



EDM PROJECT

WWW.EDM.COM

CONTENT

- 01** ABOUT US
- 02** INTRO OF FIRST DATASET
- 03** ANALYSIS
- 04** CONCLUSION
- 05** INTRO OF SECOND DATASET
- 06** FEATURE DESCRIPTION
- 07** ANALYSIS
- 08** CONCLUSION

ABOUT US



- 資工三 110590003 黃政 coding
- 資工三 110590005 蕭耕宏 coding, lead
- 資工三 110590028 黃冠鈞 coding
- 資工三 110590034 楊榮鈞 coding



INTRO OF FIRST DATESET

IBM HR Analytics Employee Attrition & Performance

Goal : predict education level

- Below College
- College
- Bachelor
- Master
- Doctor

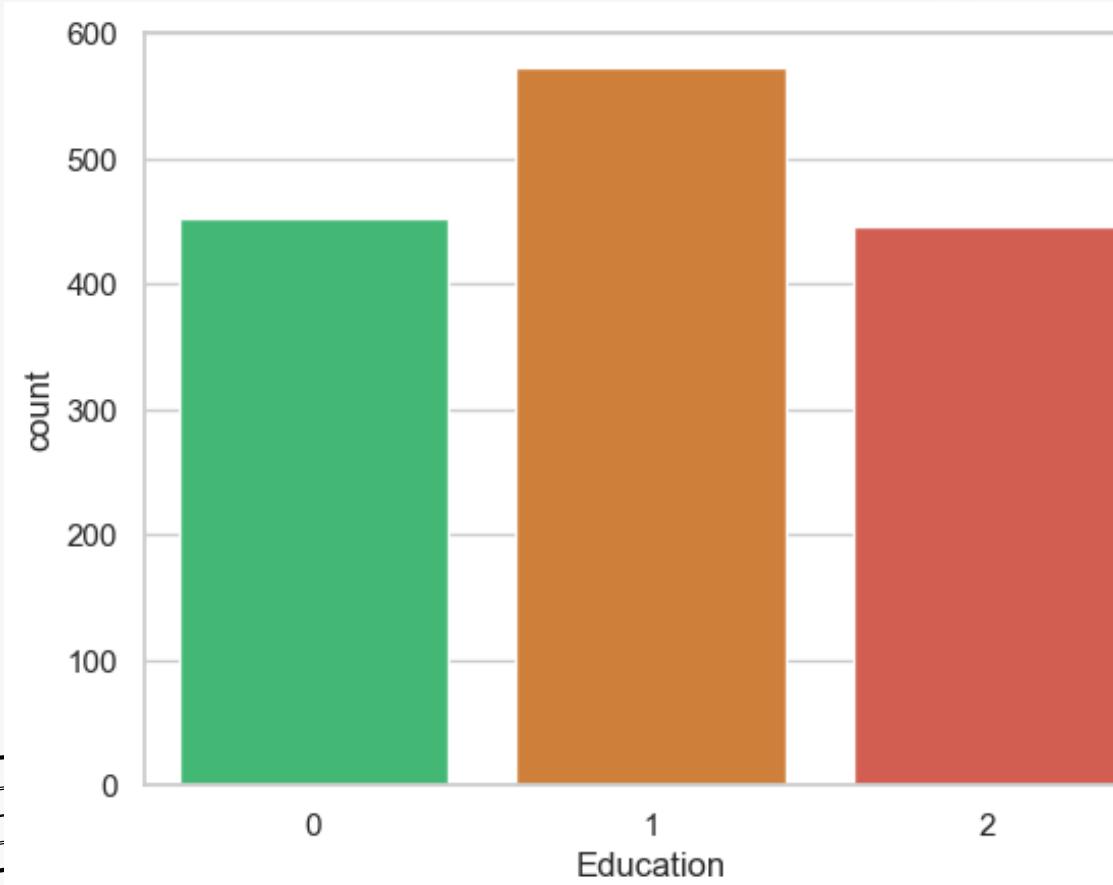
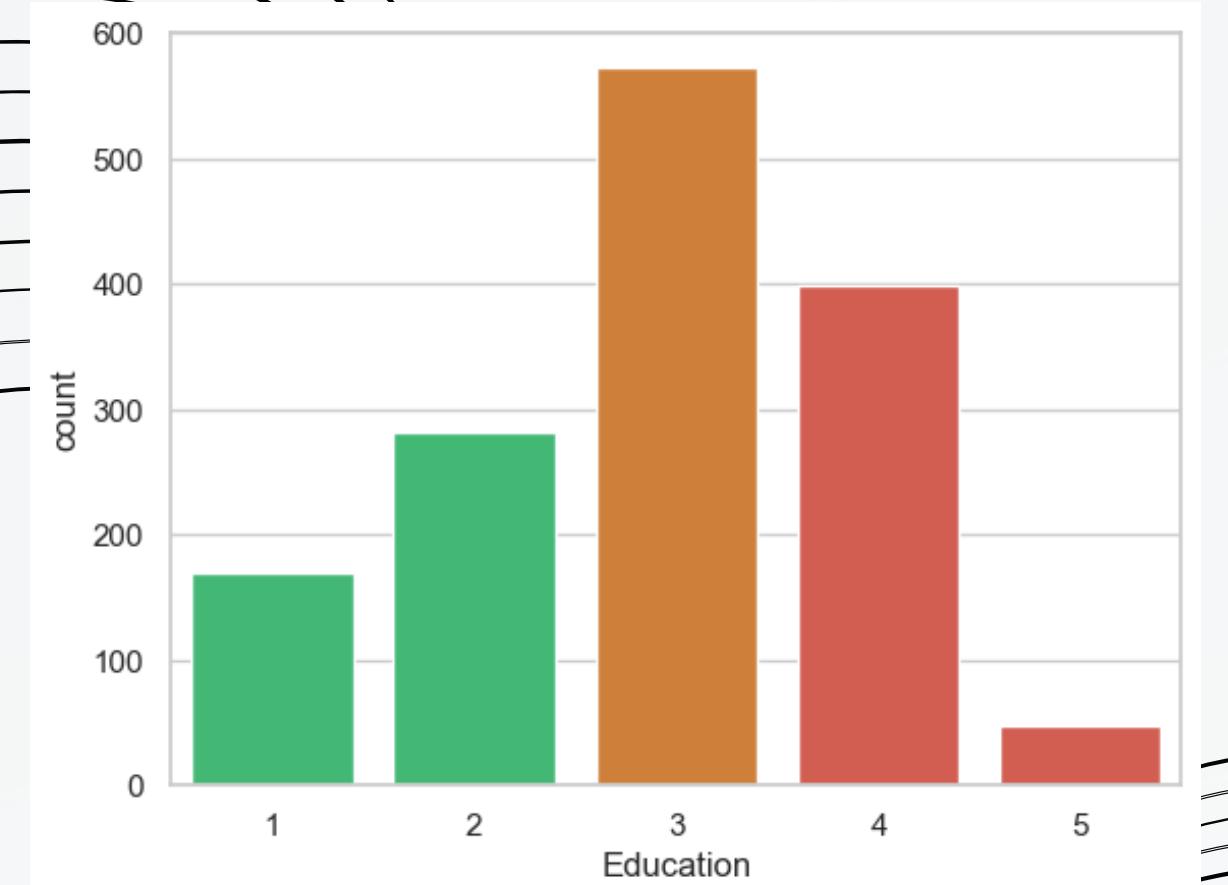
Age, Attrition, BusinessTravel, DailyRate,
Department, DistanceFromHome,
Education, EducationField, EmployeeCount,
EmployeeNumber, EnvironmentSatisfaction,
Gender, HourlyRate, JobInvolvement,
JobLevel, JobRole, JobSatisfaction,
MaritalStatus, MonthlyIncome, MonthlyRate,
NumCompaniesWorked, Over18, OverTime,
PercentSalaryHike, PerformanceRating,
RelationshipSatisfaction, StandardHours,
StockOptionLevel, TotalWorkingYears,
TrainingTimesLastYear, WorkLifeBalance,
YearsAtCompany, YearsInCurrentRole,
YearsSinceLastPromotion,
YearsWithCurrManager

INTRO OF FIRST DATASET

IBM HR Analytics Employee Attrition & Performance

simply the prediction task

- Below College 1
- College 2
- Bachelor 3
- Master 4
- Doctor 5



INTRO OF FIRST DATESET

IBM HR Analytics Employee Attrition & Performance

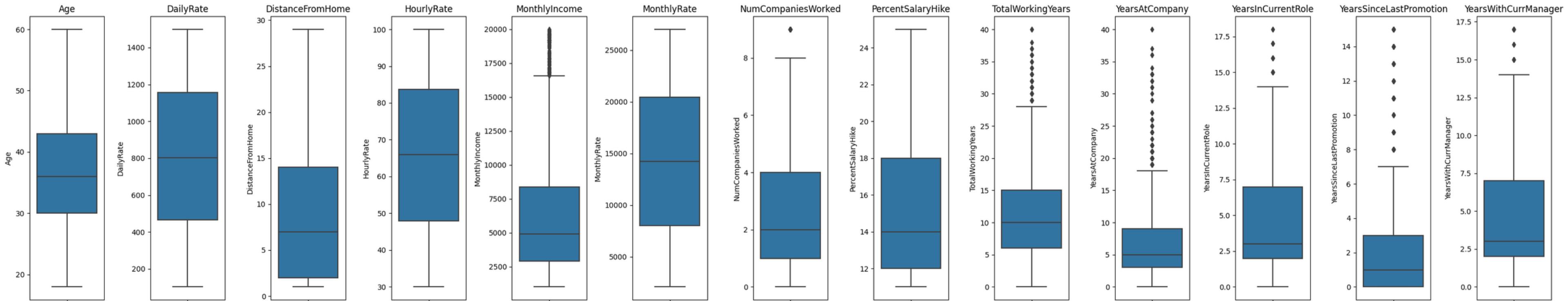
Issues to solve

- some variables contains single value
- some ordinal variables is encoded to integer and some don't
- some numerical variables have extreme outlier

INTRO OF FIRST DATASET

IBM HR Analytics Employee Attrition & Performance

numerical variables outliers



INTRO OF FIRST DATESET

IBM HR Analytics Employee Attrition & Performance

- Solution**
- categorical variables : use one-hot encoder
 - ordinal variables : use ordinal encoder
 - numerical variables : remove outlier and standardize features
removing the mean and scaling to unit variance
 - use several model to predict
 - implement k-fold validation

ANALYSIS

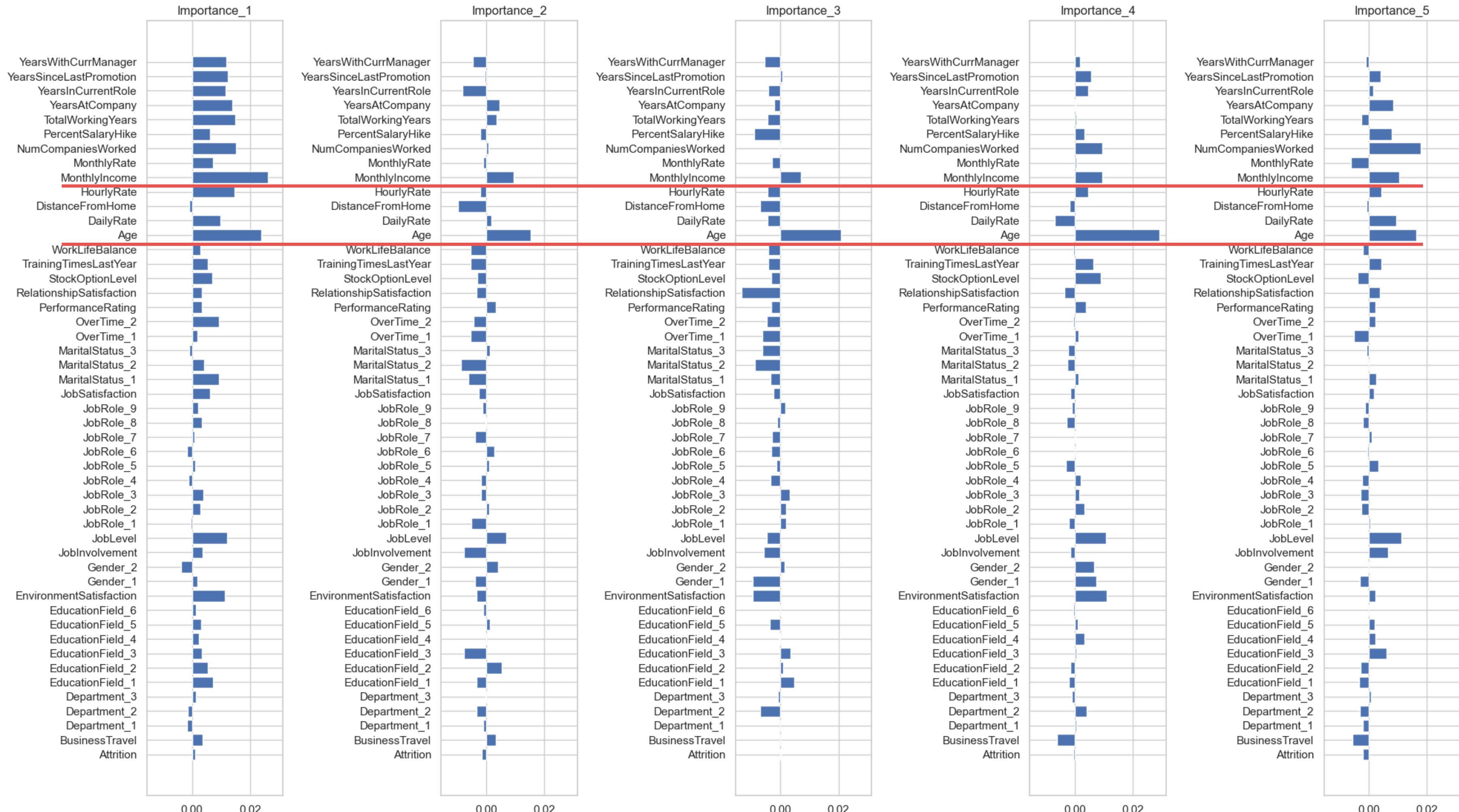
| model | acc | precision_score | recall_score | f1_score |
|----------------------------|-------|-----------------|--------------|----------|
| GradientBoostingClassifier | 0.406 | 0.4 | 0.406 | 0.396 |
| RandomForestClassifier | 0.421 | 0.42 | 0.421 | 0.413 |
| SVC | 0.386 | 0.365 | 0.386 | 0.385 |
| SVC_rbf | 0.388 | 0.365 | 0.388 | 0.377 |
| KNeighborsClassifier | 0.354 | 0.366 | 0.354 | 0.348 |
| XGBClassifier | 0.404 | 0.399 | 0.404 | 0.399 |

very poor result

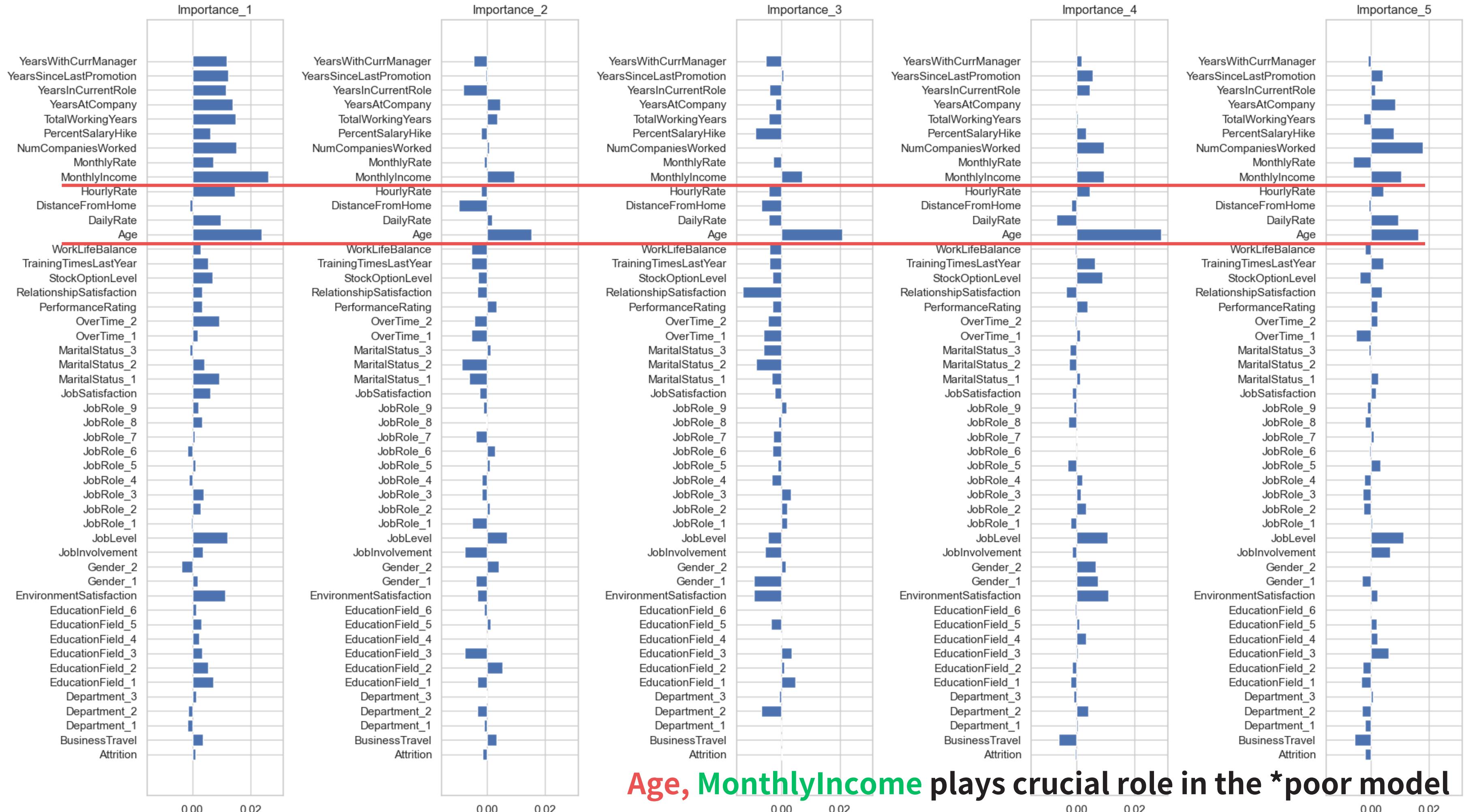
ANALYSIS

Use Permutation feature importance To find
features for RFC model

Feature Importance Comparison



Feature Importance Comparison



CONCLUSION



ADDITIONAL ANALYSIS

predict Attrition?

