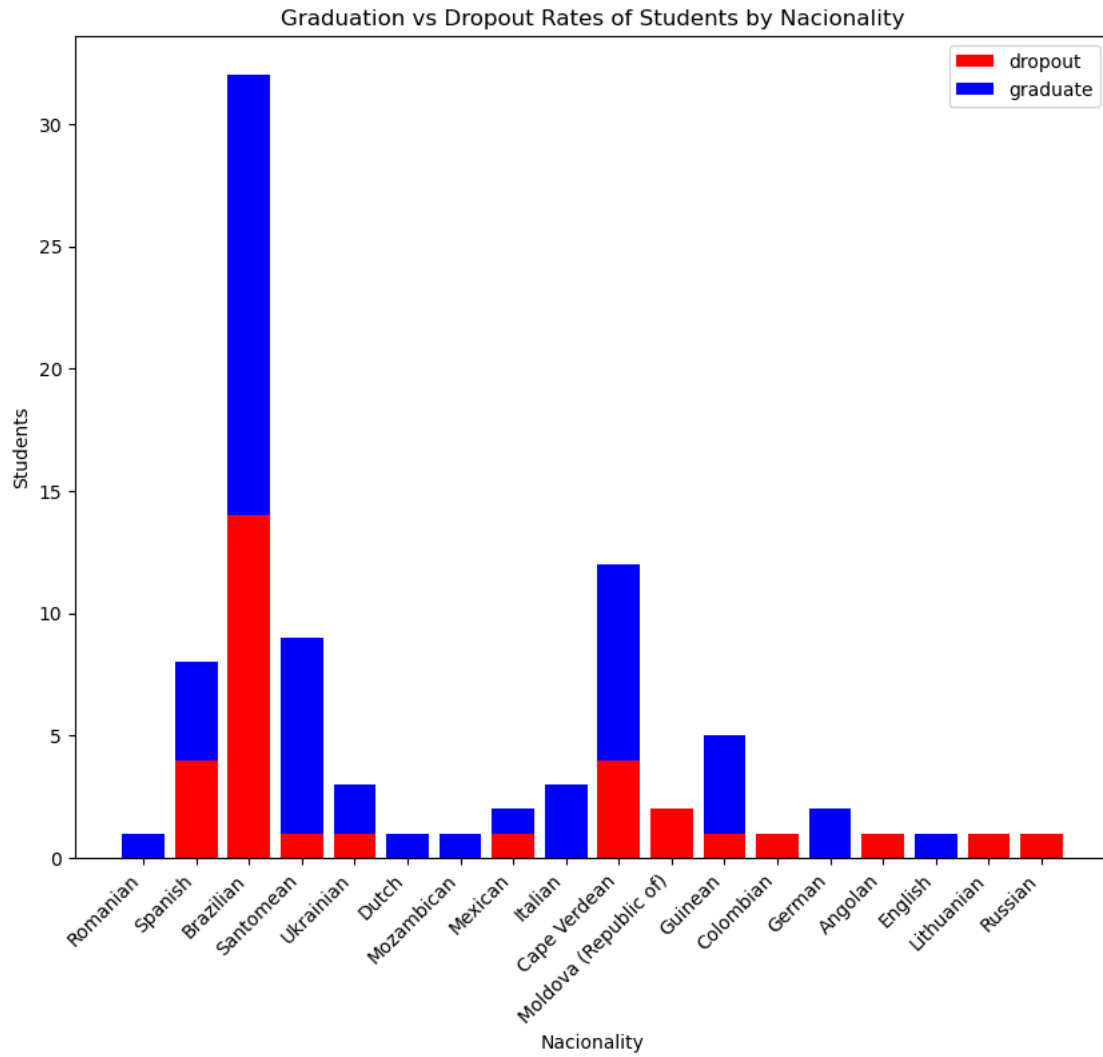
To gain a better understanding of variables responsibel for controlling the
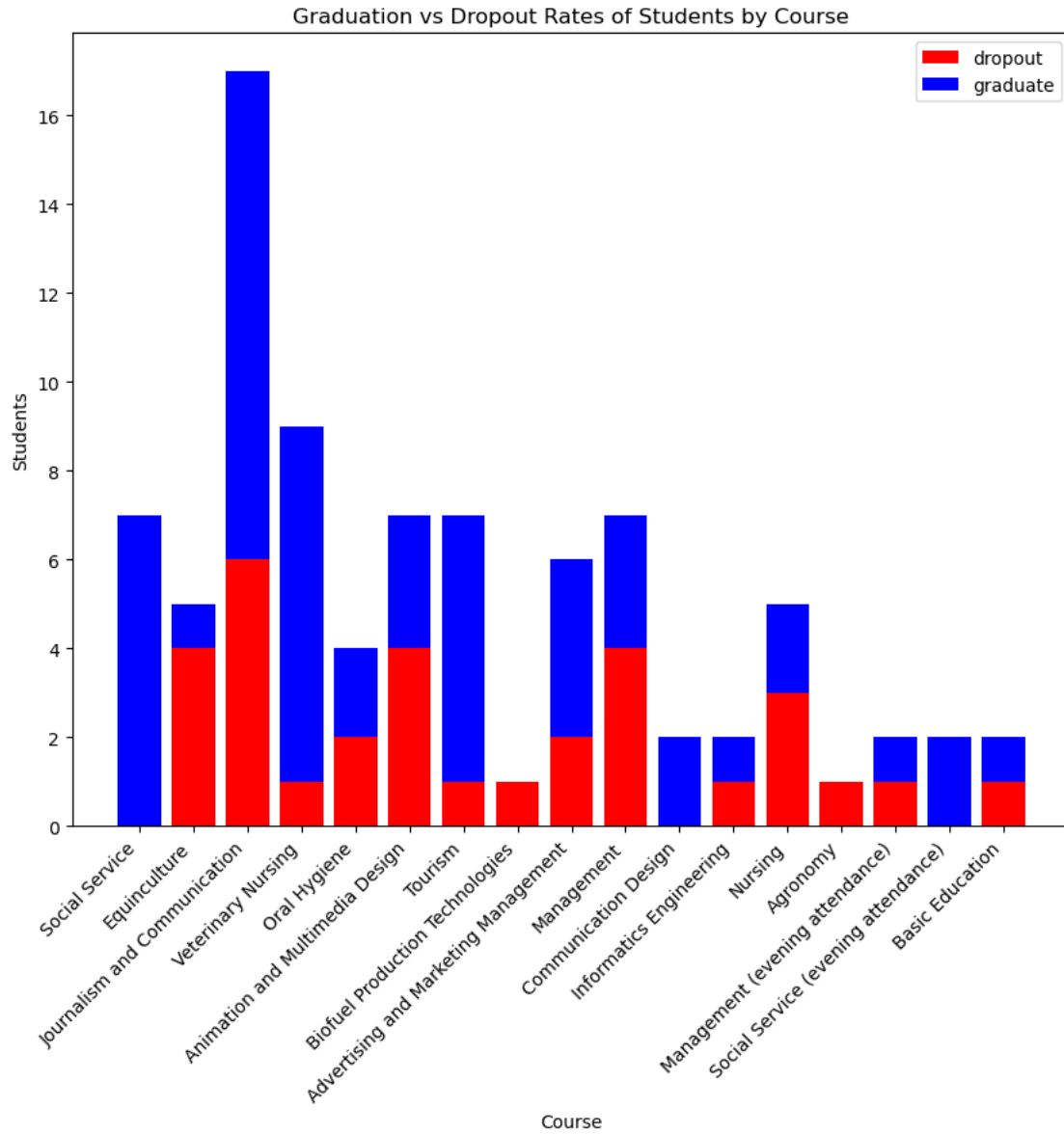
```
[5]: finished_df = df[df['dropout'].isin(['Dropout', 'Graduate'])]
     generate_general_stacked_bar_graph(finished_df, variable_map, 'Nacionality')
     finished_df = finished_df.drop(finished_df.loc[finished_df['Nacionality']==1].
       ↪index)
     generate_general_stacked_bar_graph(finished_df, variable_map, 'Nacionality')
```
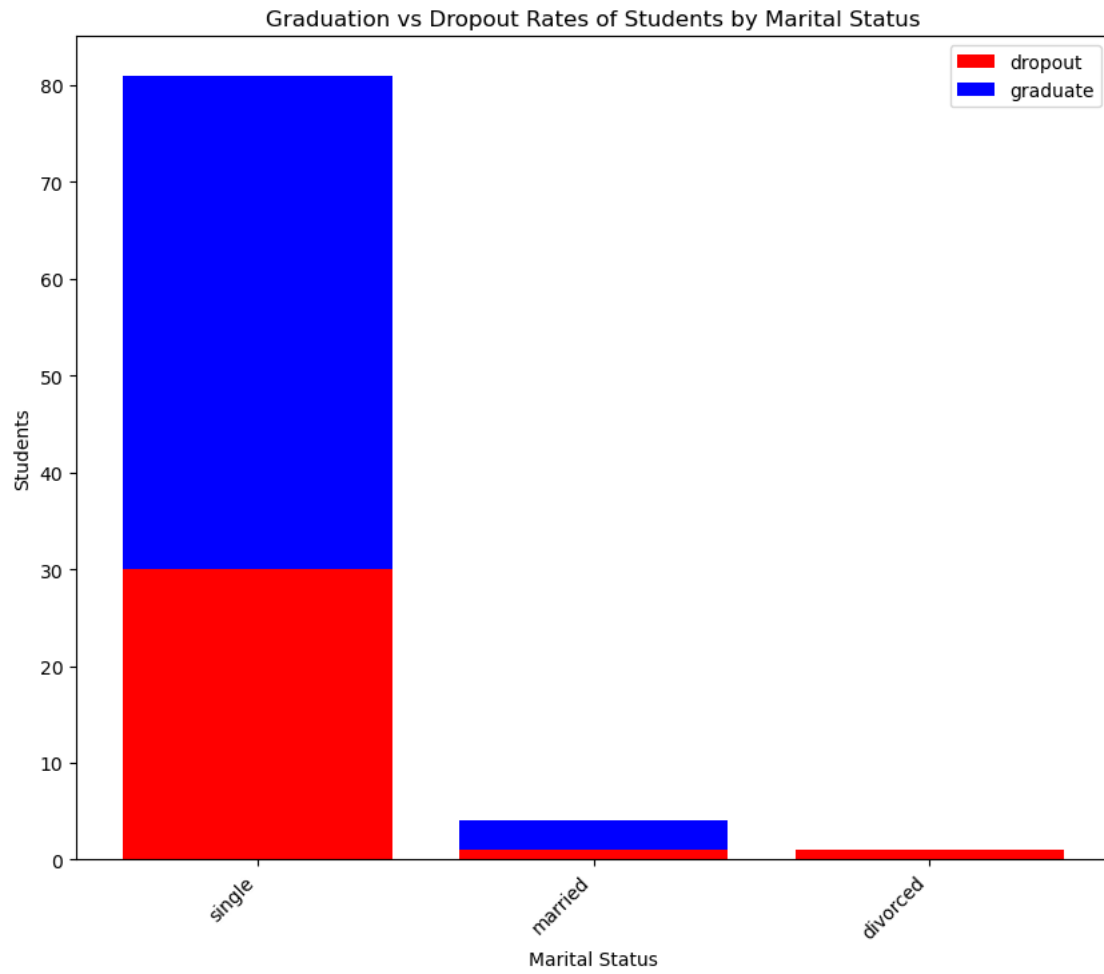
Graduation vs Dropout Rates of Students by Nacionality

