# Higher Education Grade Prediction

Gaurang Bista - u3252897

# Table of Content

# Introduction: Problem Statement

- The aim of this project is to provide students with a prediction as to whether or not they will pass their relative course studies based on various lifestyle attributes.

```
┌──────────────────┐      ┌──────────────────┐      ┌──────────────────┐
│ Exploratory Data │ ──►  │ Predictive Data  │ ──►  │ Implementation   │
│ Analysis (EDA)   │      │ Analysis (PDA)   │      │                  │
└──────────────────┘      └──────────────────┘      └──────────────────┘
```

# Dataset details

- Dataset name: *Higher Education Students Performance Evaluation* [3]
- Provided by **Nevriye Yilmaz** and **Boran Sekeroglu** in 2019.
- Data was collected from the Faculty of Engineering and Faculty of Educational Sciences students [4]

| Data Set Characteristics: | Multivariate | Number of Instances: | 145 | Area: | Social |
|---|---|---|---|---|---|
| Attribute Characteristics: | Integer | Number of Attributes: | 33 | Date Donated | 2021-01-30 |
| Associated Tasks: | Classification | Missing Values? | N/A | Number of Web Hits: | 48974 |

Source: UCI [4]

32- grade - OUTPUT Grade (0: Fail, 1: DD, 2: DC, 3: CC, 4: CB, 5: BB, 6: BA, 7: AA)

Source: Kaggle [3]

# EDA Overview

- The Exploratory Data Analysis (EDA and visualisation for this project was divided into three parts
  - Understanding the data
  - Cleaning the data
  - Creating visualisations which can be used to answer 5 questions about the dataset

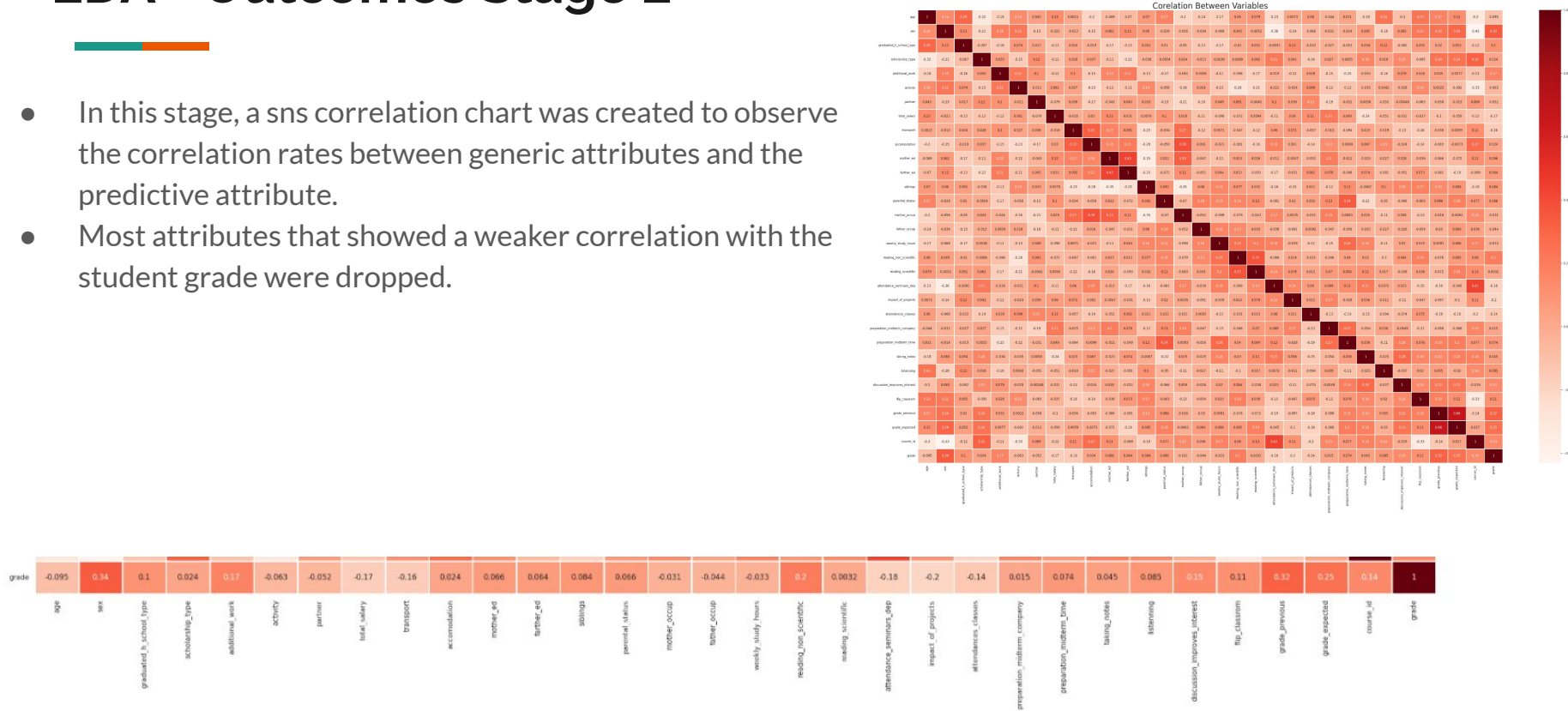| Understanding the data | → | Cleaning the data | → | Visualisations |
| --- | --- | --- | --- | --- |