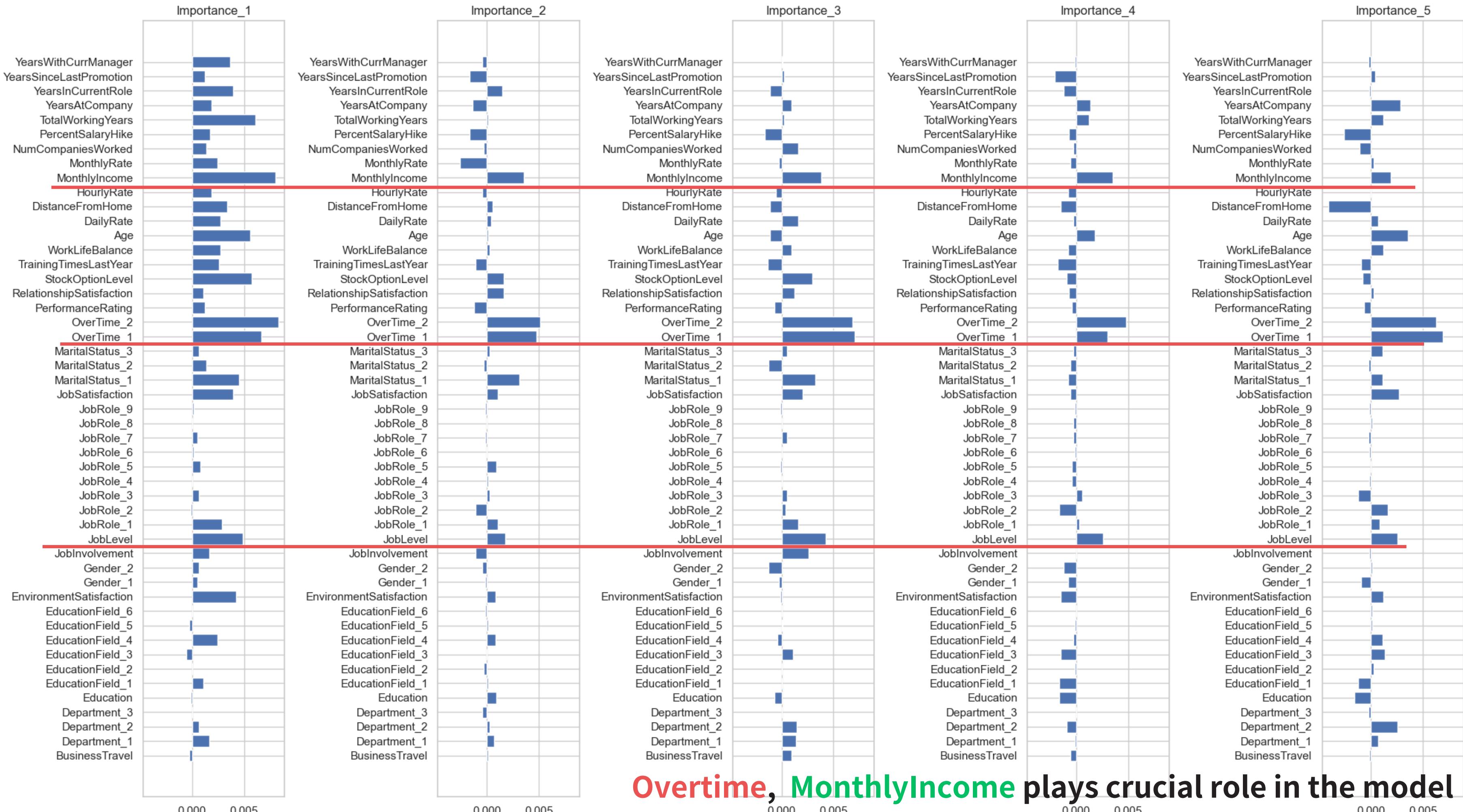
| model | acc | precision_score | recall_score | f1_score |
|----------------------------|-------|-----------------|--------------|----------|
| GradientBoostingClassifier | 0.865 | 0.848 | 0.865 | 0.843 |
| RandomForestClassifier | 0.859 | 0.849 | 0.859 | 0.82 |
| SVC | 0.867 | 0.873 | 0.867 | 0.83 |
| SVC_rbf | 0.867 | 0.873 | 0.867 | 0.83 |
| ExtraTreesClassifier | 0.859 | 0.844 | 0.859 | 0.826 |
| KNeighborsClassifier | 0.844 | 0.815 | 0.844 | 0.797 |
| XGBClassifier | 0.868 | 0.854 | 0.868 | 0.848 |

not bad!!!

ADDITIONAL ANALYSIS

Use Permutation feature importance To find features for RFC model

Feature Importance Comparison



Overtime, MonthlyIncome plays crucial role in the model



the result is unexpectedly bad,
so we find another dataset

INTRO OF SECOND DATASET

Student Stress Factors: A Comprehensive Analysis

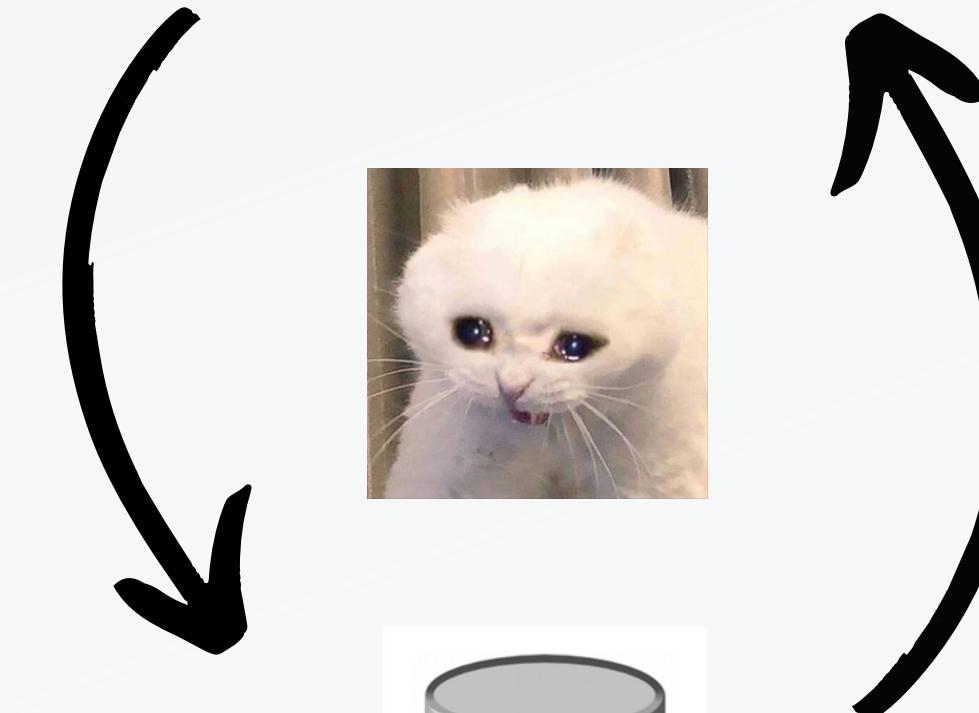
Environmental
Factors

Psychological
Factors

Academic Factors

Social Factor

Physiological Factors



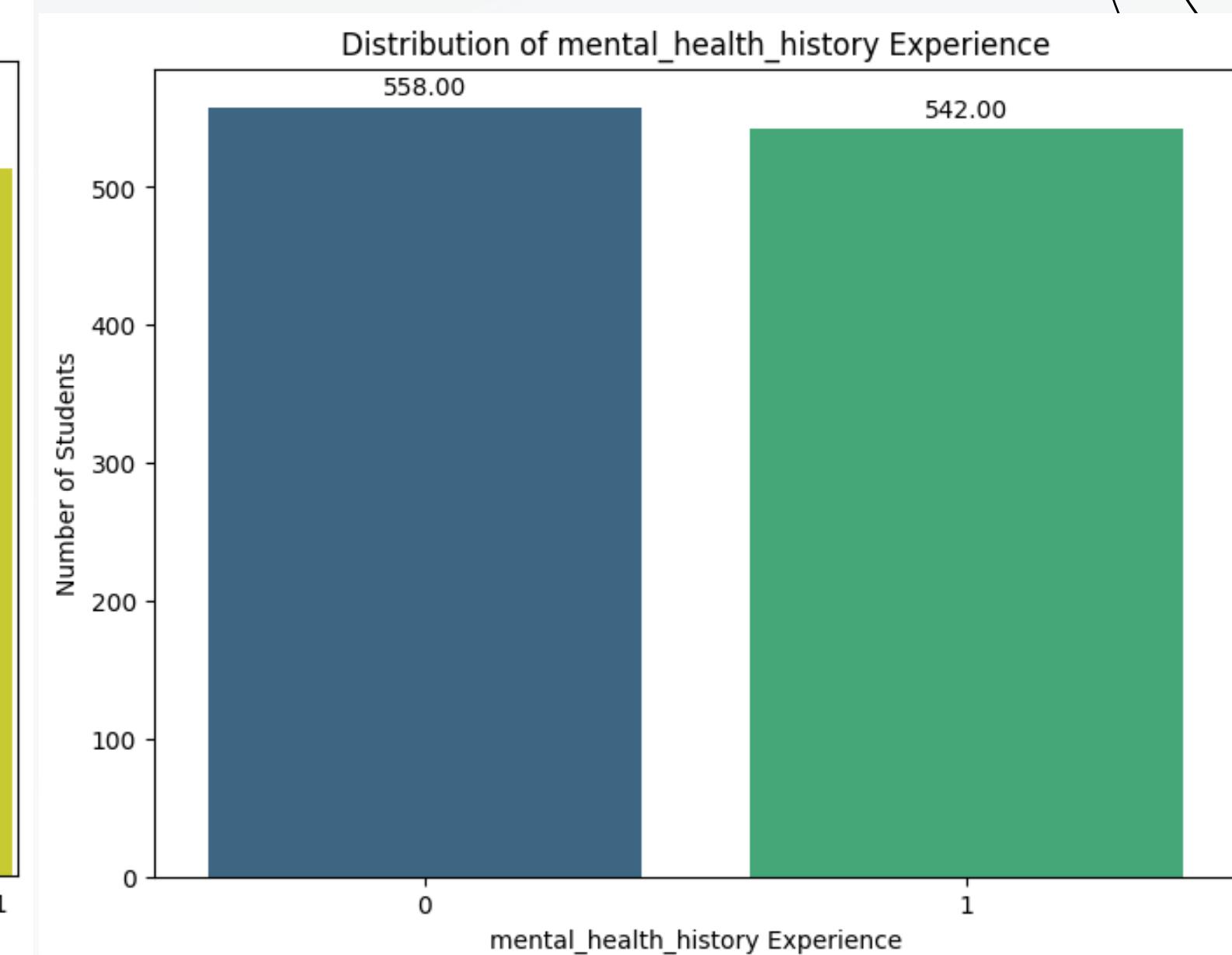
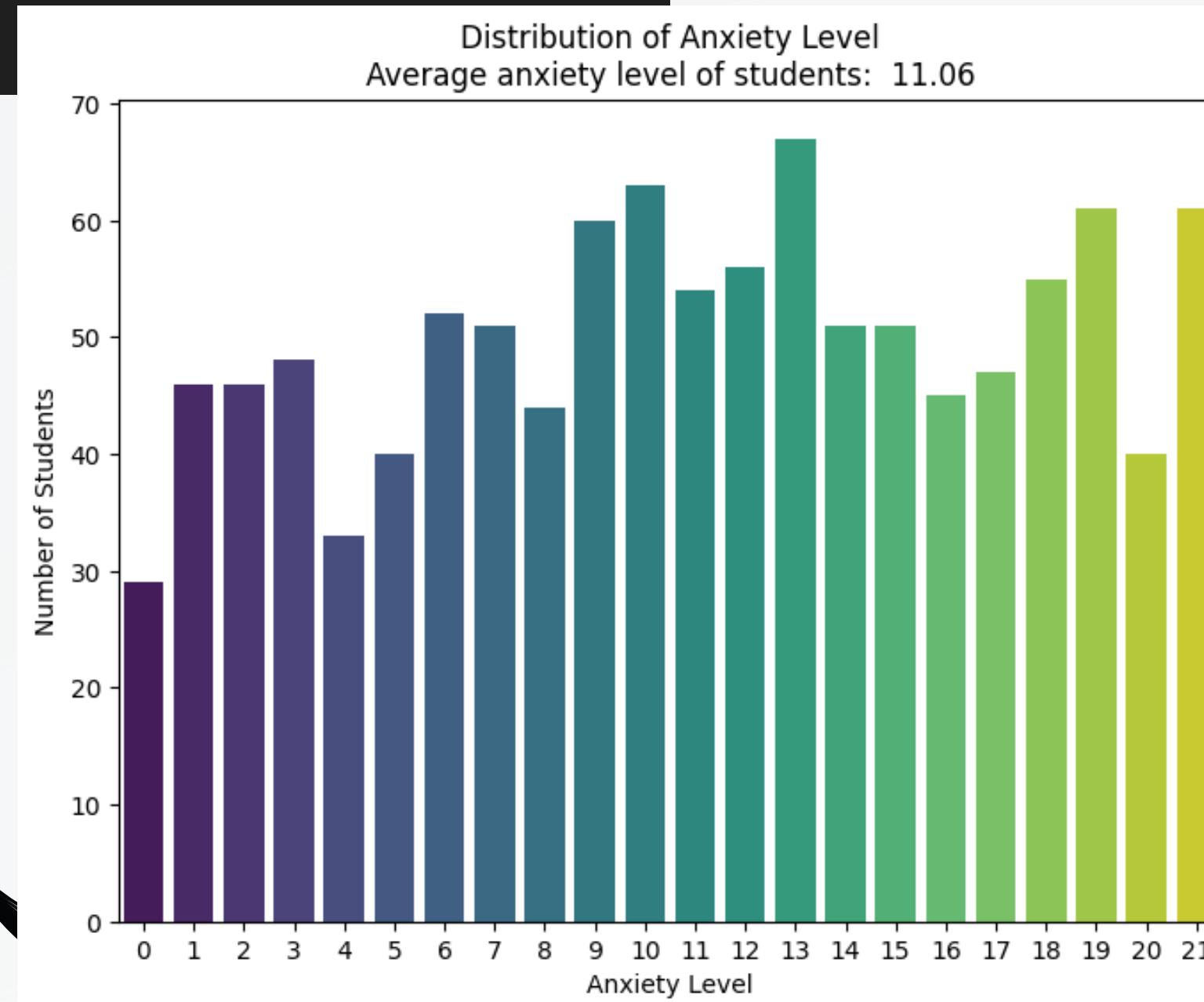
FEATURE DESCRIPTION

Descriptive Statistics

```
total_students = len(data)  
print('total students: ' + str(total_students))
```