Graduation vs Dropout Rates of Students by Nacionality

```
[6]: generate_general_stacked_bar_graph(finished_df, variable_map, 'Course')
```

Graduation vs Dropout Rates of Students by Course

```
[7]: generate_general_stacked_bar_graph(finished_df, variable_map, 'Marital Status')
```

Graduation vs Dropout Rates of Students by Marital Status

```
[8]: generate_general_stacked_bar_graph(finished_df, variable_map, 'Gender')
```

Graduation vs Dropout Rates of Students by Gender

```
[9]: generate_general_stacked_bar_graph(finished_df, variable_map, 'Daytime/evening␣
     ↪attendance')
```

Graduation vs Dropout Rates of Students by Daytime/evening attendance

```
[10]: PREV_QUAL = 'Previous qualification (grade)'
      ADMISSION_GRADE = 'Admission grade'
      tmp = df[[PREV_QUAL, ADMISSION_GRADE, 'dropout']]
      finished_df = tmp[tmp['dropout'].isin(['Dropout', 'Graduate'])]

      jitter_plot(finished_df)
```

Comparison between Admission Scores and Grades vs Dropout
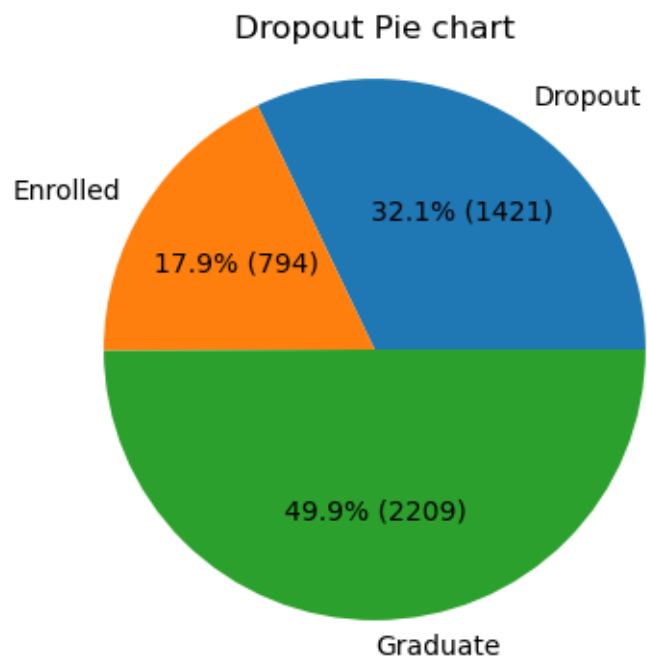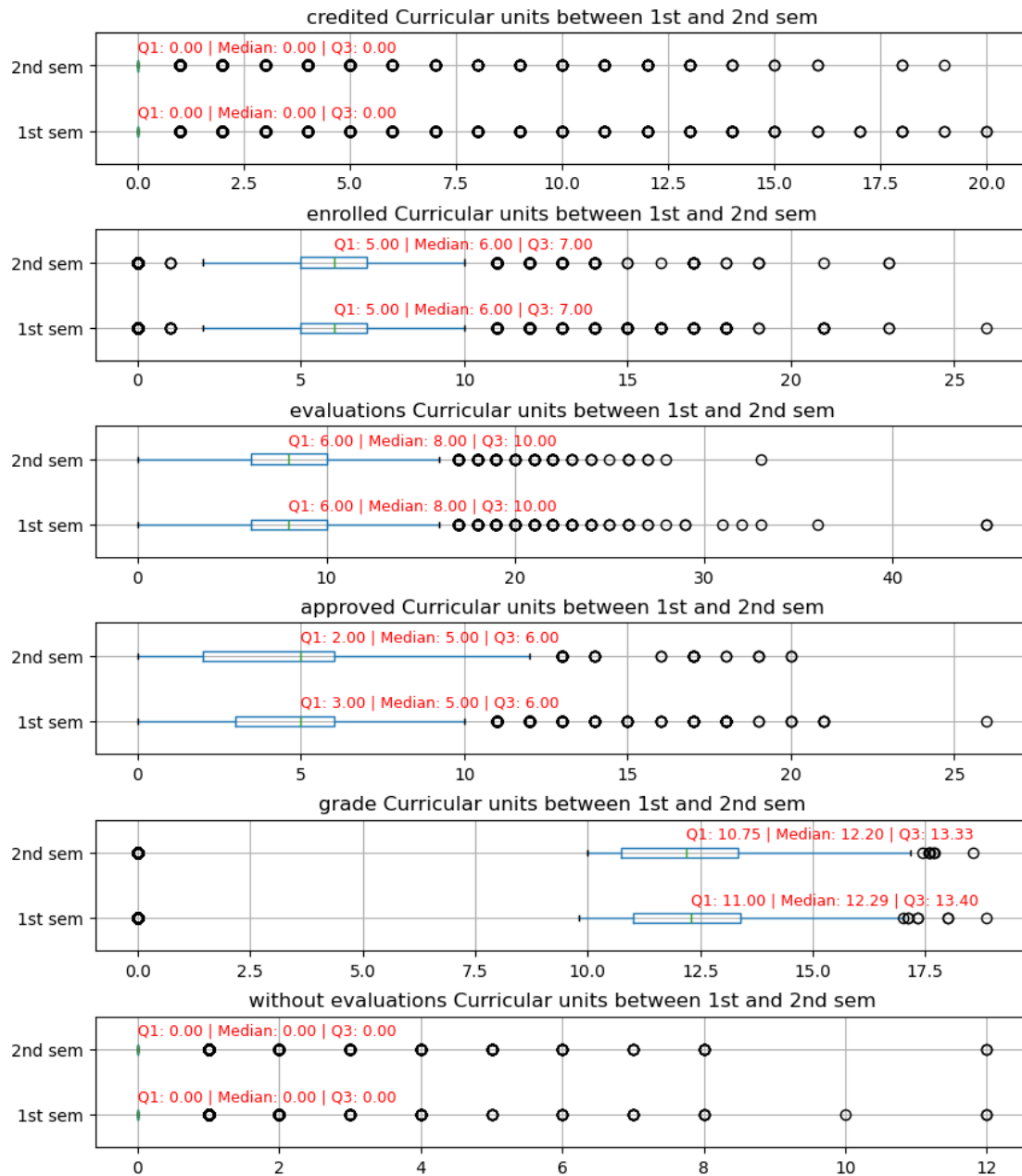
```
[11]: pie_chart(df['dropout'])
```



Dropout Pie chart
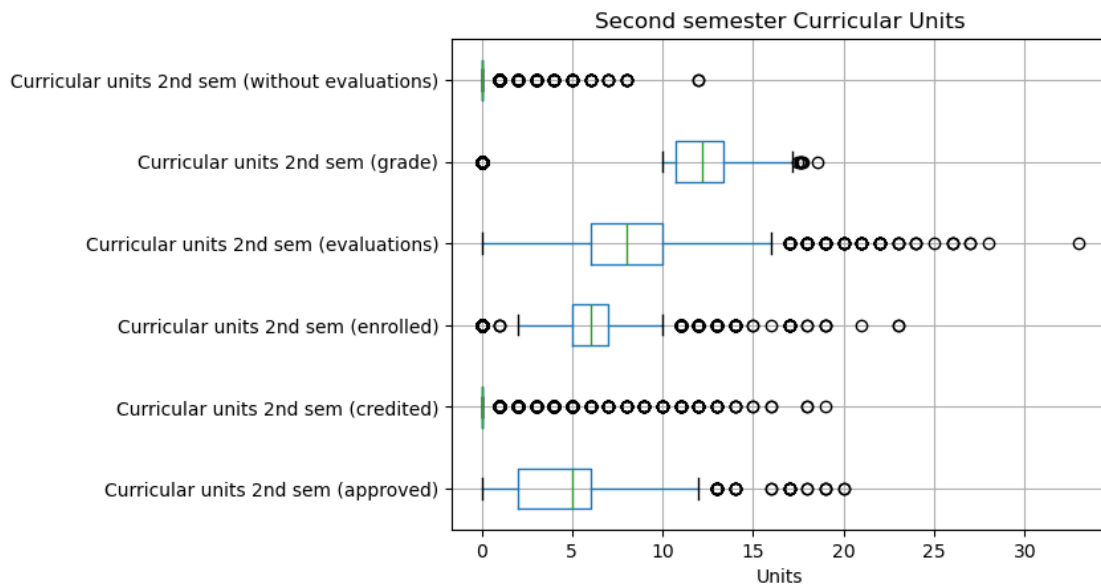
```
[12]: curriculum_units = [
          "Curricular units 1st sem (credited)",
          "Curricular units 2nd sem (credited)",
          "Curricular units 1st sem (enrolled)",
          "Curricular units 2nd sem (enrolled)",
          "Curricular units 1st sem (evaluations)",
          "Curricular units 2nd sem (evaluations)",
          "Curricular units 1st sem (approved)",
          "Curricular units 2nd sem (approved)",
          "Curricular units 1st sem (grade)",
          "Curricular units 2nd sem (grade)",
          "Curricular units 1st sem (without evaluations)",
          "Curricular units 2nd sem (without evaluations)",
      ]
      first_sem_curriculum = {unit for unit in curriculum_units if '1' in unit}
      second_sem_curriculum = set(curriculum_units) - first_sem_curriculum

      plot_semester_compare_bp(df, curriculum_units)
```