# EDA - Outcomes Stage 1

- In this stage, the dataset was mounted to google drive and the operations below were conducted to gain a better understanding of the data.
- The dataset was checked for any null-values and the data types were also examined.

[ ] df.shape

    (145, 33)

[ ] df.columns

    Index(['student_id', 'age', 'sex', 'graduated_h_school_type',
           'scholarship_type', 'additional_work', 'activity', 'partner',
           'total_salary', 'transport', 'accomodation', 'mother_ed', 'farther_ed',
           'siblings', 'parental_status', 'mother_occup', 'father_occup',
           'weekly_study_hours', 'reading_non_scientific', 'reading_scientific',
           'attendance_seminars_dep', 'impact_of_projects', 'attendances_classes',
           'preparation_midterm_company', 'preparation_midterm_time',
           'taking_notes', 'listening', 'discussion_improves_interest',
           'flip_classrom', 'grade_previous', 'grade_expected', 'course_id',
           'grade'],
          dtype='object')

df.nunique()

    student_id                       145
    age                                3
    sex                                2
    graduated_h_school_type            3
    scholarship_type                   5
    additional_work                    2
    activity                           2
    partner                            2
    total_salary                       5
    transport                          4
    accomodation                       4
    mother_ed                          6
    farther_ed                         6
    siblings                           5
    parental_status                    3
    mother_occup                       5
    father_occup                       5
    weekly_study_hours                 5
    reading_non_scientific             3
    reading_scientific                 3
    attendance_seminars_dep            2
    impact_of_projects                 3
    attendances_classes                3
    preparation_midterm_company        3
    preparation_midterm_time           3
    taking_notes                       3
    listening                          3
    discussion_improves_interest       3
    flip_classrom                      3
    grade_previous                     5
    grade_expected                     4
    course_id                          9
    grade                              8
    dtype: int64

df.info()

    <class 'pandas.core.frame.DataFrame'>
    RangeIndex: 145 entries, 0 to 144
    Data columns (total 33 columns):
     #   Column                         Non-Null Count   Dtype
    ---  ------                         --------------   -----
     0   student_id                     145 non-null     object
     1   age                            145 non-null     int64
     2   sex                            145 non-null     int64
     3   graduated_h_school_type        145 non-null     int64
     4   scholarship_type               145 non-null     int64
     5   additional_work                145 non-null     int64
     6   activity                       145 non-null     int64
     7   partner                        145 non-null     int64
     8   total_salary                   145 non-null     int64
     9   transport                      145 non-null     int64
     10  accomodation                   145 non-null     int64
     11  mother_ed                      145 non-null     int64
     12  farther_ed                     145 non-null     int64
     13  siblings                       145 non-null     int64
     14  parental_status                145 non-null     int64
     15  mother_occup                   145 non-null     int64
     16  father_occup                   145 non-null     int64
     17  weekly_study_hours             145 non-null     int64
     18  reading_non_scientific         145 non-null     int64
     19  reading_scientific             145 non-null     int64
     20  attendance_seminars_dep        145 non-null     int64
     21  impact_of_projects             145 non-null     int64
     22  attendances_classes            145 non-null     int64
     23  preparation_midterm_company    145 non-null     int64
     24  preparation_midterm_time       145 non-null     int64
     25  taking_notes                   145 non-null     int64
     26  listening                      145 non-null     int64
     27  discussion_improves_interest   145 non-null     int64
     28  flip_classrom                  145 non-null     int64
     29  grade_previous                 145 non-null     int64
     30  grade_expected                 145 non-null     int64
     31  course_id                      145 non-null     int64
     32  grade                          145 non-null     int64
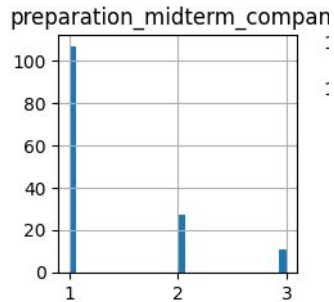
# EDA - Outcomes Stage 2