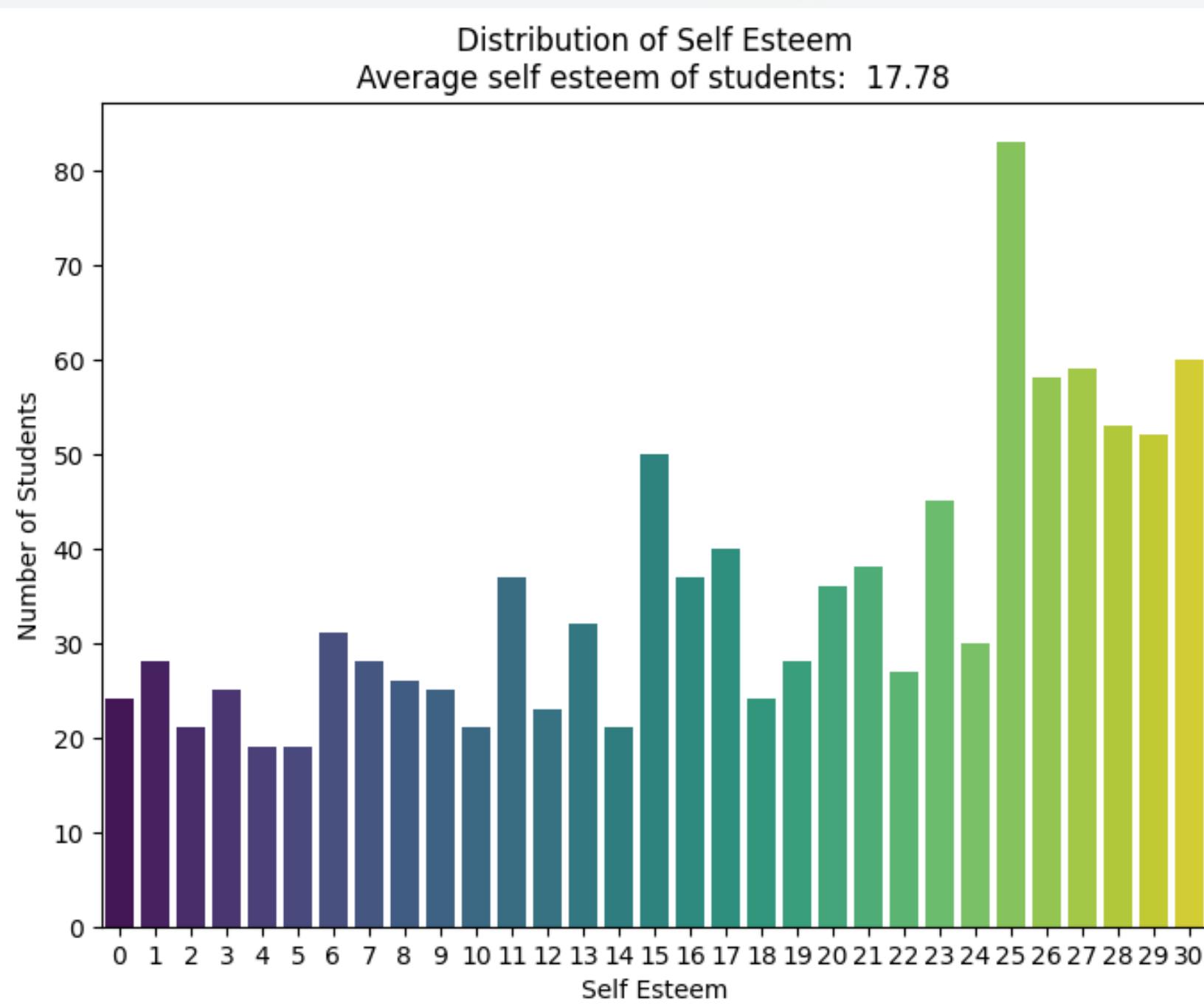
✓ 0.0s

total students: 1100



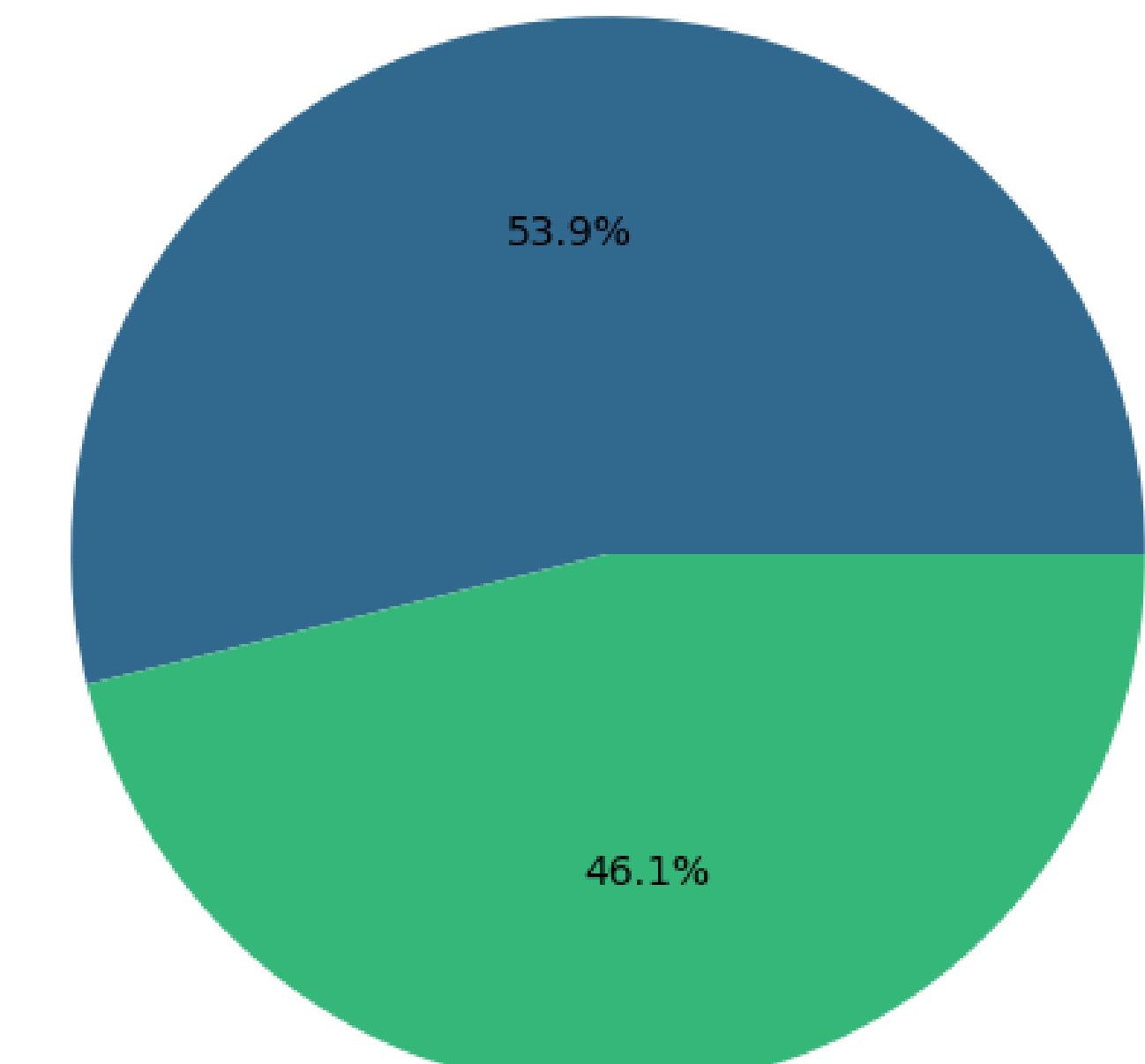
FEATURE DESCRIPTION

PSYCHOLOGICAL FACTORS



Distribution of Self Esteem

Above Average

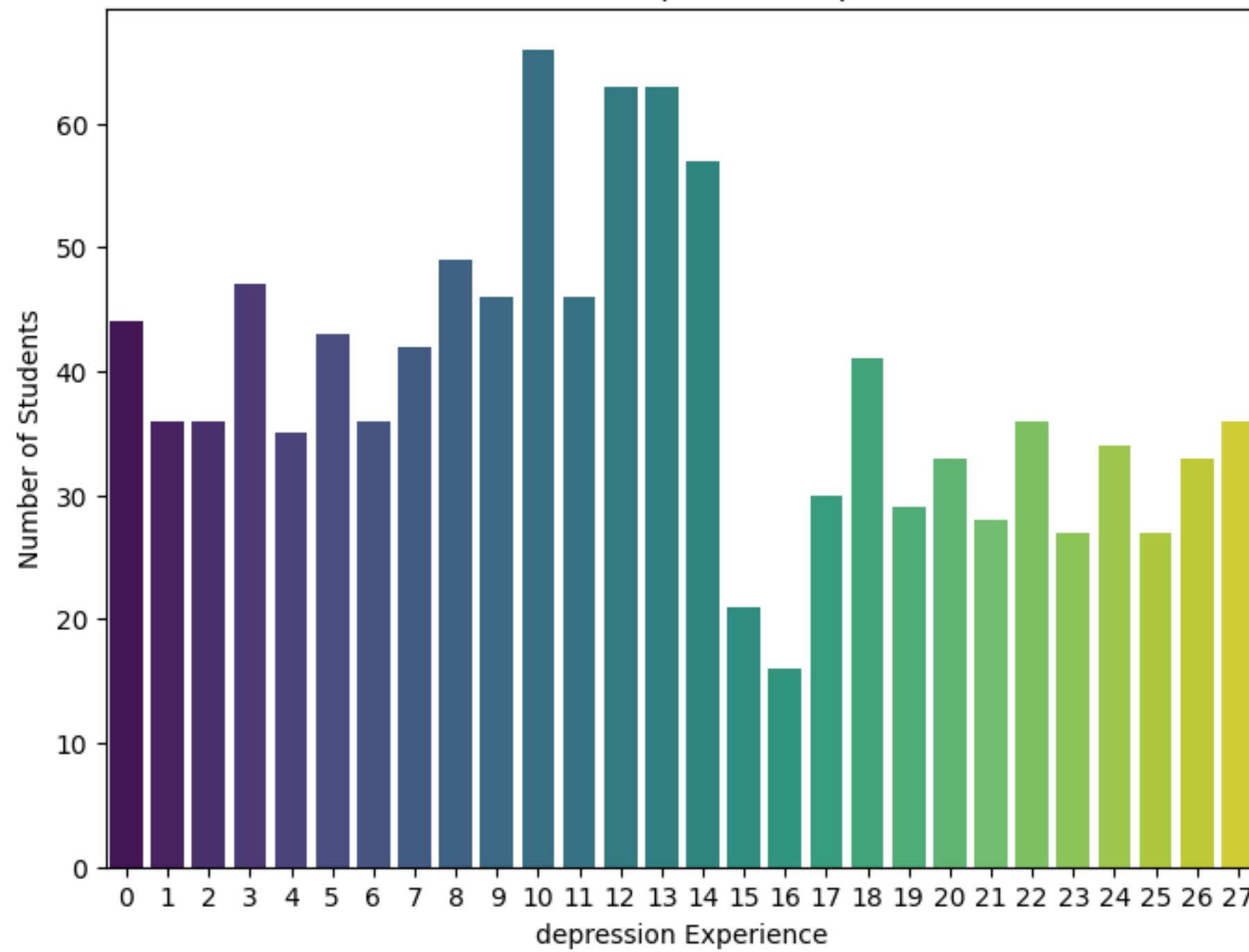


Below or Equal to Average

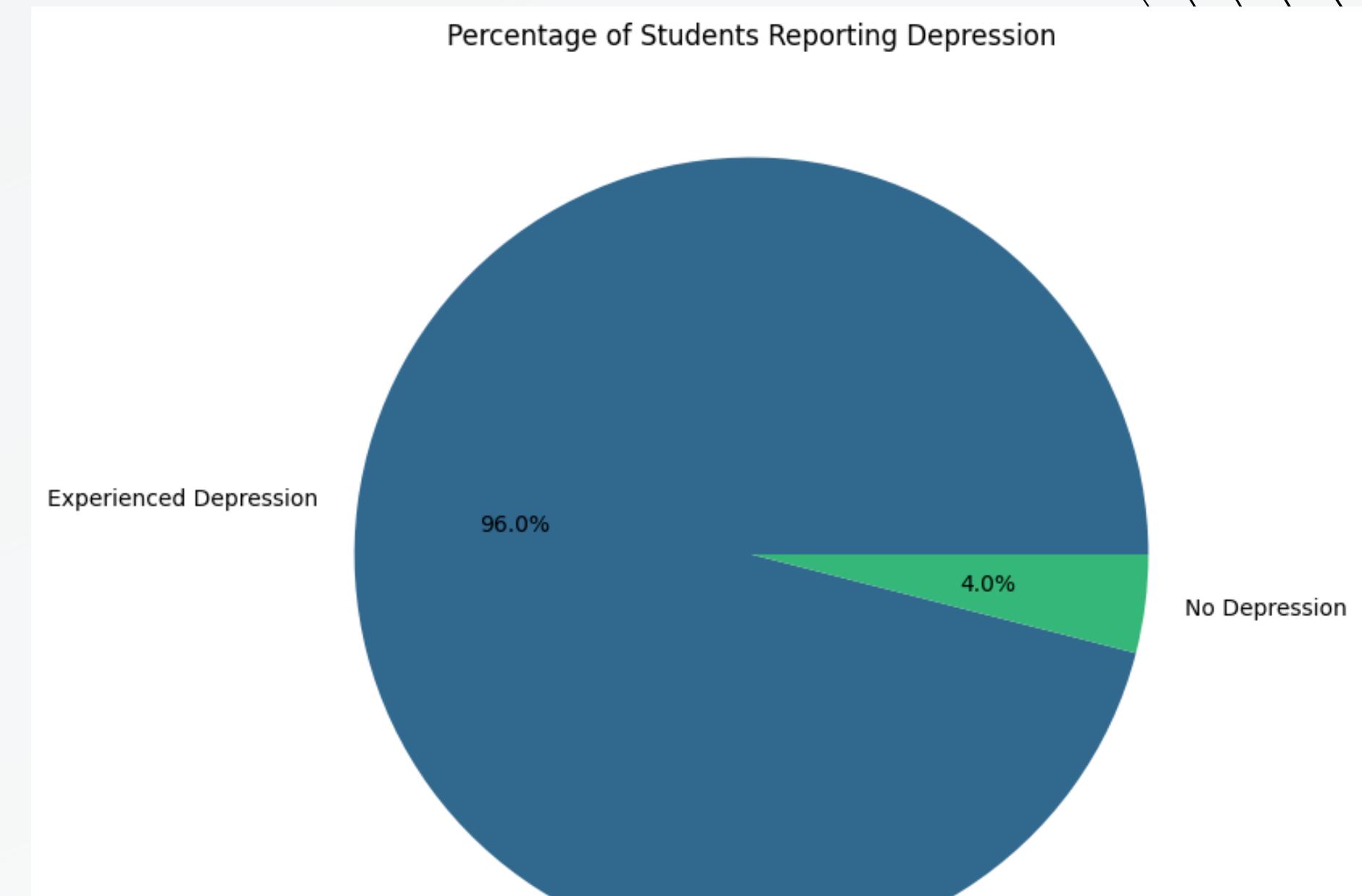
FEATURE DESCRIPTION

PSYCHOLOGICAL FACTORS

Distribution of depression Experience

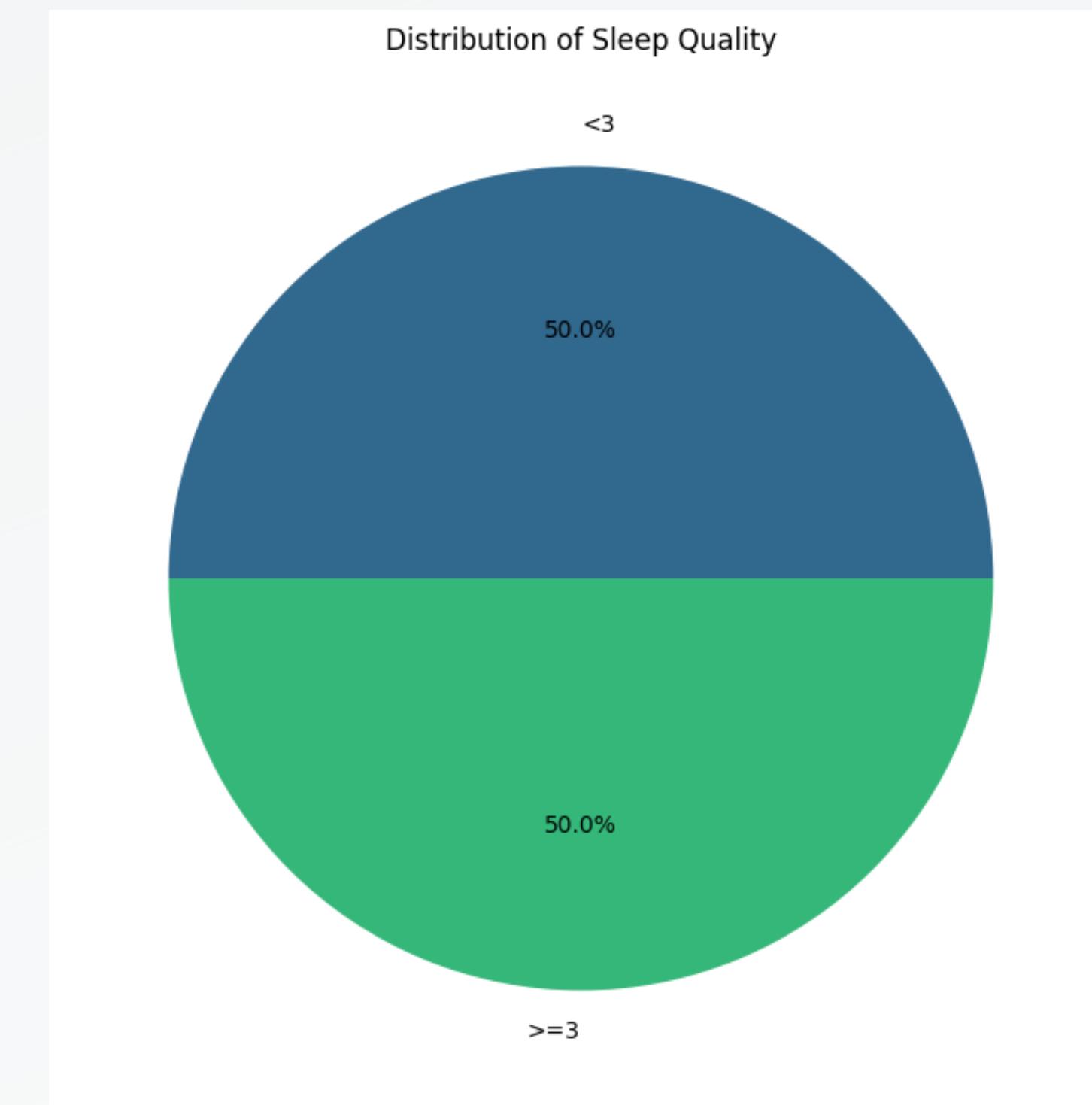
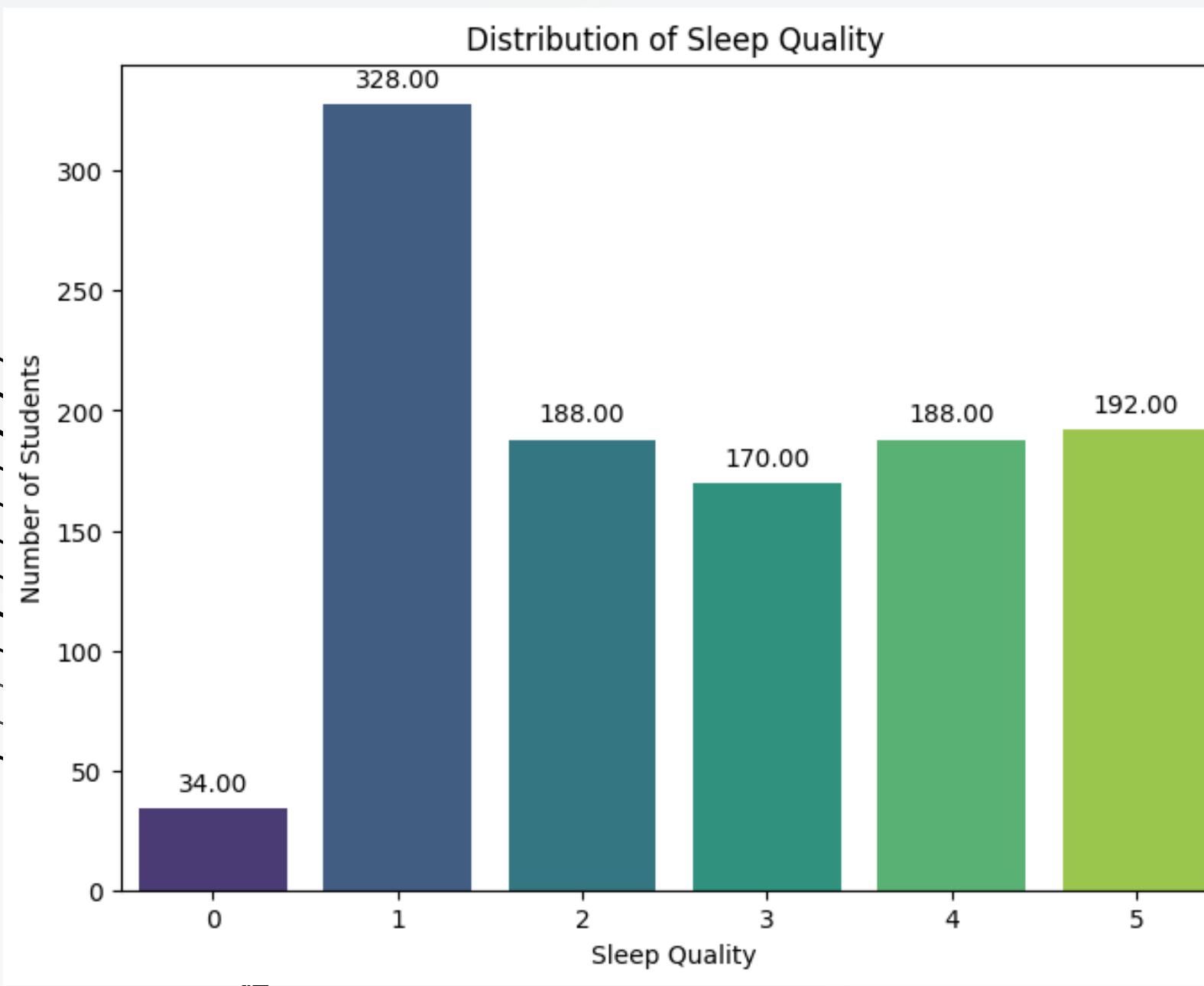