# Curricular Units between 1st and 2nd Semester

### credited Curricular units between 1st and 2nd sem



### enrolled Curricular units between 1st and 2nd sem



### evaluations Curricular units between 1st and 2nd sem



### approved Curricular units between 1st and 2nd sem



### grade Curricular units between 1st and 2nd sem



### without evaluations Curricular units between 1st and 2nd sem



```
[13]:  # will not plot first and semester comparison
       df[sorted(list(first_sem_curriculum))].boxplot(vert=False)
       plt.title("First semester Curricular Units")
       plt.xlabel("Units")
       plt.show()
```

```
df[sorted(list(second_sem_curriculum))].boxplot(vert=False)
plt.title("Second semester Curricular Units")
plt.xlabel("Units")
plt.show()
```
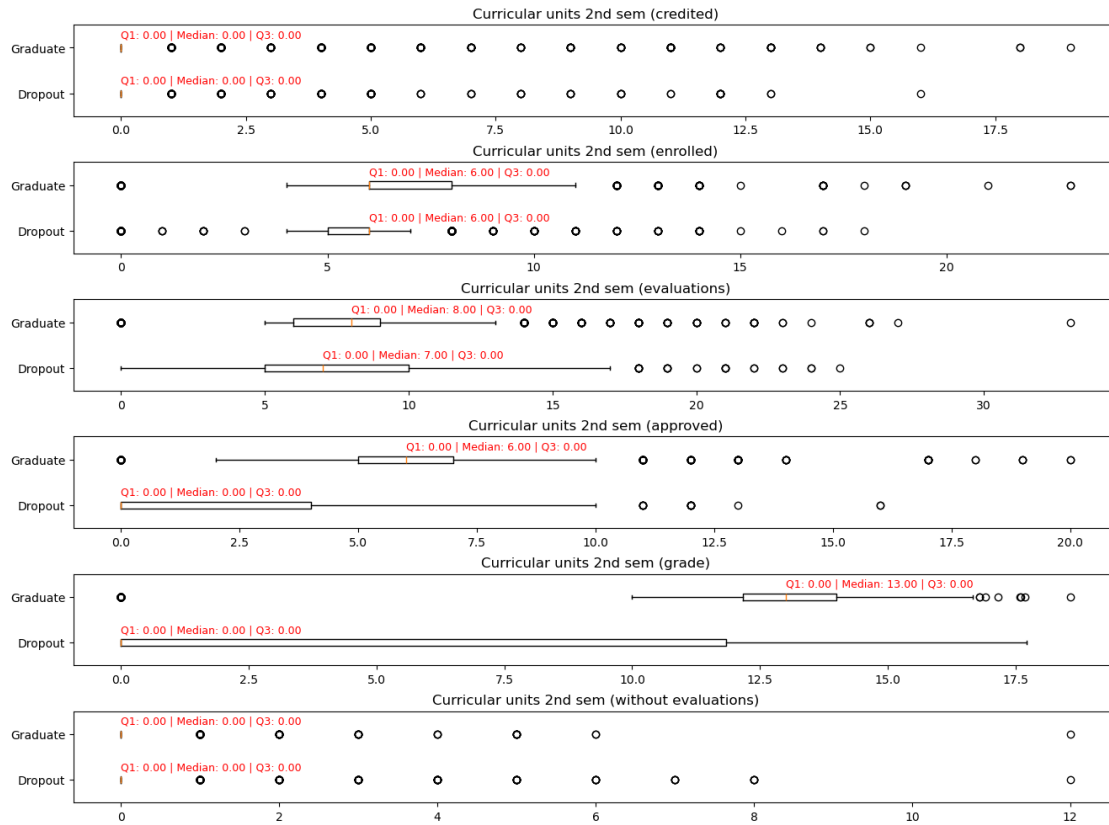


First semester Curricular Units



Second semester Curricular Units

[14]:
```
plot_grad_drop_compare_bp(df, curriculum_units, '1')
plot_grad_drop_compare_bp(df, curriculum_units, '2')
```

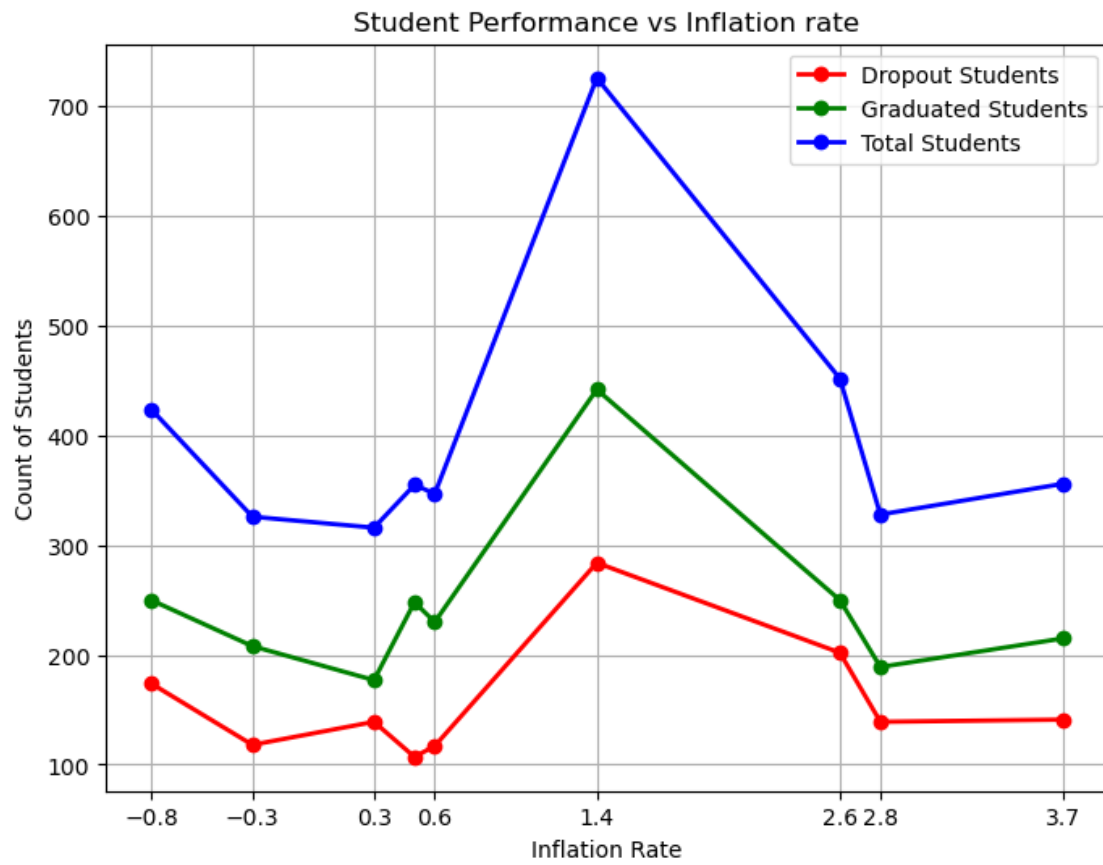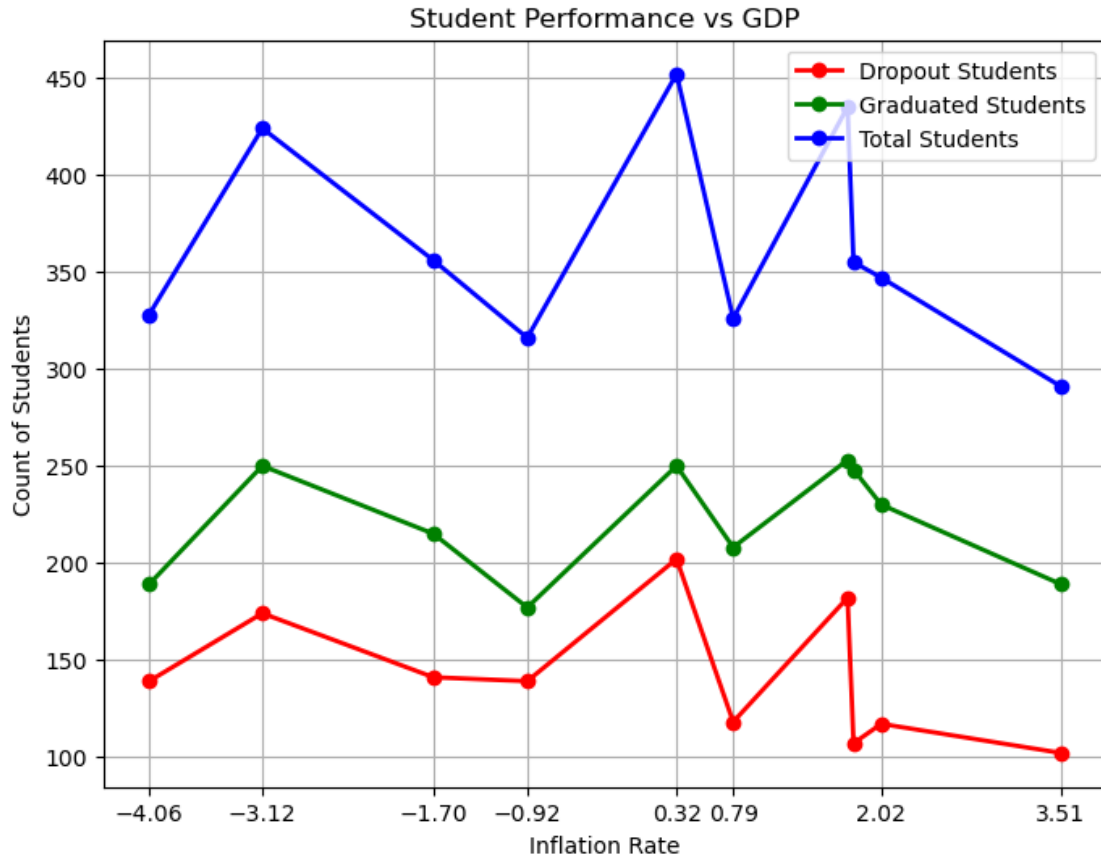# 1st Semester Curricular Units by credits