- In this stage, a sns correlation chart was created to observe the correlation rates between generic attributes and the predictive attribute.
- Most attributes that showed a weaker correlation with the student grade were dropped.
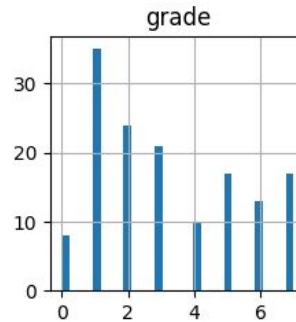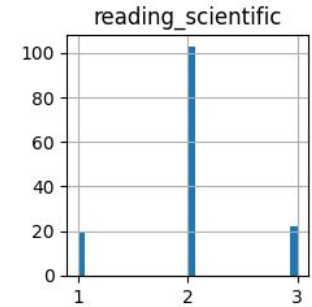


Correlation Between Variables

# EDA - Outcome Stage 3

○ Q1. Which grade is most commonly achieved by students in the dataset?
○ Q2. What are the most common weekly study hours?
○ Q3. What is the most common reading frequency for scientific text?
○ Q4. How do most students prepare for midterm? Alone? With friends?
○ Q5. What is the correlation coefficient between previous grade and current grade?



reading_scientific

**1**: None, **2**: Sometimes, **3**: Often



preparation_midterm_compan...

**1**: alone, **2**: with friends, **3**: n/a



grade

**0**: Fail, **1**: DD, **2**: DC, **3**: CC, **4**: CB, **5**: BB, **6**: BA, **7**: AA



weekly_study_hours

**1**: None, **2**: <5 hours, **3**: 6-10 hours, **4**: 11-20 hours, **5**: more than 20 hours

# EDA - Q5


Corelation Between Variables

# PDA Outcomes

- The PDA was divided into three stages:
  - Initial dataset prediction
  - Pre -processed dataset prediction
  - Further preprocessed dataset prediction

| Prediction of student grade using initial dataset | → | Prediction of student grade using preprocessed dataset | → | Prediction of student grade using further preprocessed dataset |

# PDA Outcomes - Stage 1

- This stage involved testing four predictive models Naive Bayes, Support Vector Machines, Gradient Boosting, and Random Forest on the initial dataset without preprocessing.
- Due to the dataset being small and multiple grade types, the accuracy rates of each model were low.

```
Performance on Training set
NB: 0.156818 (0.085713)

SVM: 0.217424 (0.108365)

GBM: 0.252273 (0.142257)

RF: 0.337879 (0.136582)
```

| Fail | DD | DC | CC | CB | BB | BA | AA |
|------|----|----|----|----|----|----|----|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

# PDA Outcomes - Stage 2

Performance on Training set
NB: 0.493182 (0.115661)

SVM: 0.779545 (0.121753)

GBM: 0.656061 (0.088035)

RF: 0.668182 (0.144250)

- To improve the model prediction accuracy, the following categorisation process was conducted.

Fail

| 0 |

Satisfactory

| 1 | 2 | 3 |

Above Average

| 4 | 5 | 6 | 7 |

# PDA Outcomes - Stage 3

Performance on Training set
NB: 0.747727 (0.112980)

SVM: 0.949242 (0.055778)

- To further improve the model prediction accuracy, the following categorisation process was conducted.

GBM: 0.915152 (0.074597)

- The predictive outcome grade was made to ba binary

RF: 0.940909 (0.053761)

Fail

Pass

| 0 | | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

# Implementation/Deployment Plan and Status update

- I plan to use Tkinter for the implementation/deployment phase of the ST1 Capstone project
- The window will ask the user for their inputs for each attribute and provide a prediction as to whether or not they will pass.

# Bibliography

- **[1]** J. Brownlee, "How to Choose a Feature Selection Method For Machine Learning," Machine Learning Mastery, Nov. 26, 2019. https://machinelearningmastery.com/feature-selection-with-real-and-categorical-data/
- **[2]** T. Bush, "Predictive Analysis: Definition, Tools, and Examples," pestleanalysis.com, Jun. 01, 2020. https://pestleanalysis.com/predictive-analysis/
- **[3]** "Higher Education Students Performance Evaluation," www.kaggle.com. https://www.kaggle.com/datasets/mariazhokhova/higher-education-students-performance-evaluation (accessed May 02, 2023).
- **[4]** "UCI Machine Learning Repository: Higher Education Students Performance Evaluation Dataset Data Set," archive.ics.uci.edu. https://archive.ics.uci.edu/ml/datasets/Higher+Education+Students+Performance+Evaluation+Dataset
- **[5]** C. Cote, "What Is Predictive Analytics? 5 Examples | HBS Online," Business Insights - Blog, Oct. 26, 2021. https://online.hbs.edu/blog/post/predictive-analytics