Percentage of Students Reporting Depression



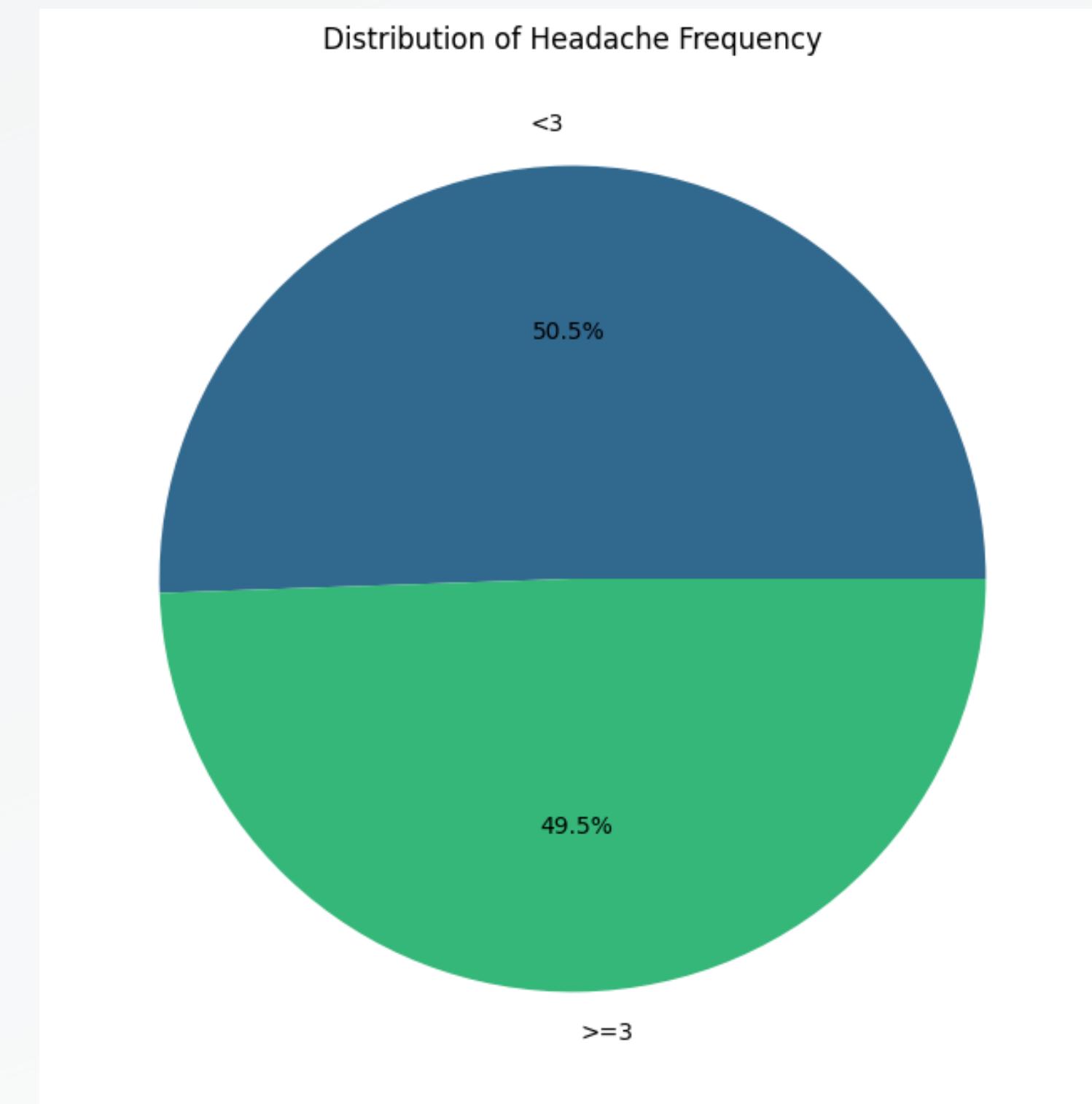
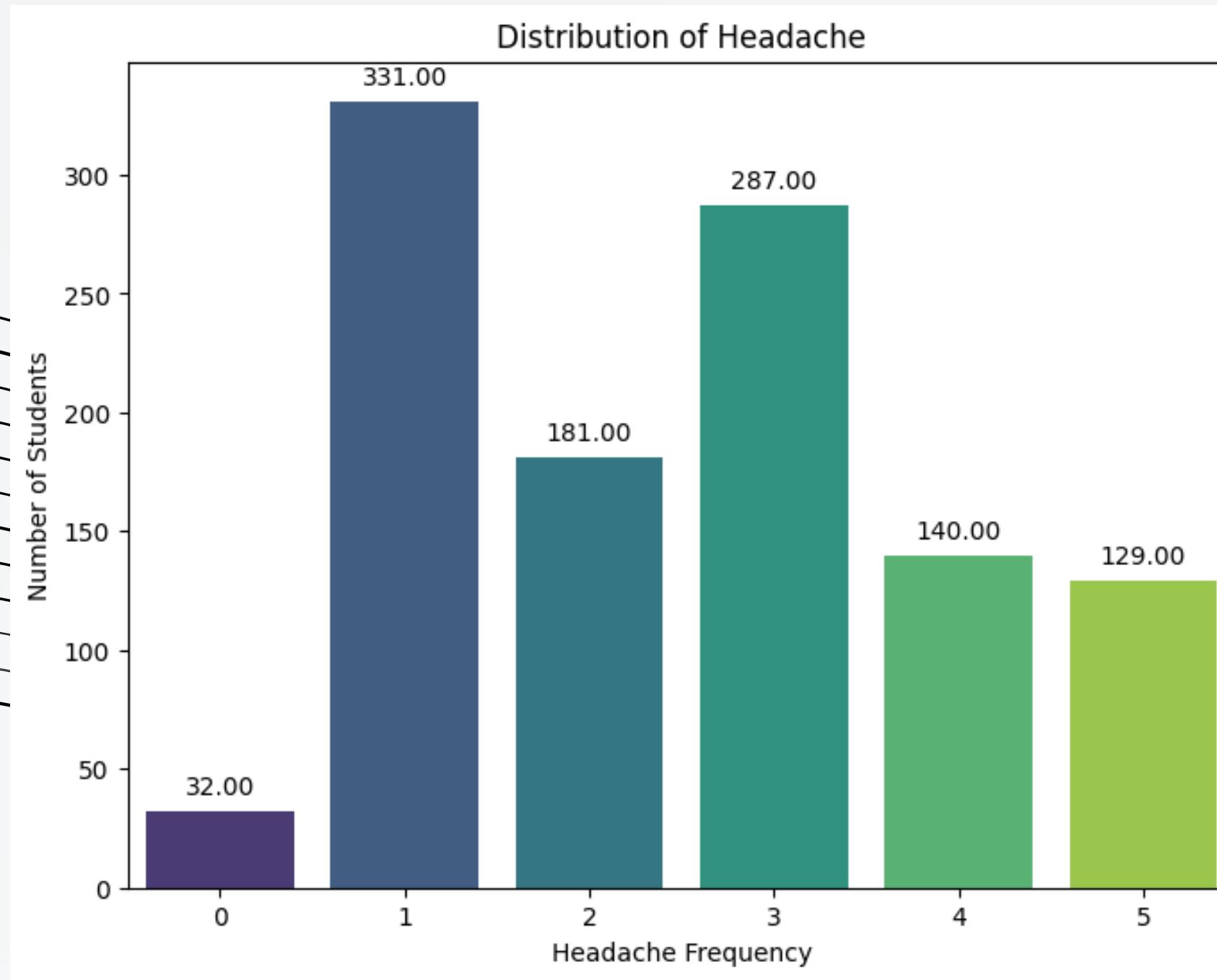
FEATURE DESCRIPTION

PHYSIOLOGICAL FACTORS



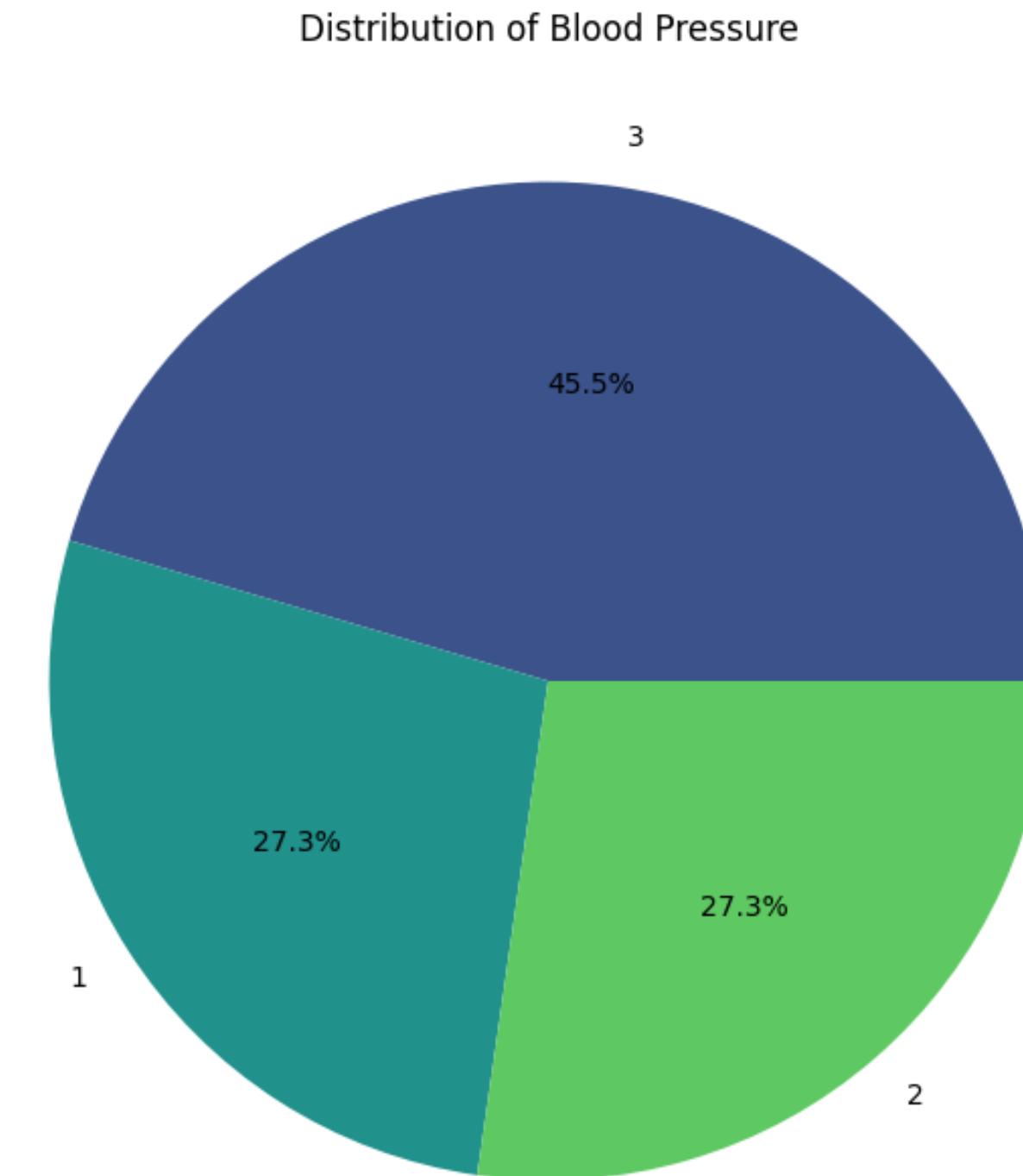
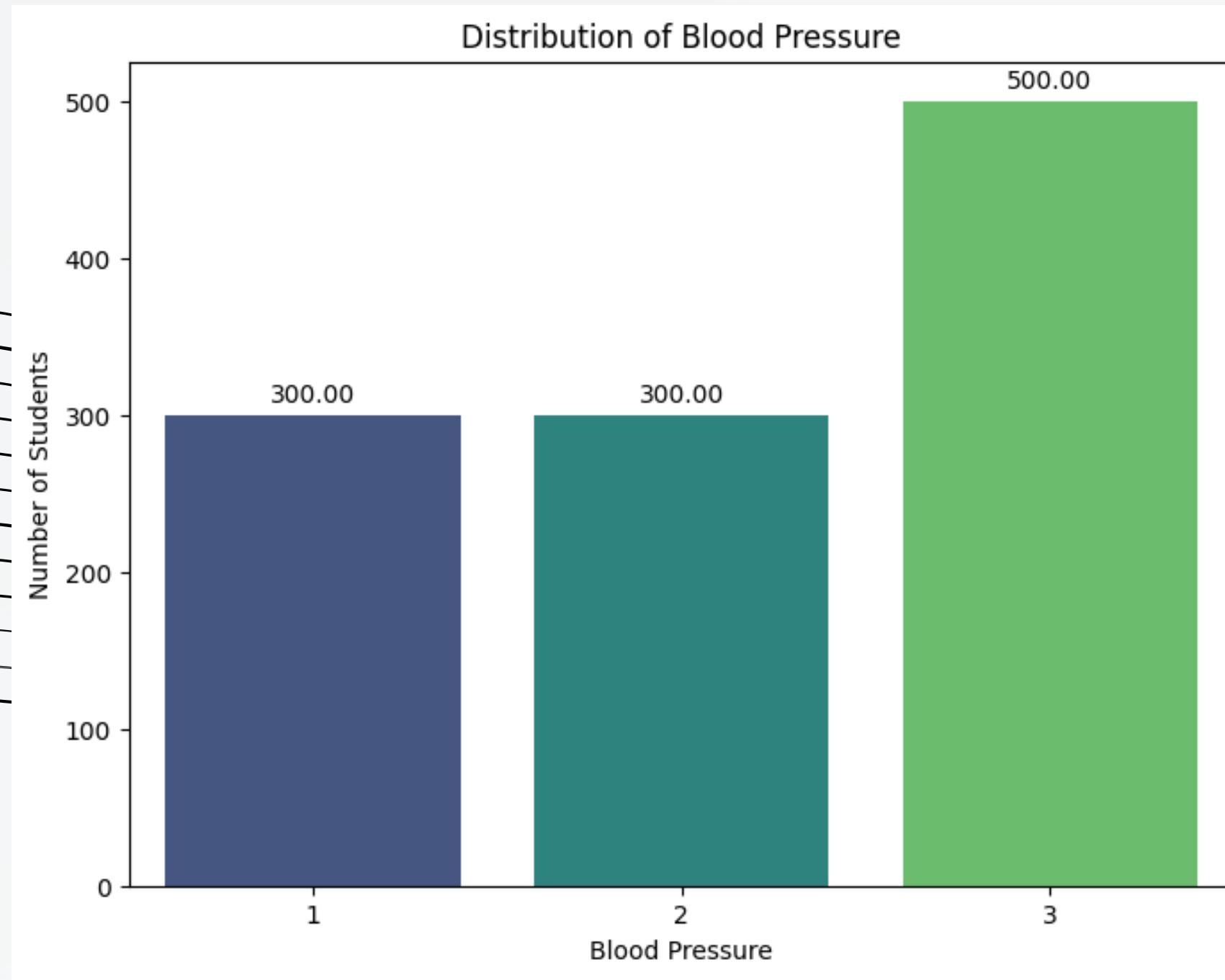
FEATURE DESCRIPTION