## Curricular units 1st sem (credited)

Graduate — Q1: 0.00 | Median: 0.00 | Q3: 0.00

Dropout — Q1: 0.00 | Median: 0.00 | Q3: 0.00

0.0    2.5    5.0    7.5    10.0    12.5    15.0    17.5    20.0

## Curricular units 1st sem (enrolled)

Graduate — Q1: 0.00 | Median: 6.00 | Q3: 0.00

Dropout — Q1: 0.00 | Median: 6.00 | Q3: 0.00

0    5    10    15    20    25

## Curricular units 1st sem (evaluations)

Graduate — Q1: 0.00 | Median: 8.00 | Q3: 0.00

Dropout — Q1: 0.00 | Median: 8.00 | Q3: 0.00

0    10    20    30    40

## Curricular units 1st sem (approved)

Graduate — Q1: 0.00 | Median: 6.00 | Q3: 0.00

Dropout — Q1: 0.00 | Median: 2.00 | Q3: 0.00

0    5    10    15    20    25

## Curricular units 1st sem (grade)

Graduate — Q1: 0.00 | Median: 13.00 | Q3: 0.00

Dropout — Q1: 0.00 | Median: 10.93 | Q3: 0.00

0.0    2.5    5.0    7.5    10.0    12.5    15.0    17.5

## Curricular units 1st sem (without evaluations)

Graduate — Q1: 0.00 | Median: 0.00 | Q3: 0.00

Dropout — Q1: 0.00 | Median: 0.00 | Q3: 0.00

0    2    4    6    8    10    12

13

## 2nd Semester Curricular Units by credits

### Curricular units 2nd sem (credited)



### Curricular units 2nd sem (enrolled)



### Curricular units 2nd sem (evaluations)



### Curricular units 2nd sem (approved)



### Curricular units 2nd sem (grade)



### Curricular units 2nd sem (without evaluations)



```
[15]: plot_line_graph(df, 'Inflation rate')
      plot_line_graph(df, 'GDP')
```

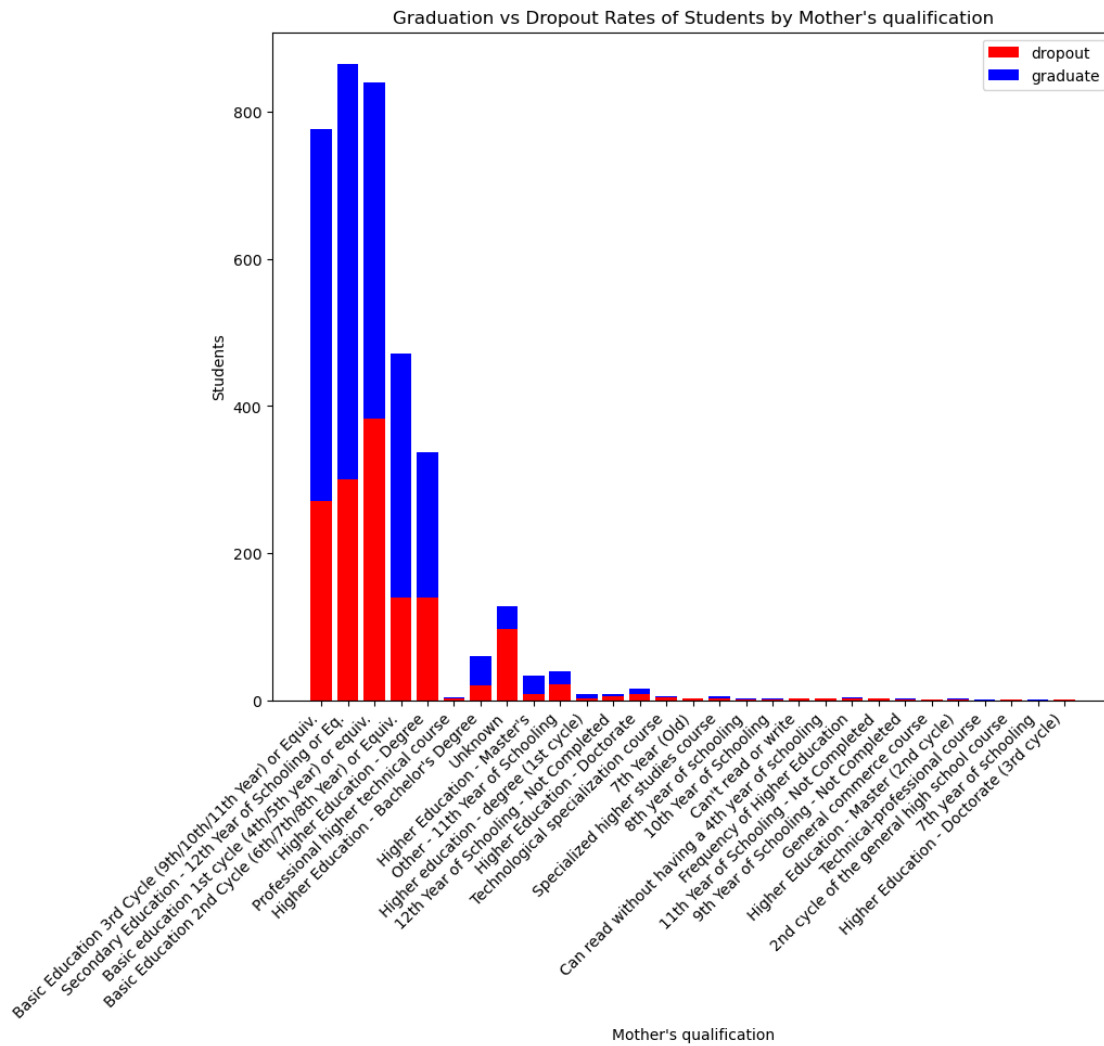Student Performance vs Inflation rate

## 3  Affect of Parental Education on Student Performance
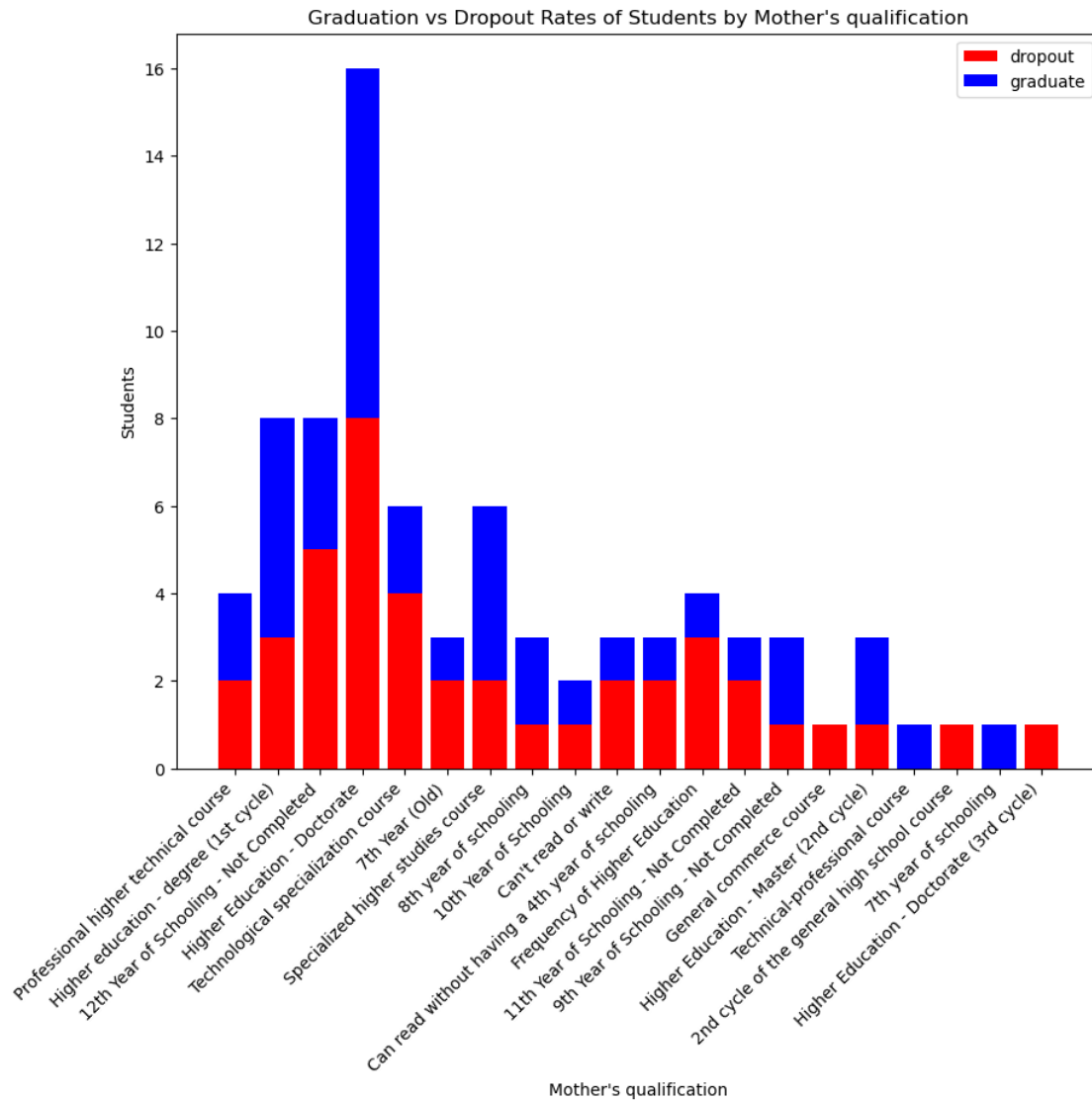
```
[16]: finished_df = df[df['dropout'].isin(['Dropout', 'Graduate'])]
      # filtering out sample categories
      mother_quals = finished_df['Mother\'s qualification'].value_counts()
      low_mother_quals = list(filter(lambda qual: mother_quals[qual] > 20,␣
        ↪mother_quals.keys()))
      high_mother_quals = list(filter(lambda qual: mother_quals[qual] < 20,␣
        ↪mother_quals.keys()))