PHYSIOLOGICAL FACTORS



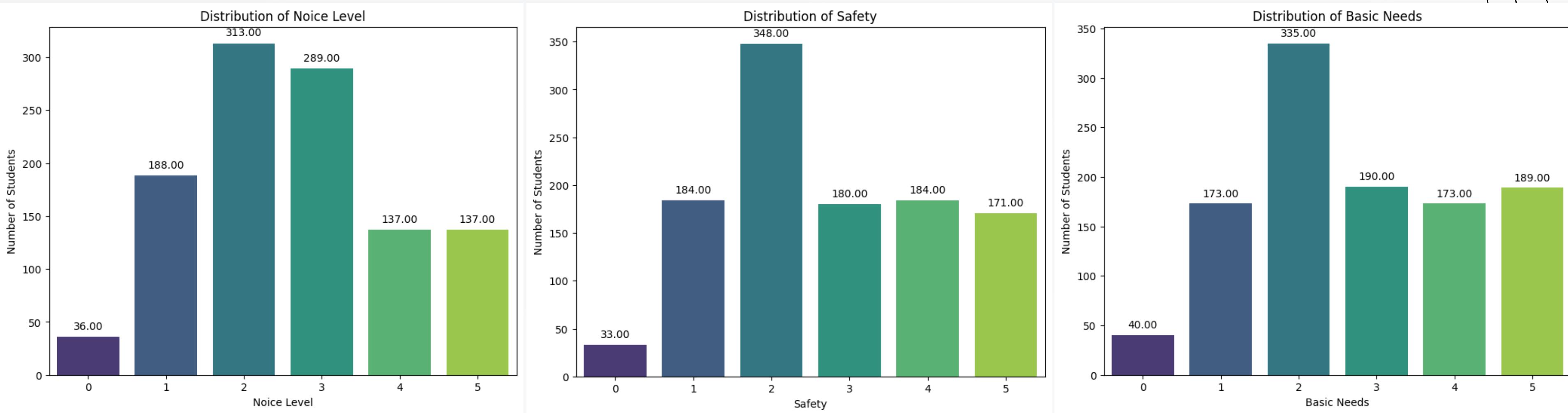
FEATURE DESCRIPTION

PHYSIOLOGICAL FACTORS



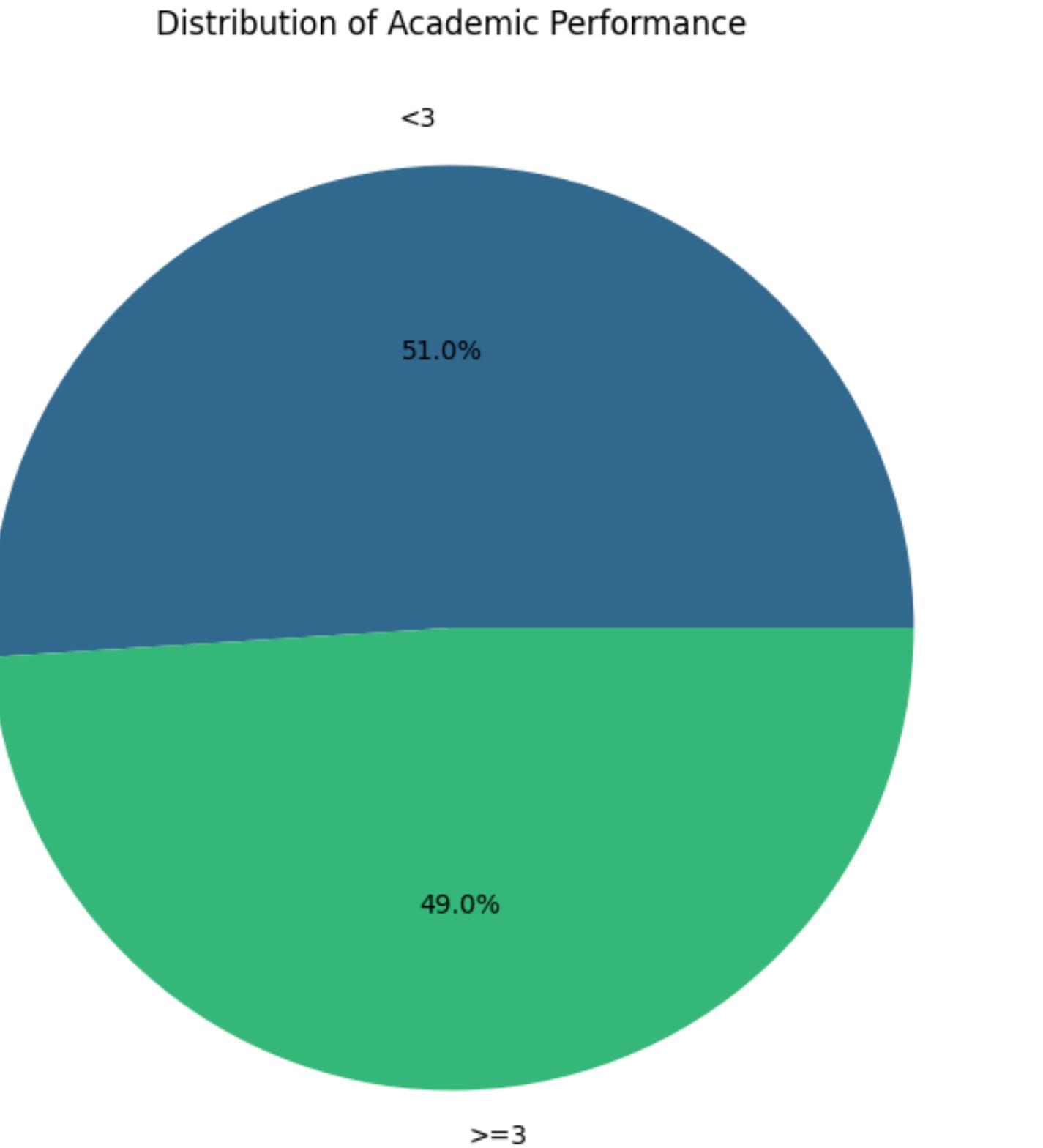
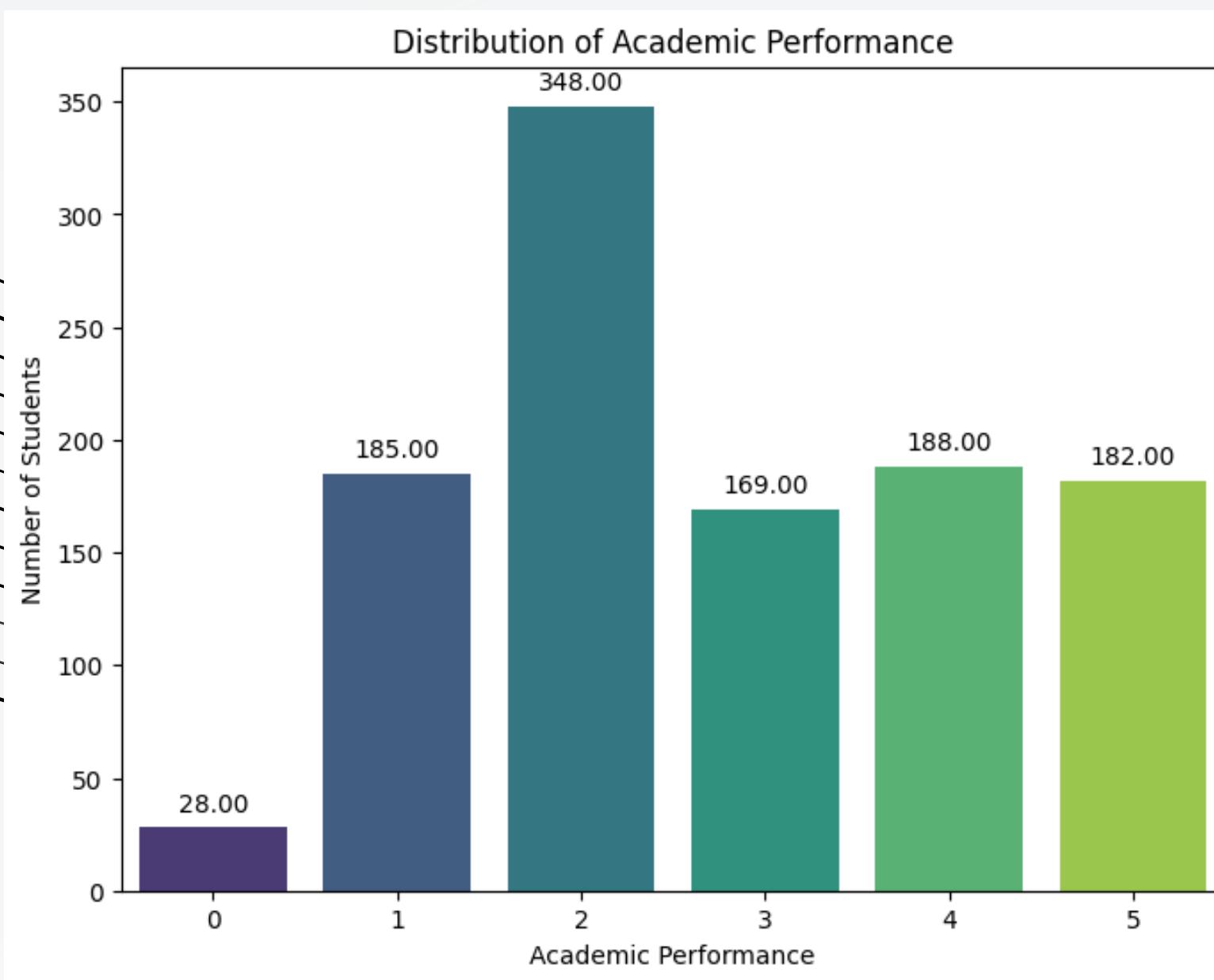
FEATURE DESCRIPTION