      low_finished_df = finished_df.loc[~(finished_df['Mother\'s qualification'].
        ↪isin(low_mother_quals))]
      high_finished_df = finished_df.loc[~(finished_df['Mother\'s qualification'].
        ↪isin(high_mother_quals))]

      generate_general_stacked_bar_graph(finished_df, variable_map, 'Mother\'s␣
        ↪qualification')
      generate_general_stacked_bar_graph(low_finished_df, variable_map, 'Mother\'s␣
        ↪qualification')
```
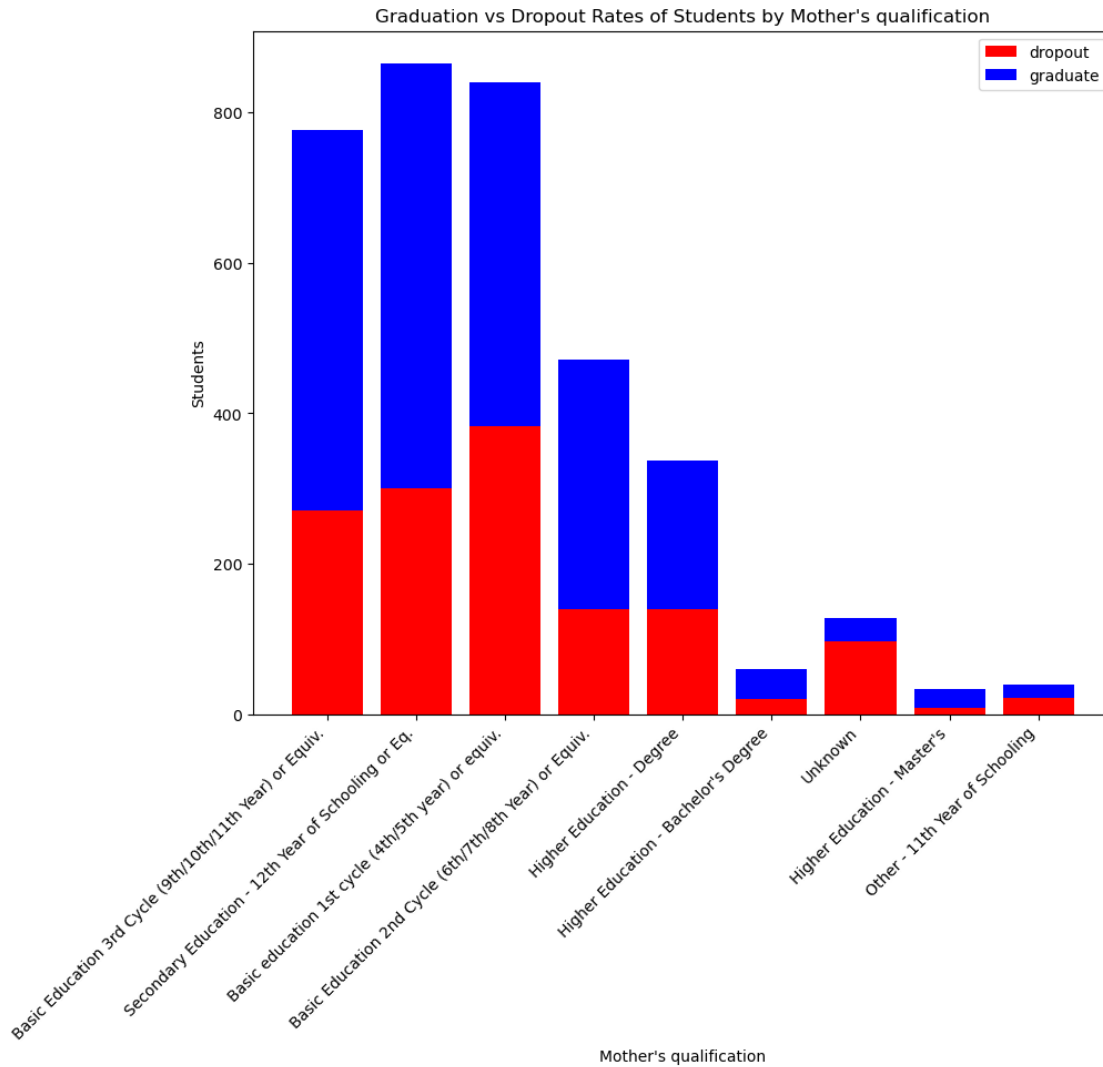
```
generate_general_stacked_bar_graph(high_finished_df, variable_map, 'Mother\'s␣
↪qualification')
```



Graduation vs Dropout Rates of Students by Mother's qualification

Graduation vs Dropout Rates of Students by Mother's qualification

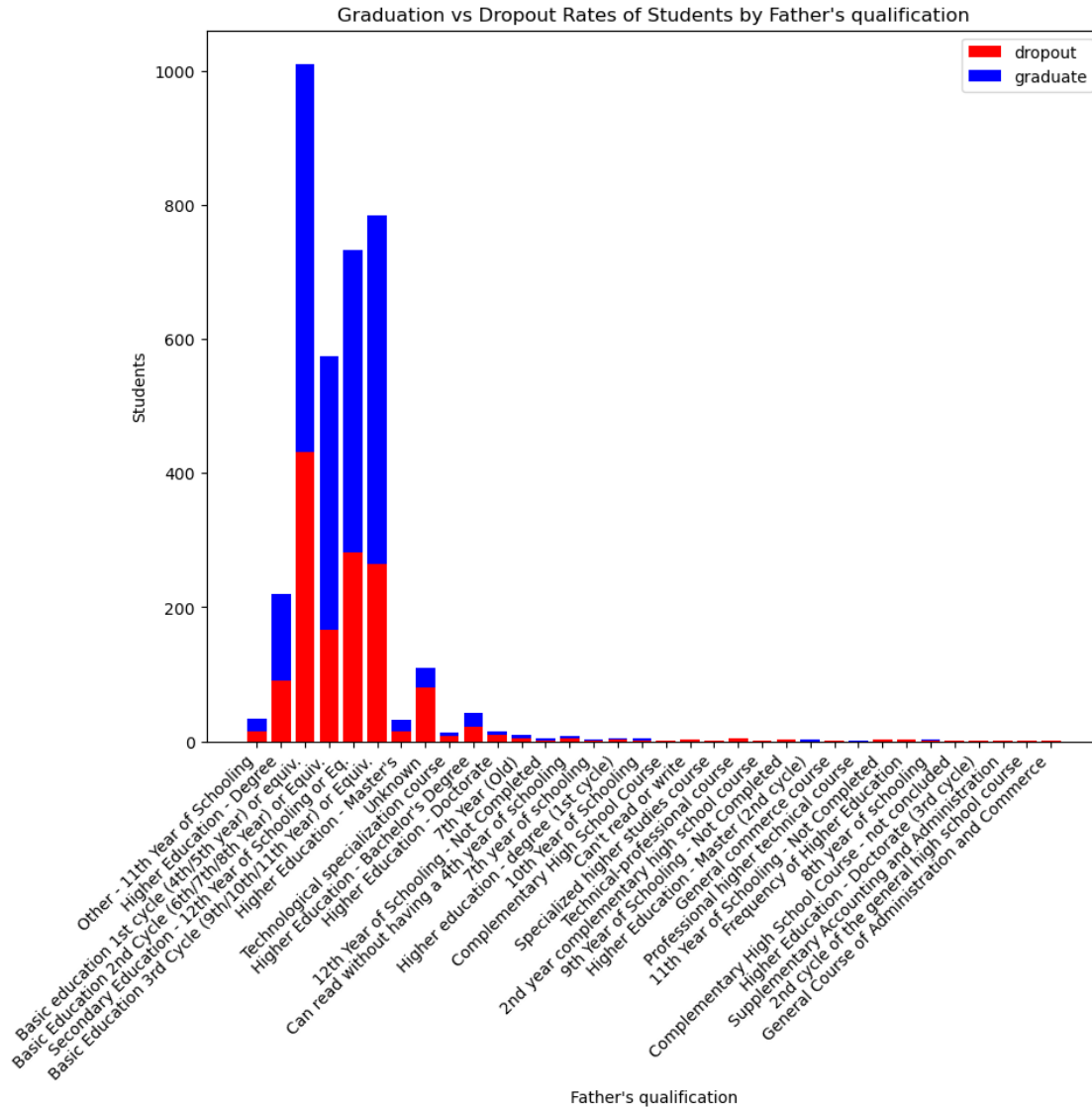Graduation vs Dropout Rates of Students by Mother's qualification

```
[17]: finished_df = df[df['dropout'].isin(['Dropout', 'Graduate'])]
      # filtering out sample categories
      father_quals = finished_df['Father\'s qualification'].value_counts()
      low_mother_quals = list(filter(lambda qual: father_quals[qual] > 20,␣
        ↪father_quals.keys()))
      high_mother_quals = list(filter(lambda qual: father_quals[qual] < 20,␣
        ↪father_quals.keys()))
      low_finished_df = finished_df.loc[~(finished_df['Father\'s qualification'].