ENVIRONMENTAL FACTORS



FEATURE DESCRIPTION

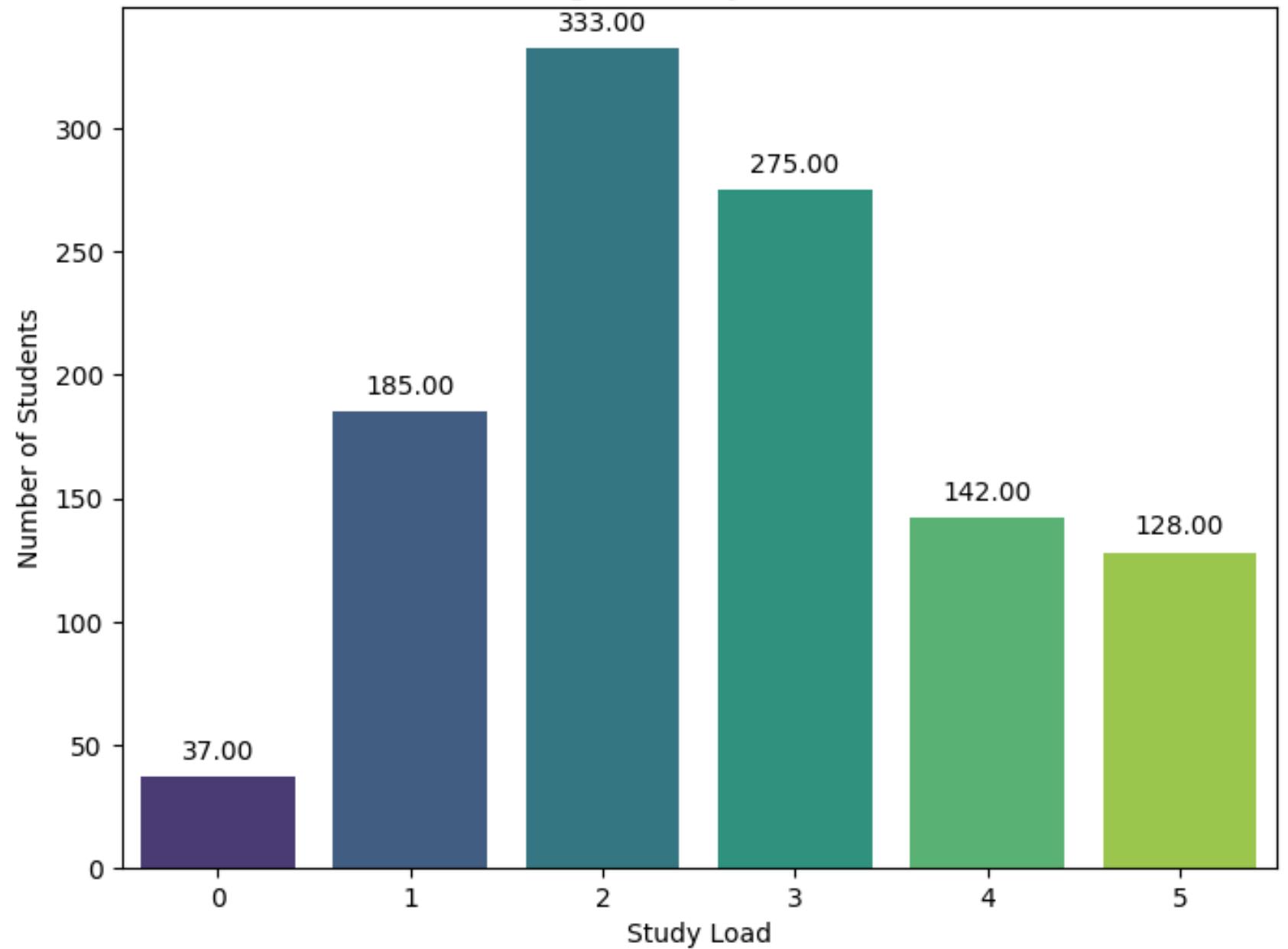
ACADEMIC FACTORS



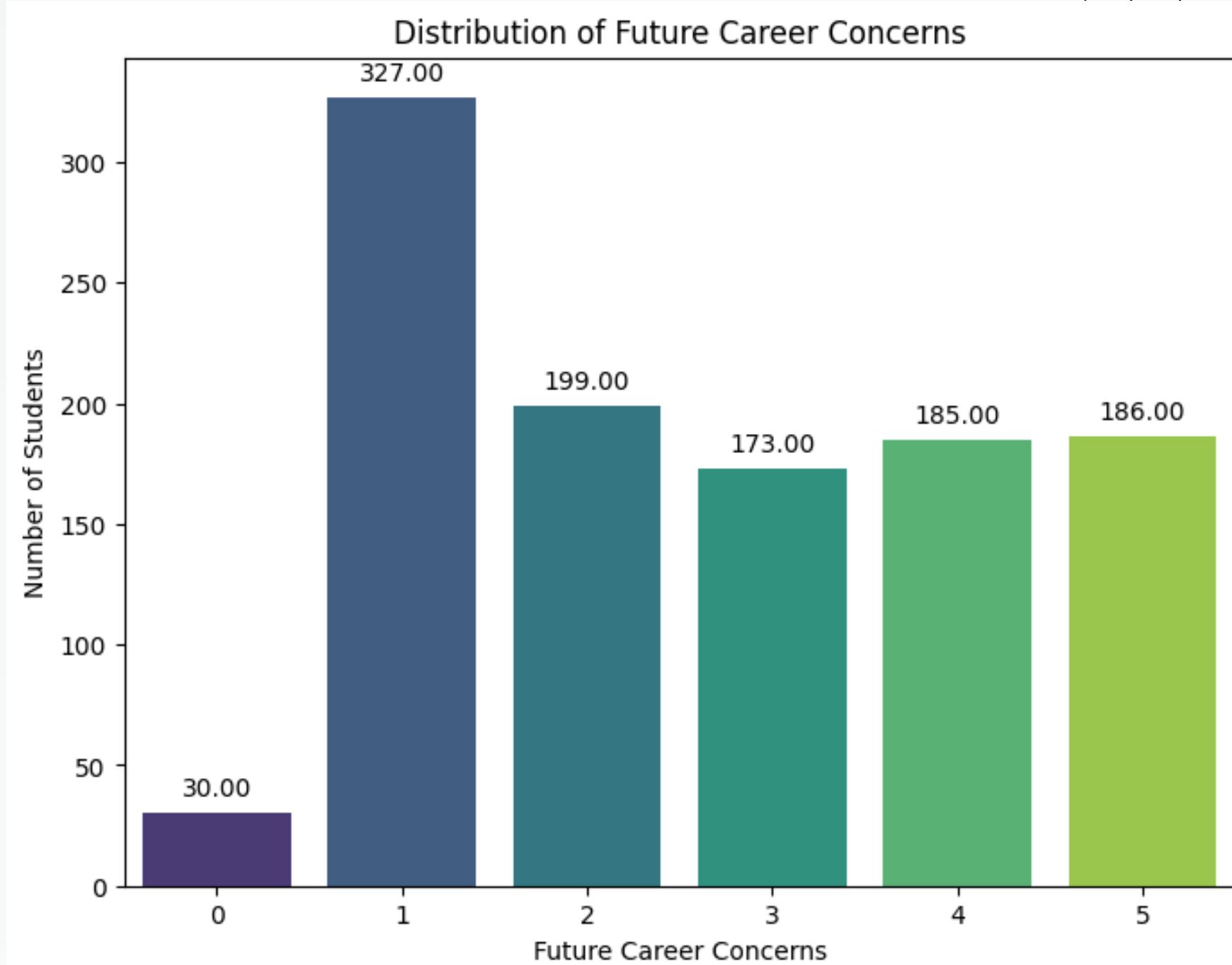
FEATURE DESCRIPTION

ACADEMIC FACTORS

Distribution of Study Load
Average of Study Load: 2.62



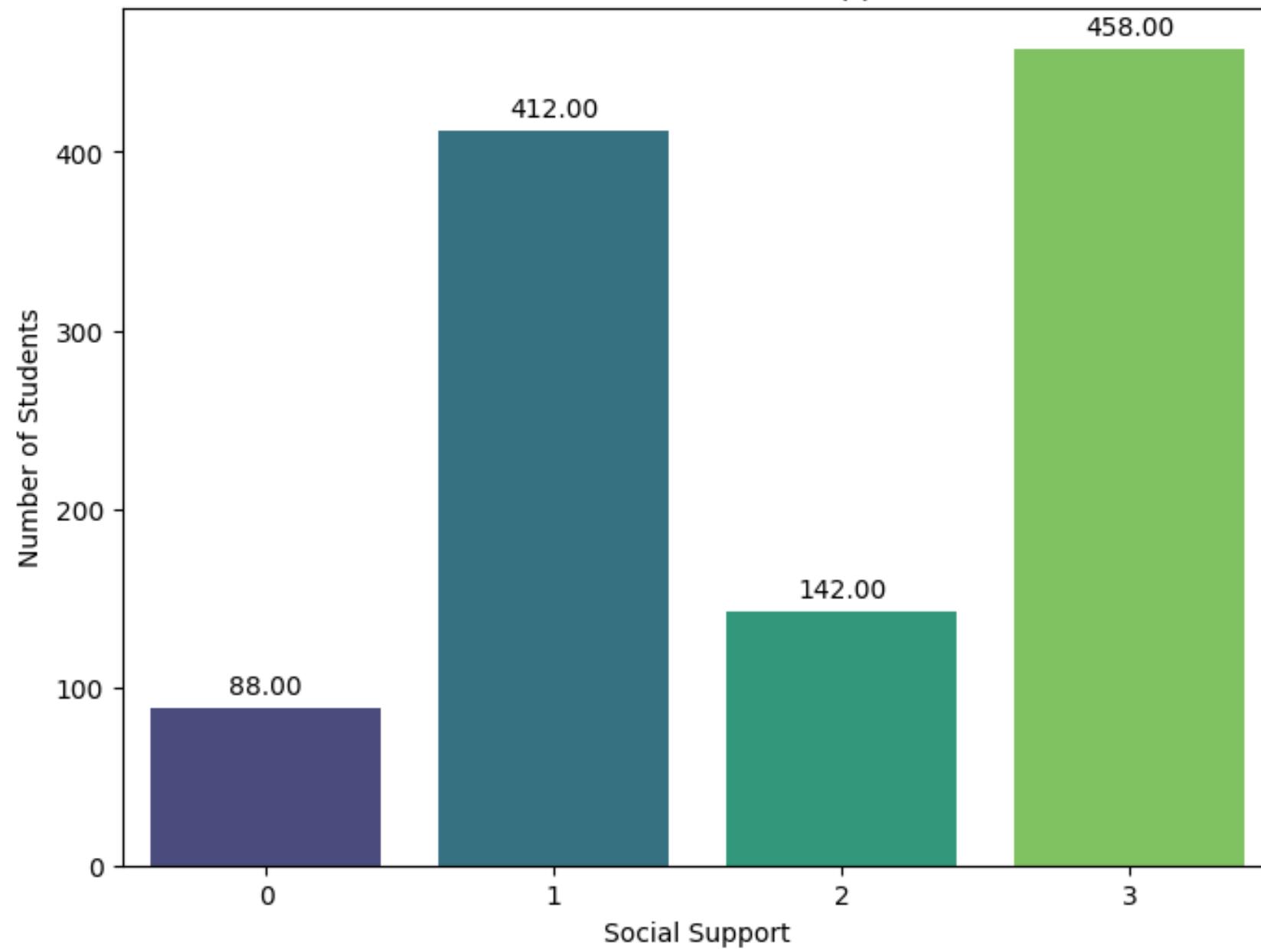
Distribution of Future Career Concerns



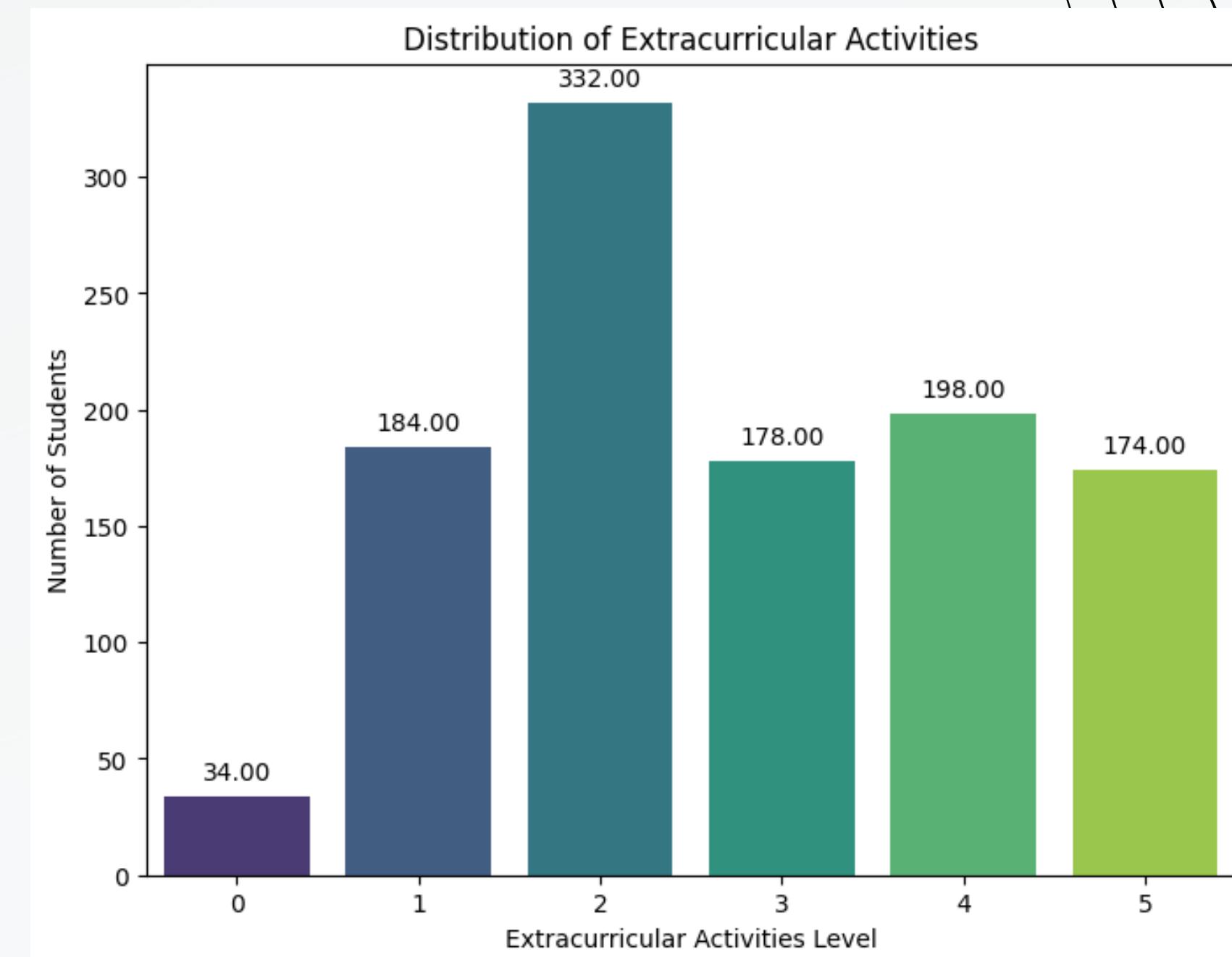
FEATURE DESCRIPTION

SOCIAL FACTORS

Distribution of Social Support

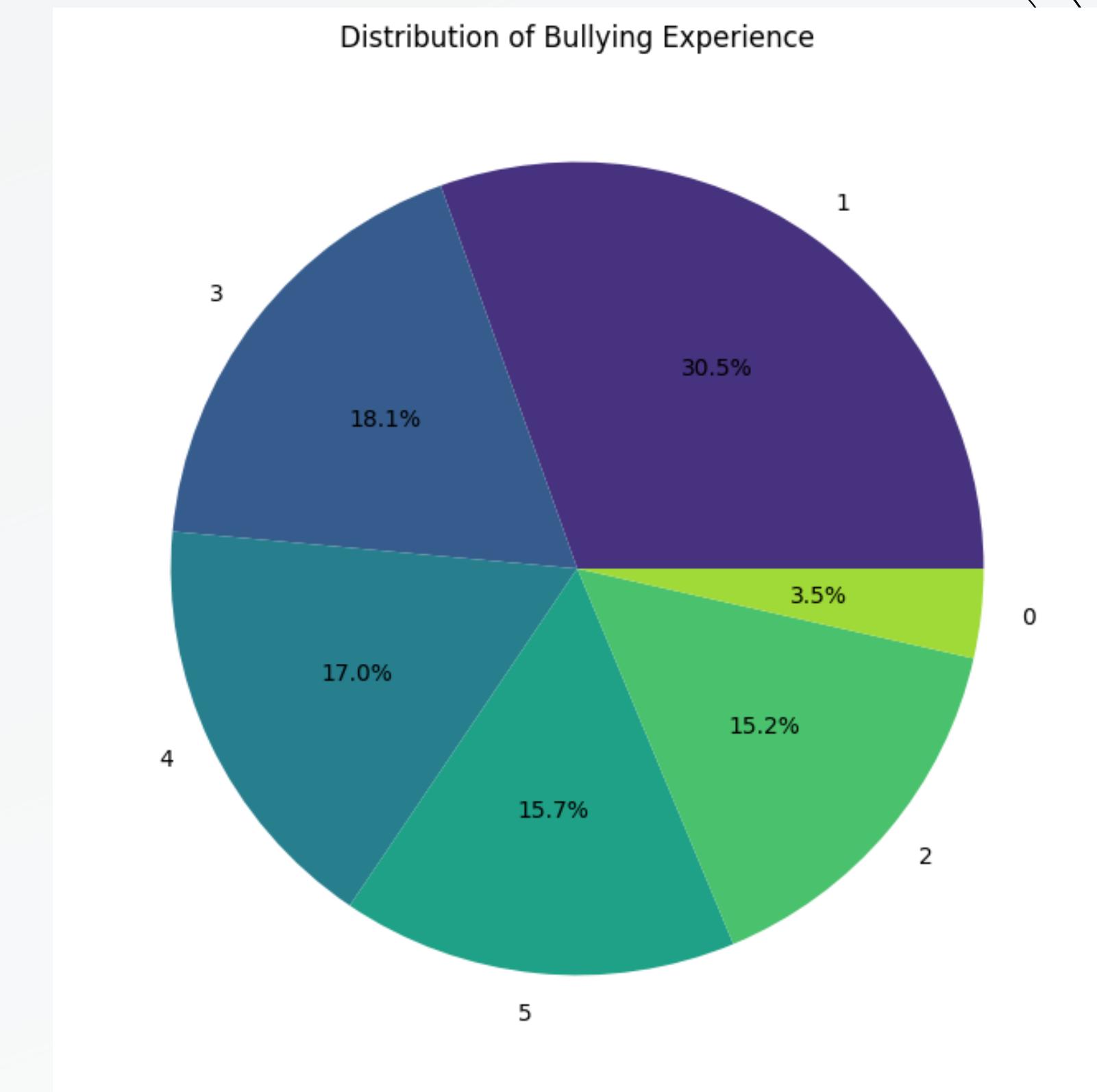
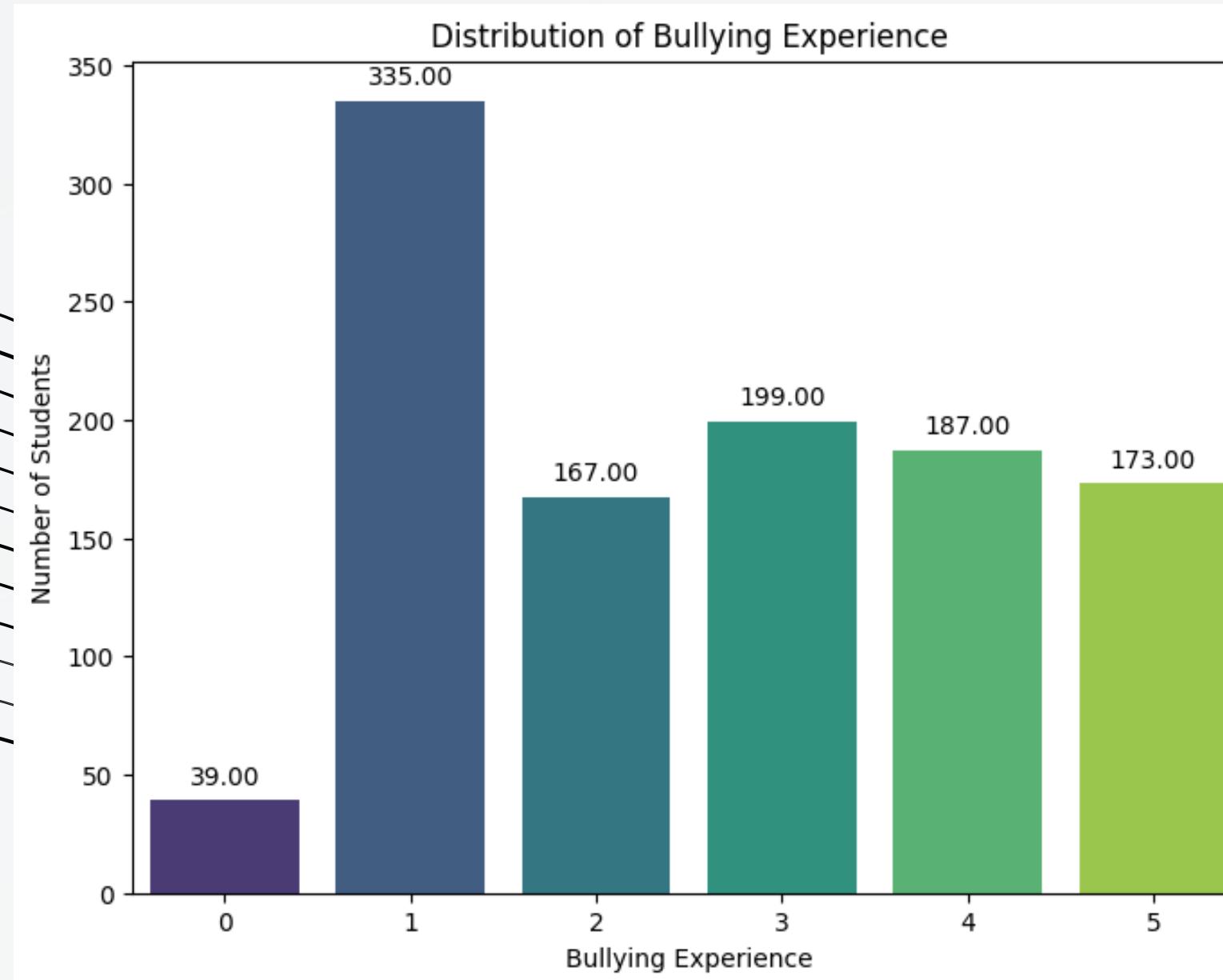


Distribution of Extracurricular Activities



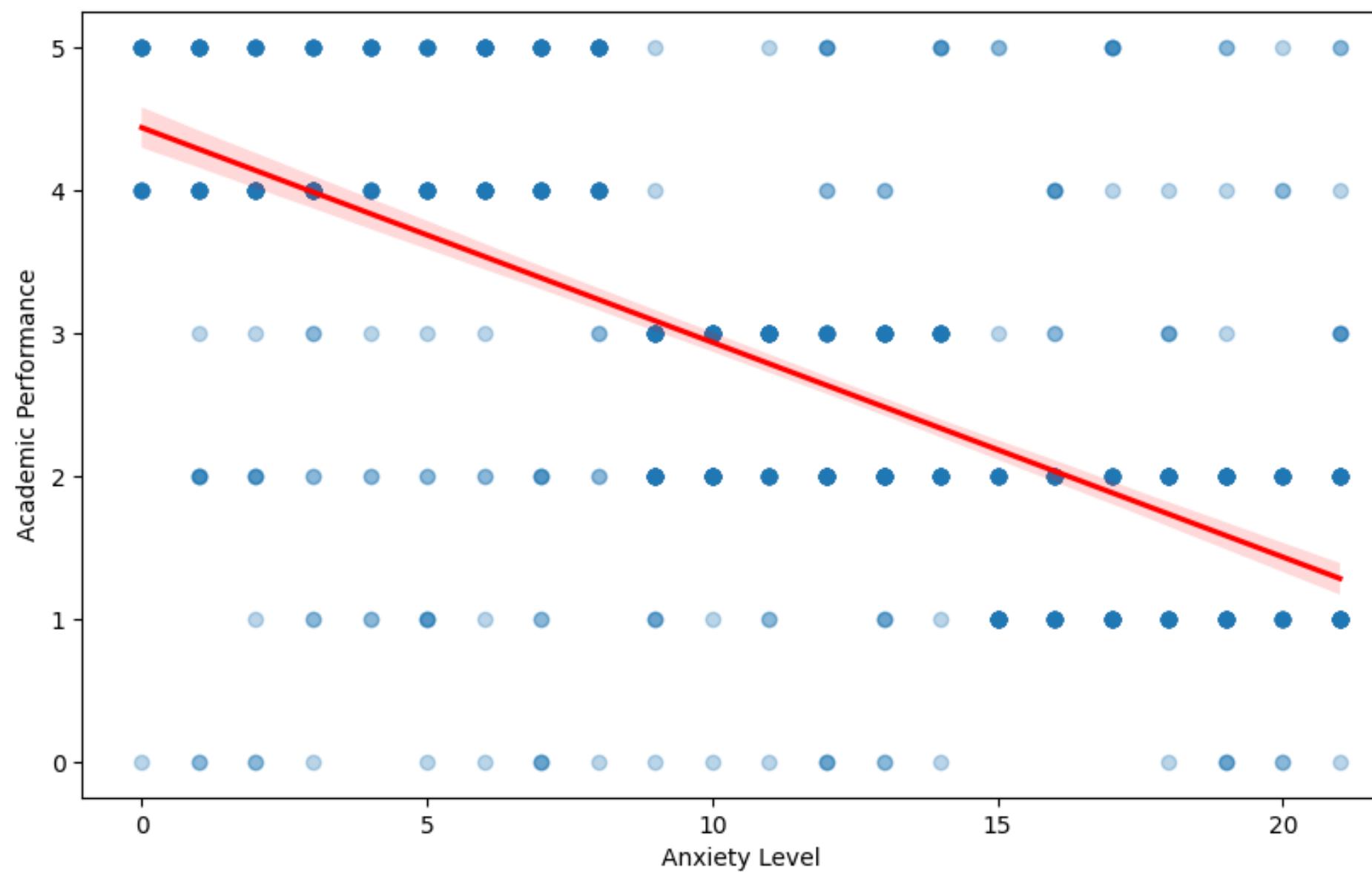
FEATURE DESCRIPTION

SOCIAL FACTORS

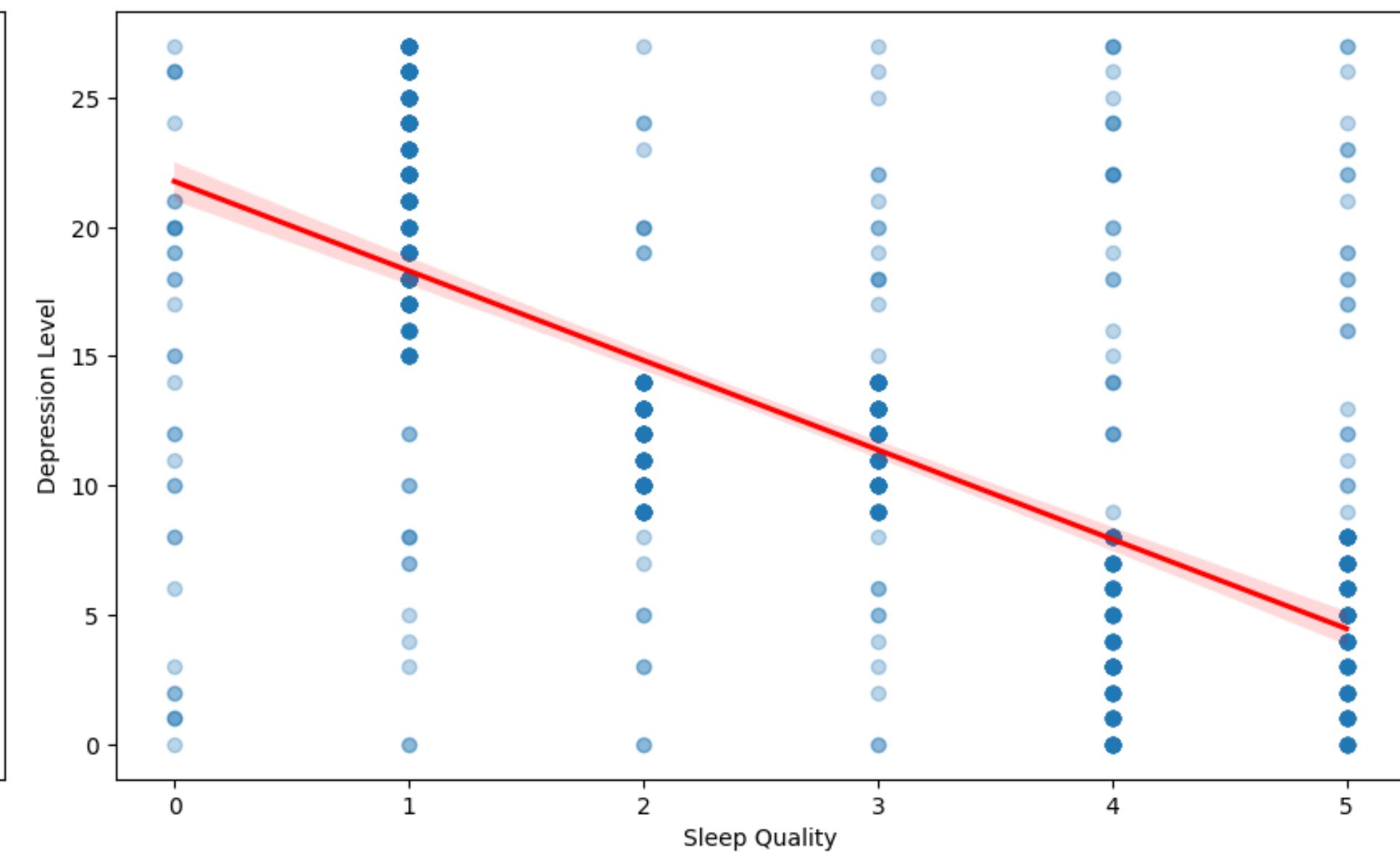


FEATURE DESCRIPTION COMPARATIVE ANALYSIS

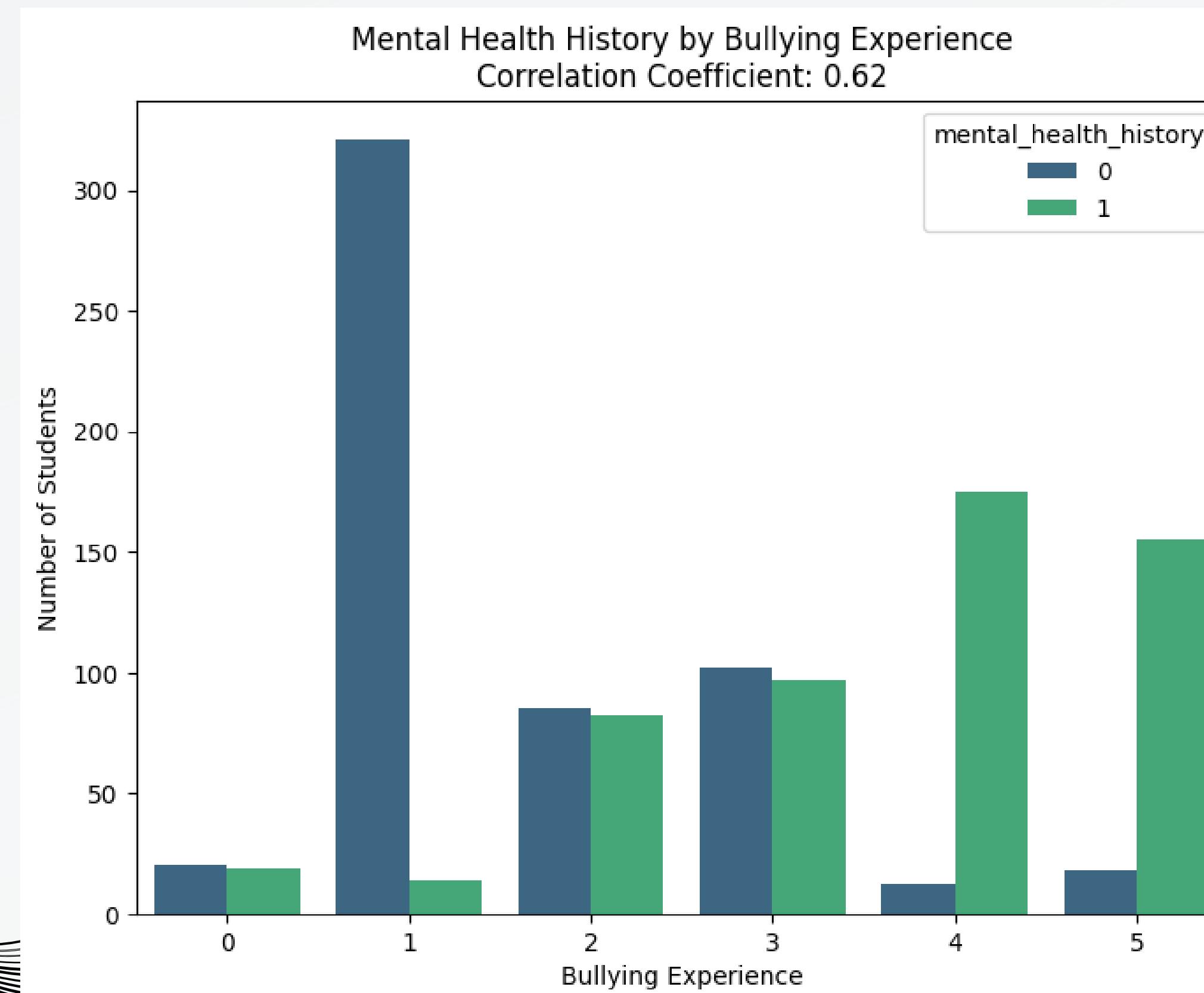
Correlation between Anxiety Level and Academic Performance
Correlation Coefficient: -0.65



Correlation between Sleep Quality and Depression
Correlation Coefficient: -0.69

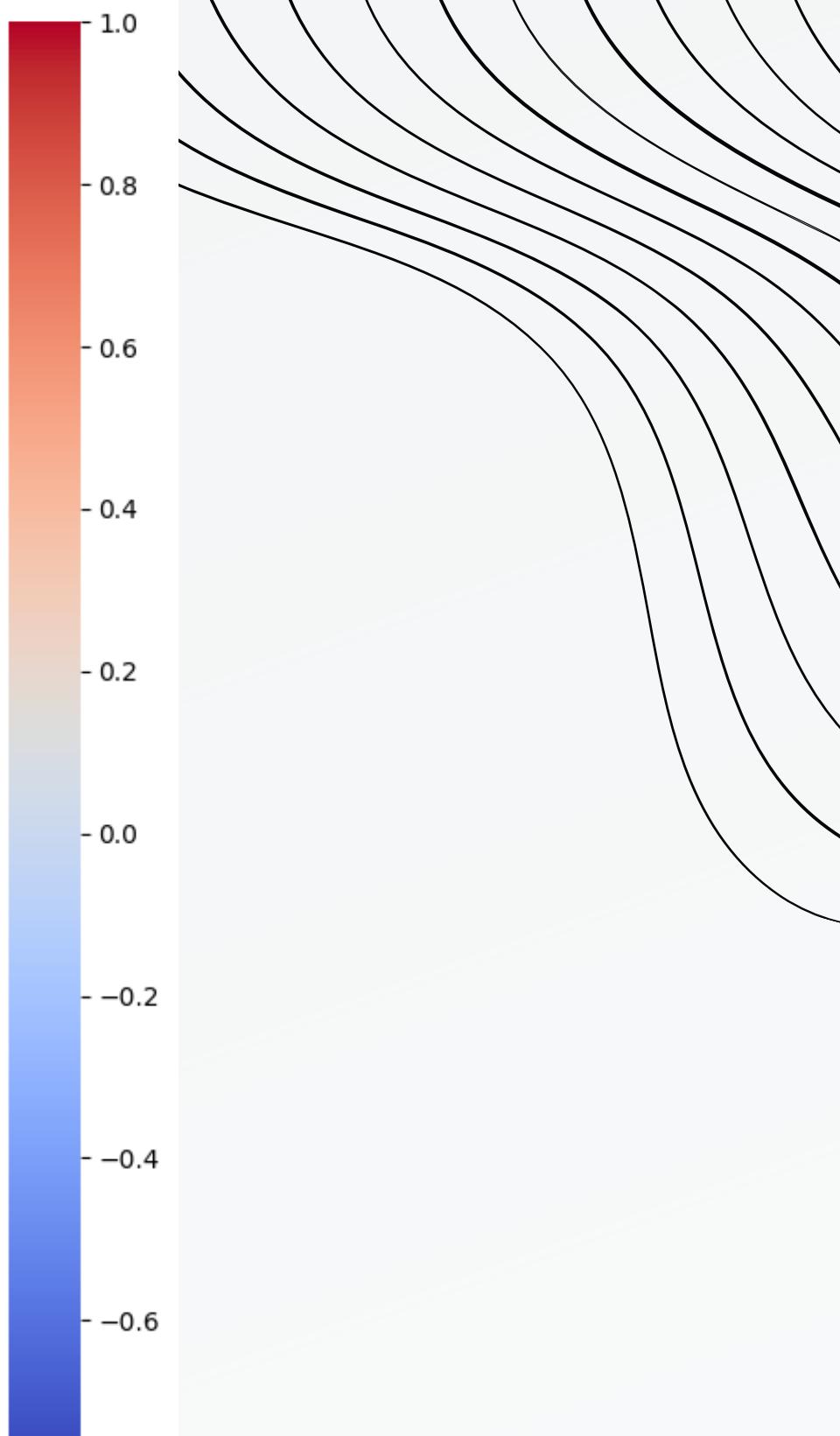
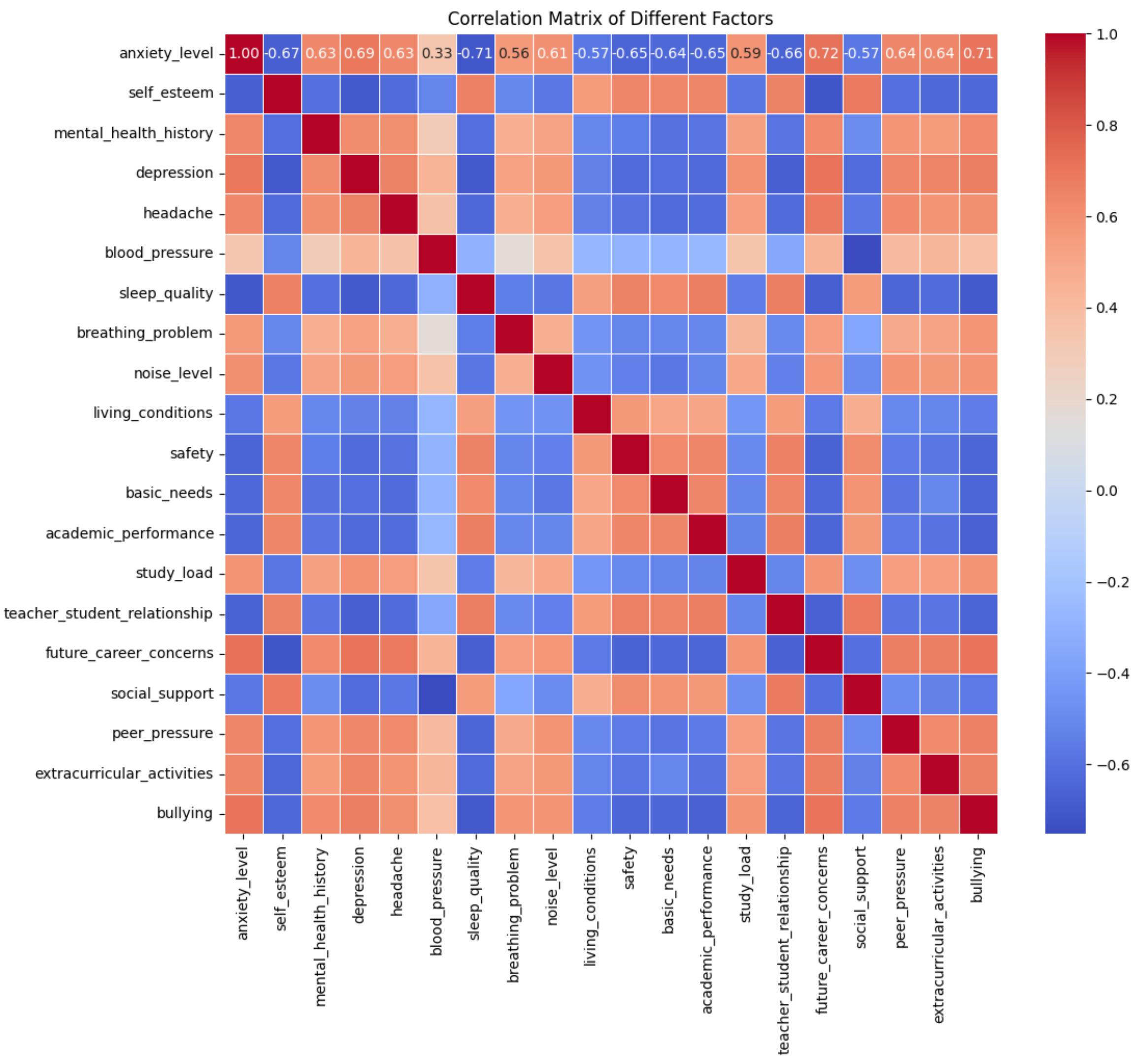
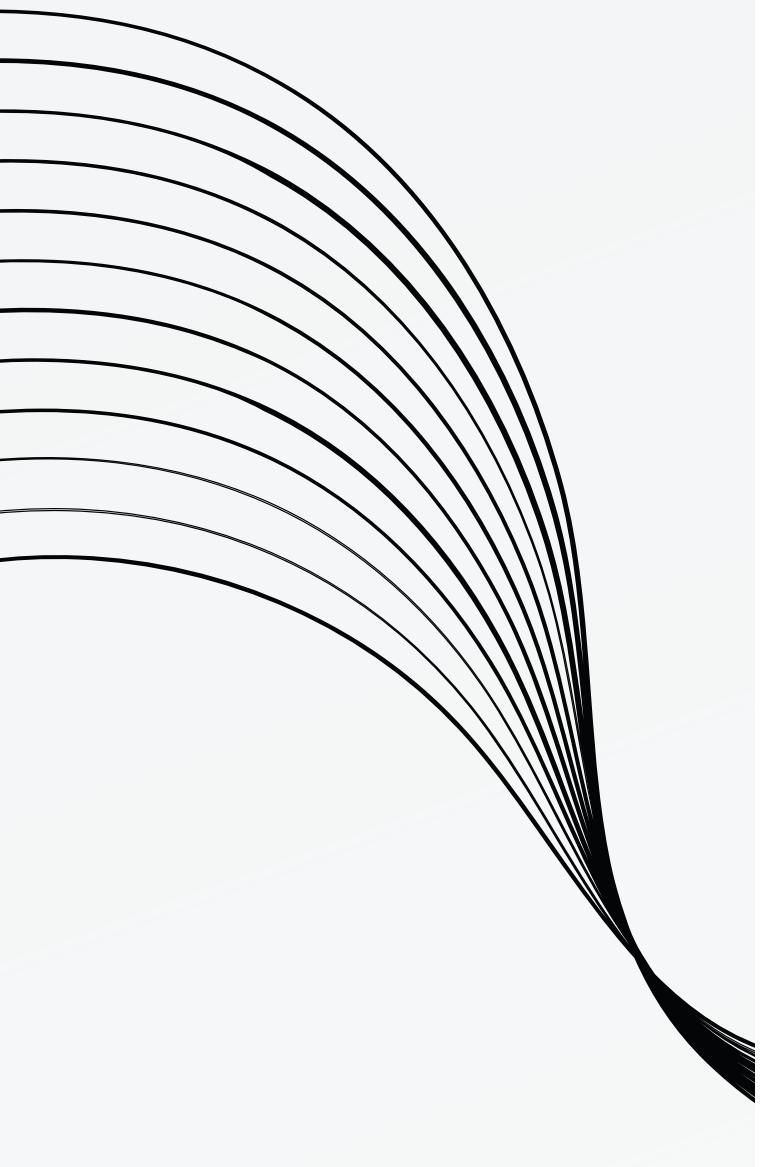