        ↪isin(low_mother_quals))]
      high_finished_df = finished_df.loc[~(finished_df['Father\'s qualification'].
        ↪isin(high_mother_quals))]

      generate_general_stacked_bar_graph(finished_df, variable_map, 'Father\'s␣
        ↪qualification')
```
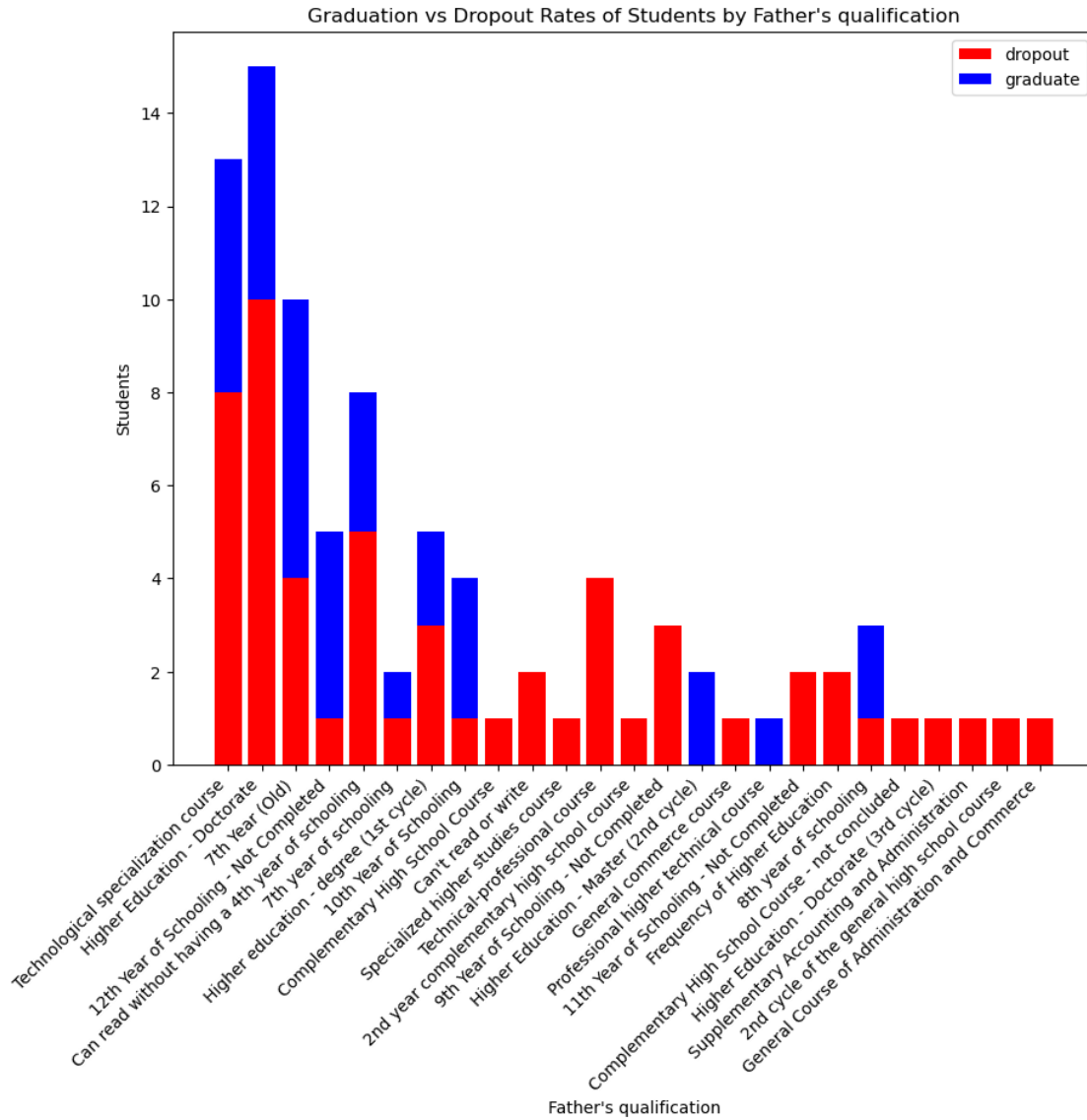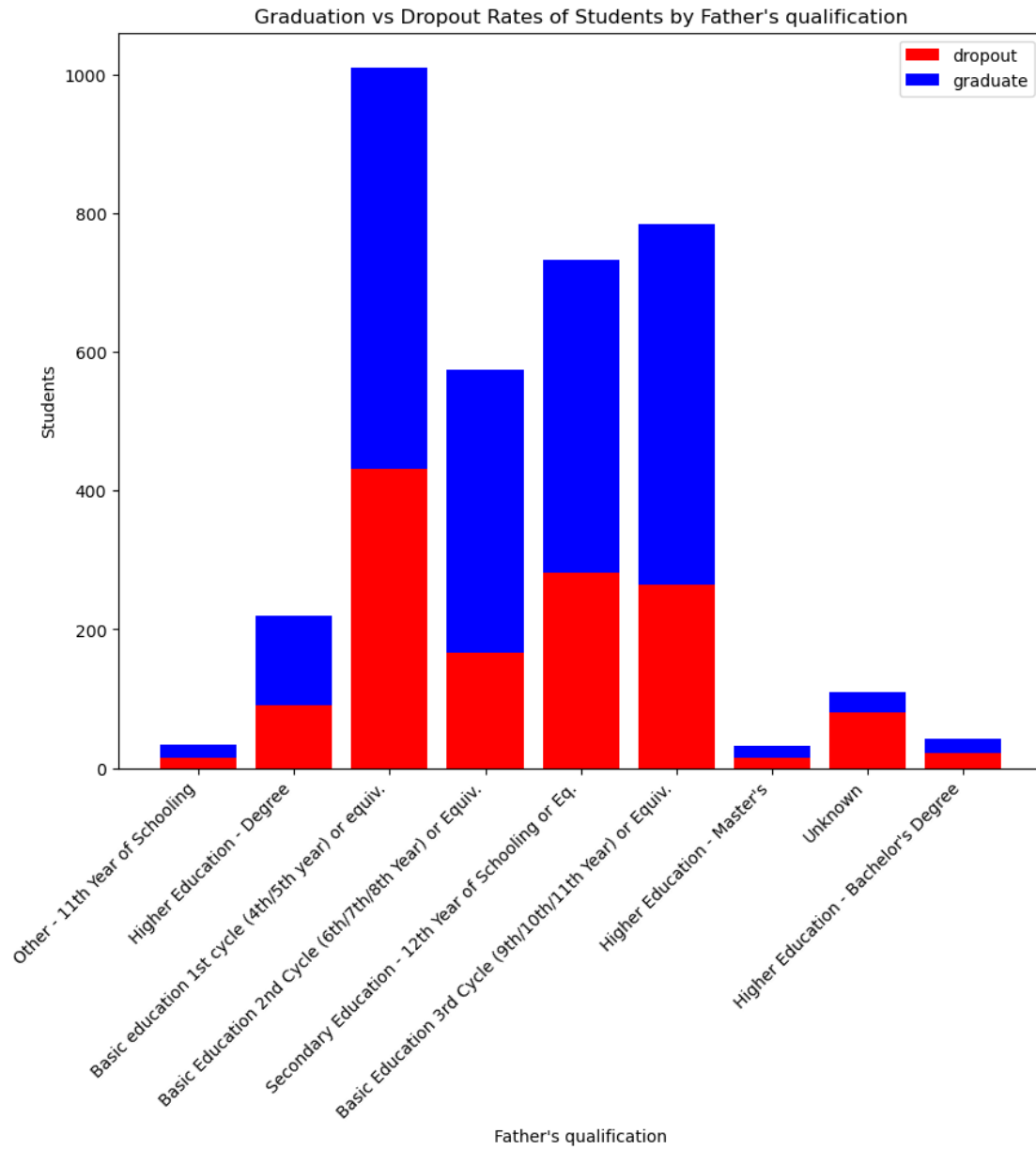
```
generate_general_stacked_bar_graph(low_finished_df, variable_map, 'Father\'s␣
␣→qualification')
generate_general_stacked_bar_graph(high_finished_df, variable_map, 'Father\'s␣
␣→qualification')
```



Graduation vs Dropout Rates of Students by Father's qualification

Graduation vs Dropout Rates of Students by Father's qualification

Graduation vs Dropout Rates of Students by Father's qualification