FEATURE DESCRIPTION COMPARATIVE ANALYSIS

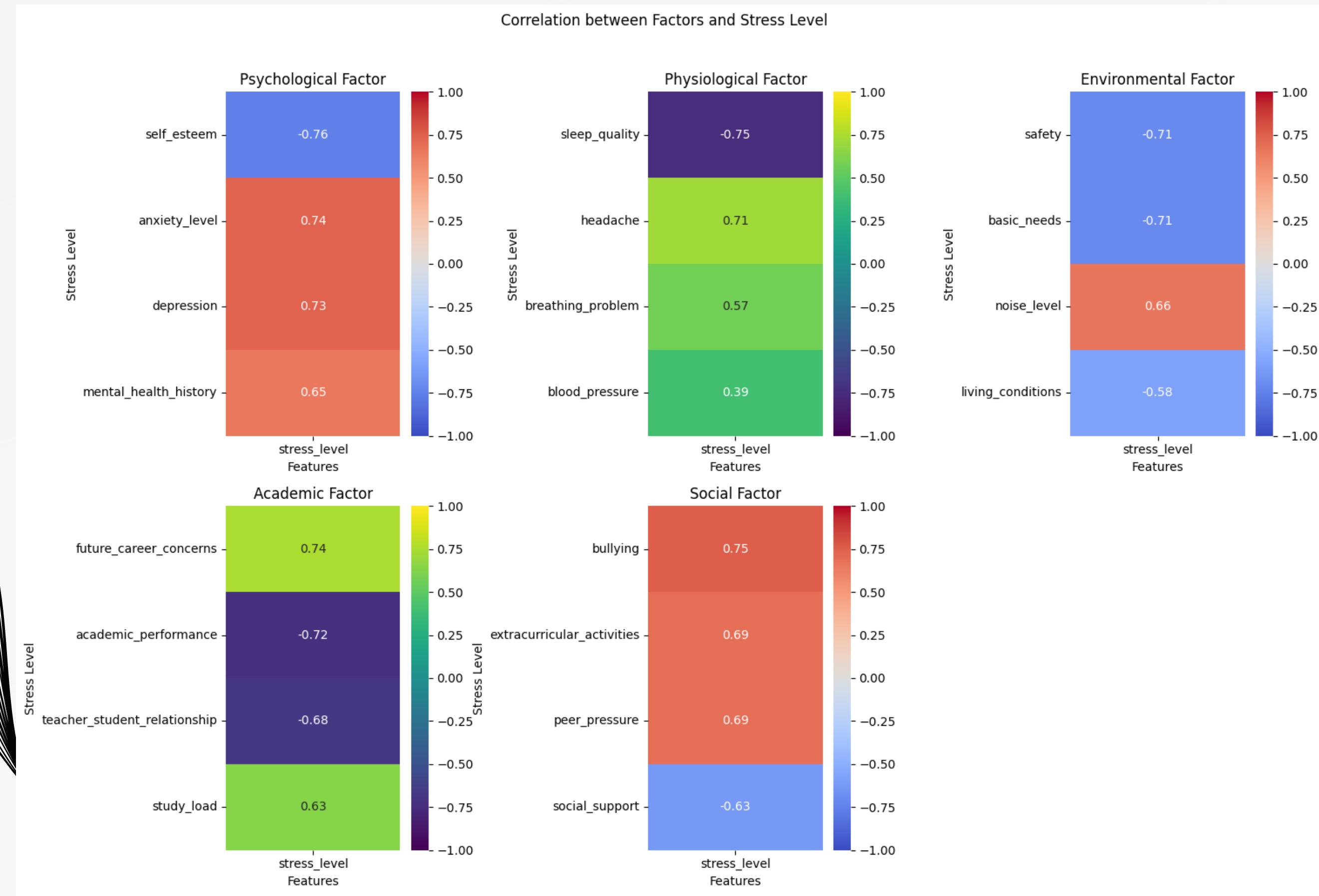


ANALYSIS



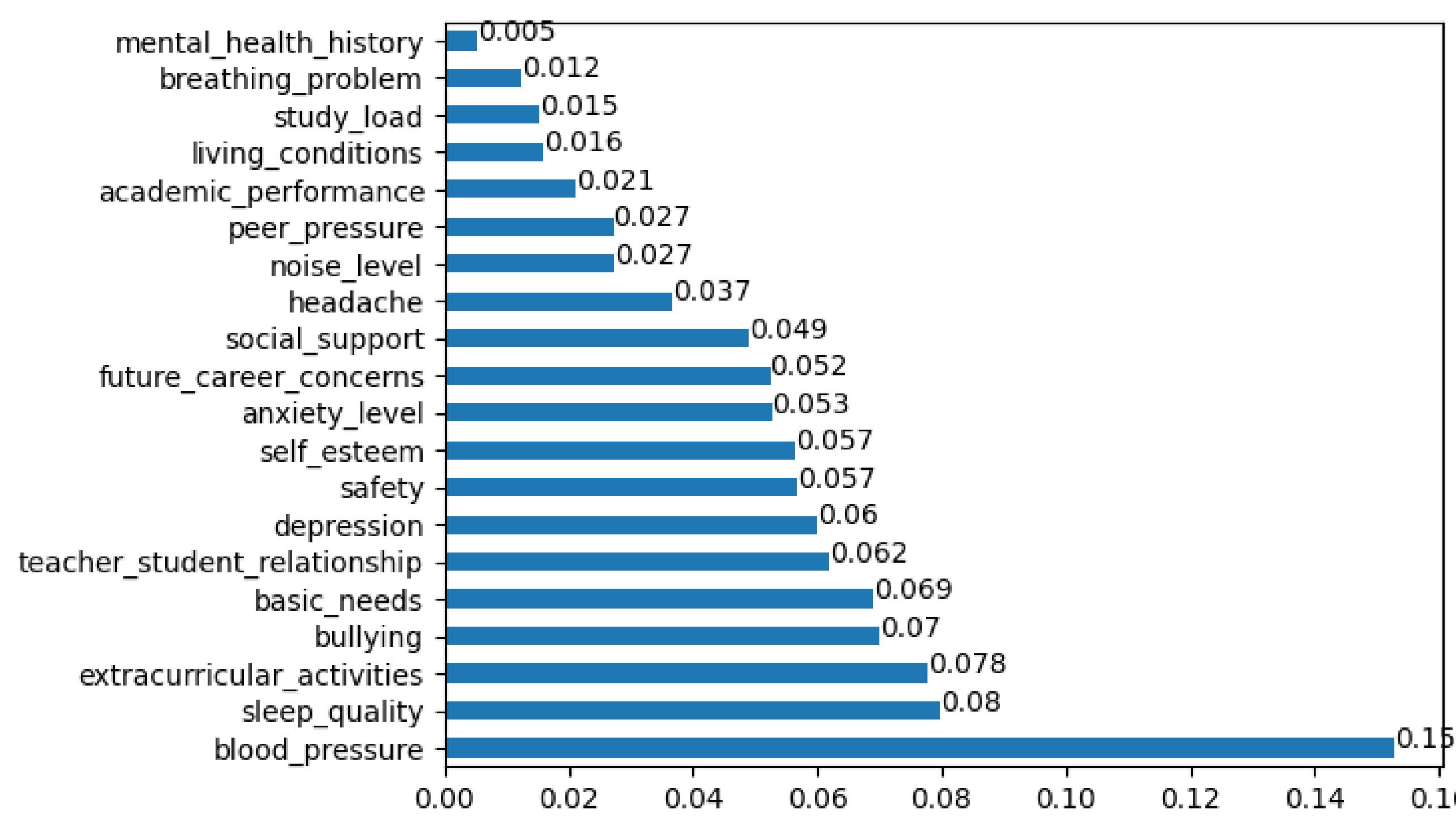


ANALYSIS



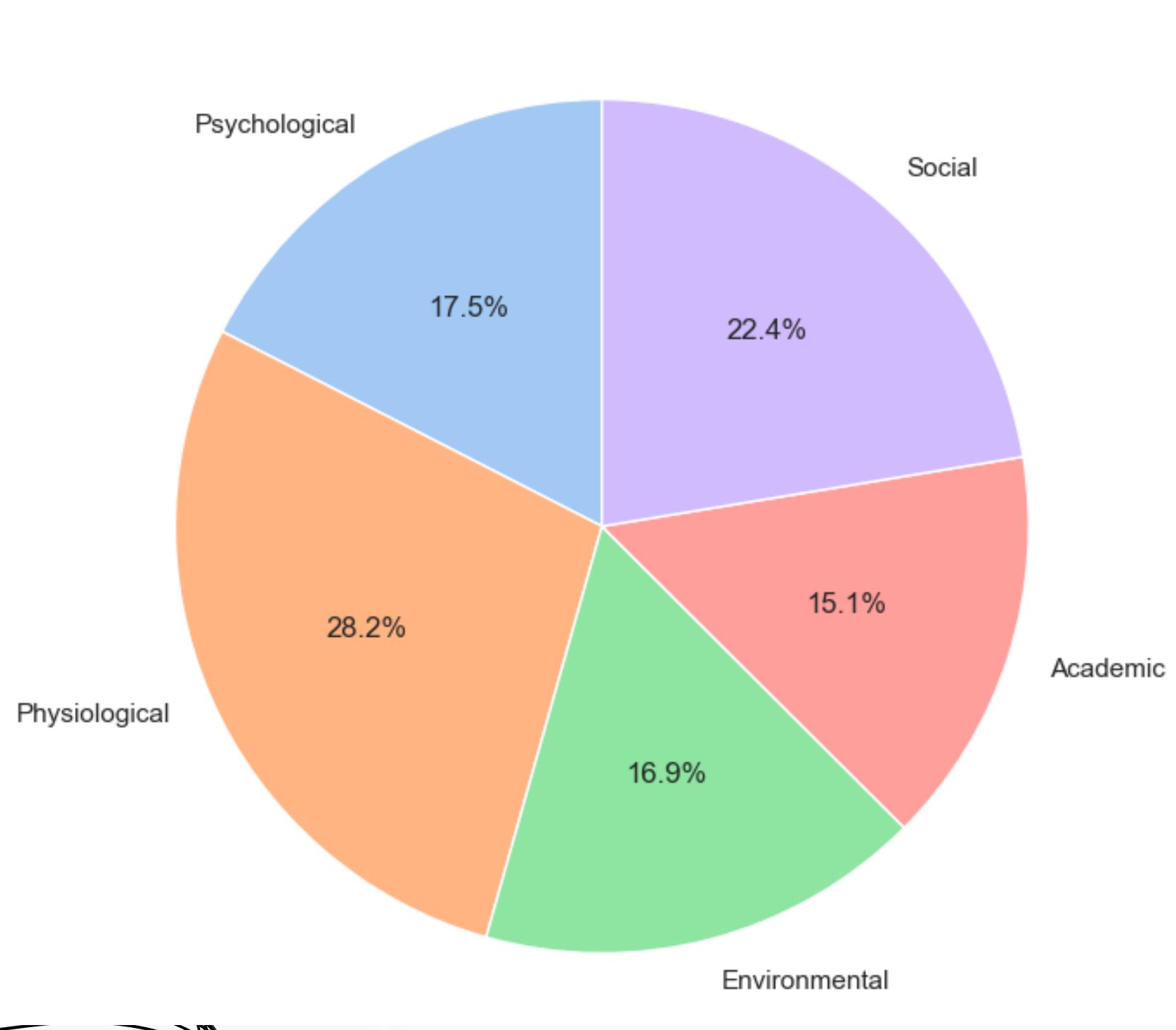
ANALYSIS

IMPORTANCE OF RANDOM FOREST



ANALYSIS

IMPORTANCE OF RANDOM FOREST



ANALYSIS

RANDOM FOREST RESULT



Model Accuracy: 0.8727272727272727

Confusion Matrix:

```
[[68  4  4]
 [ 6 63  4]
 [ 7  3 61]]
```

Classification Report:

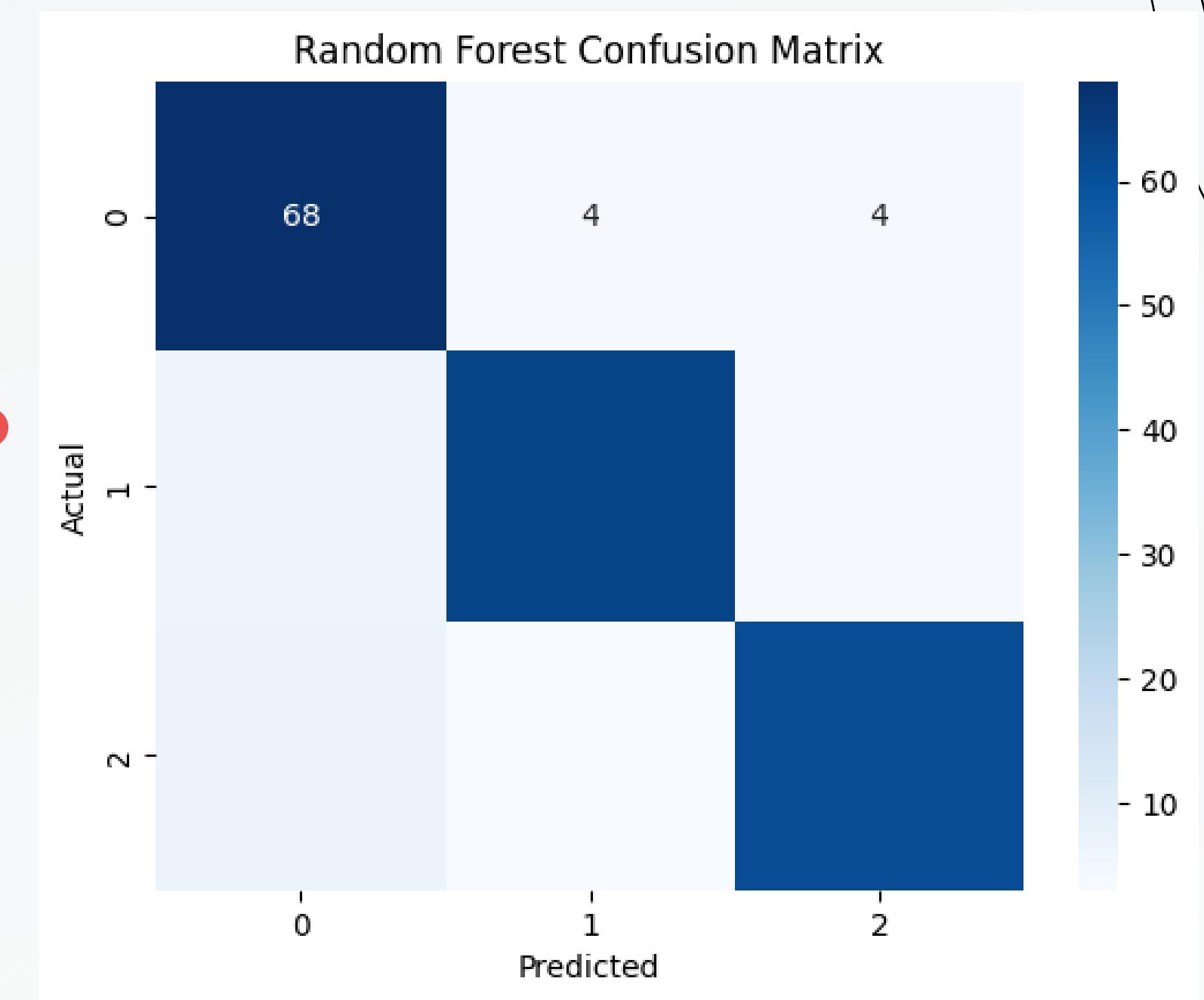
| | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0 | 0.84 | 0.89 | 0.87 | 76 |
| 1 | 0.90 | 0.86 | 0.88 | 73 |
| 2 | 0.88 | 0.86 | 0.87 | 71 |
| accuracy | | | 0.87 | 220 |
| macro avg | 0.87 | 0.87 | 0.87 | 220 |
| weighted avg | 0.87 | 0.87 | 0.87 | 220 |

ANALYSIS

RANDOM FOREST RESULT



0.87



ANALYSIS

REPO

DATASET

Dataset 1

<https://www.kaggle.com/datasets/pavansubhasht/ibm-hr-analytics-attrition-dataset>

Dataset 2

<https://www.kaggle.com/datasets/rxnach/student-stress-factors-a-comprehensive-analysis>

THANKS

