# Project 4: Analysis on Students' Performance

Team 9 – Rachel Schoen, Andrea Wu, Mahind Rao, Mahalel Peter

# Agenda

**01** → **Problem Statement**

**02** → **Data Description & Preparation**

**03** → **Exploratory Analysis**

**04** → **Database Design**

**05** → **Model Designs & Evaluations**
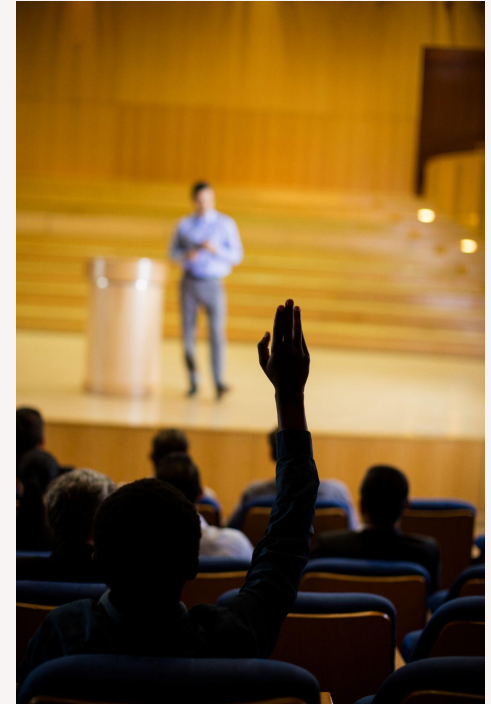
**06** → **Limitations & Future Considerations**

# Problem Statement

**What are the key factors that influence the success of a student?**

- How do students with additional work responsibilities fare academically compared to those without?
- How does the type of accommodation affect students' study habits and performance?
- How does discussion and group work influence students' interest and success in their courses?
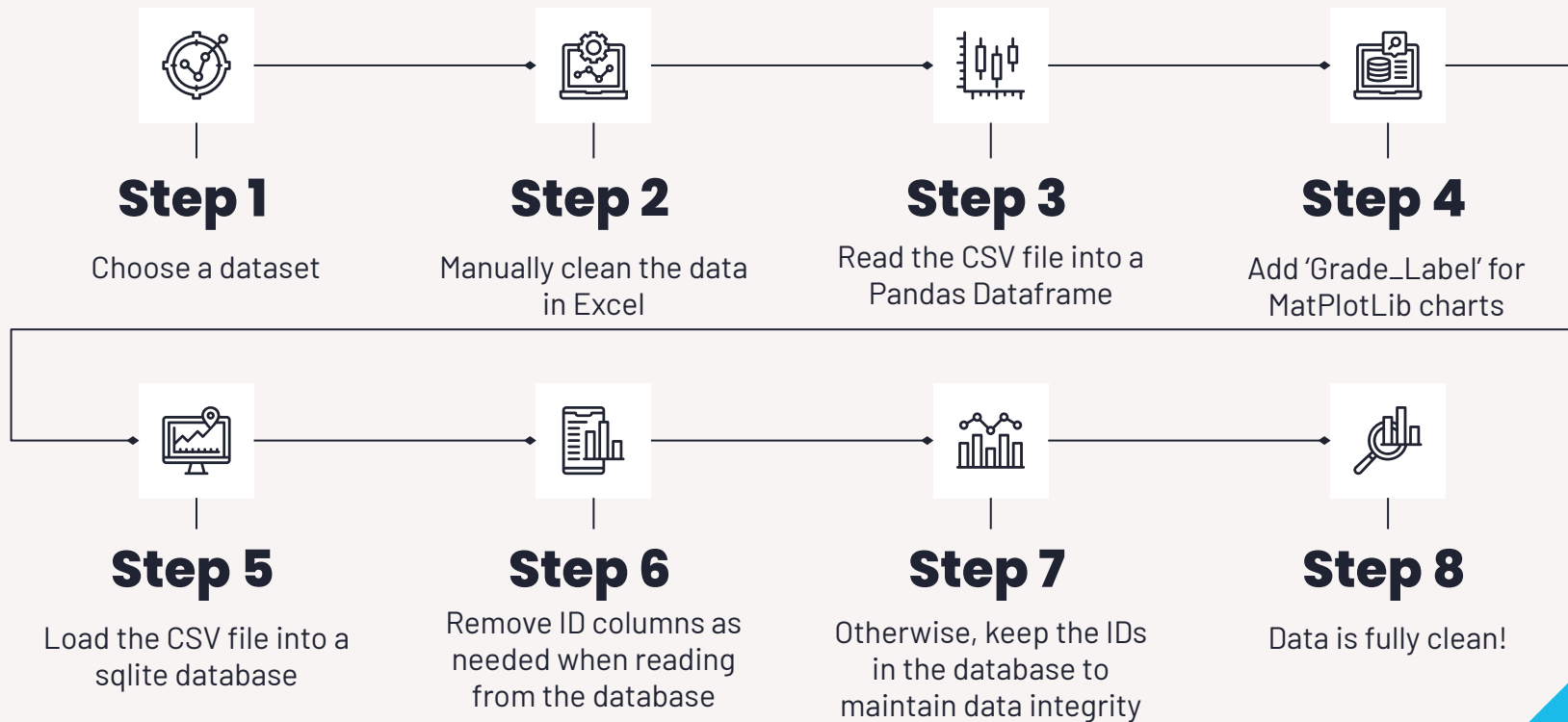
**And most importantly:**
- **Can we predict a student's academic performance (e.g., grade point average or success in courses) based on their personal background, study habits, and extracurricular activities?**

# Data Description

- **Dataset: Students performance**
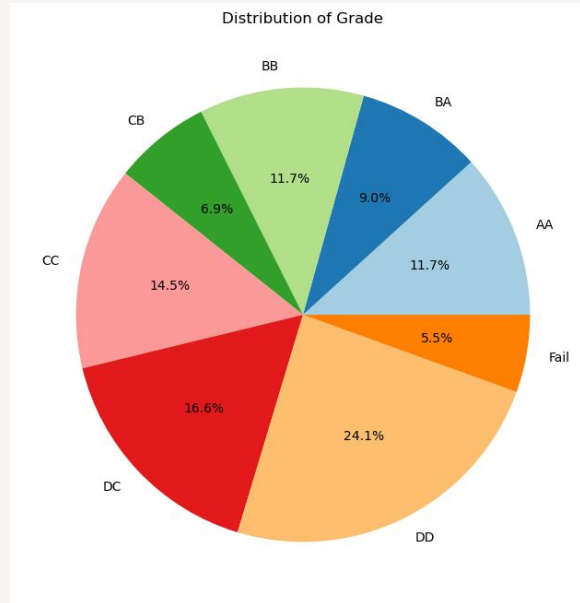
- **Format: CSV file**

- **Data Source: Dataset from Kaggle**
  https://www.kaggle.com/datasets/joebeachcapital/students-performance/data

- **Data Structure**
    - Data of 145 students
    - 33 Columns
    - Majority is categorical data

# Data Preparation

## Step 1
Choose a dataset

## Step 2
Manually clean the data in Excel

## Step 3
Read the CSV file into a Pandas Dataframe

## Step 4
Add 'Grade_Label' for MatPlotLib charts

## Step 5
Load the CSV file into a sqlite database

## Step 6
Remove ID columns as needed when reading from the database

## Step 7
Otherwise, keep the IDs in the database to maintain data integrity
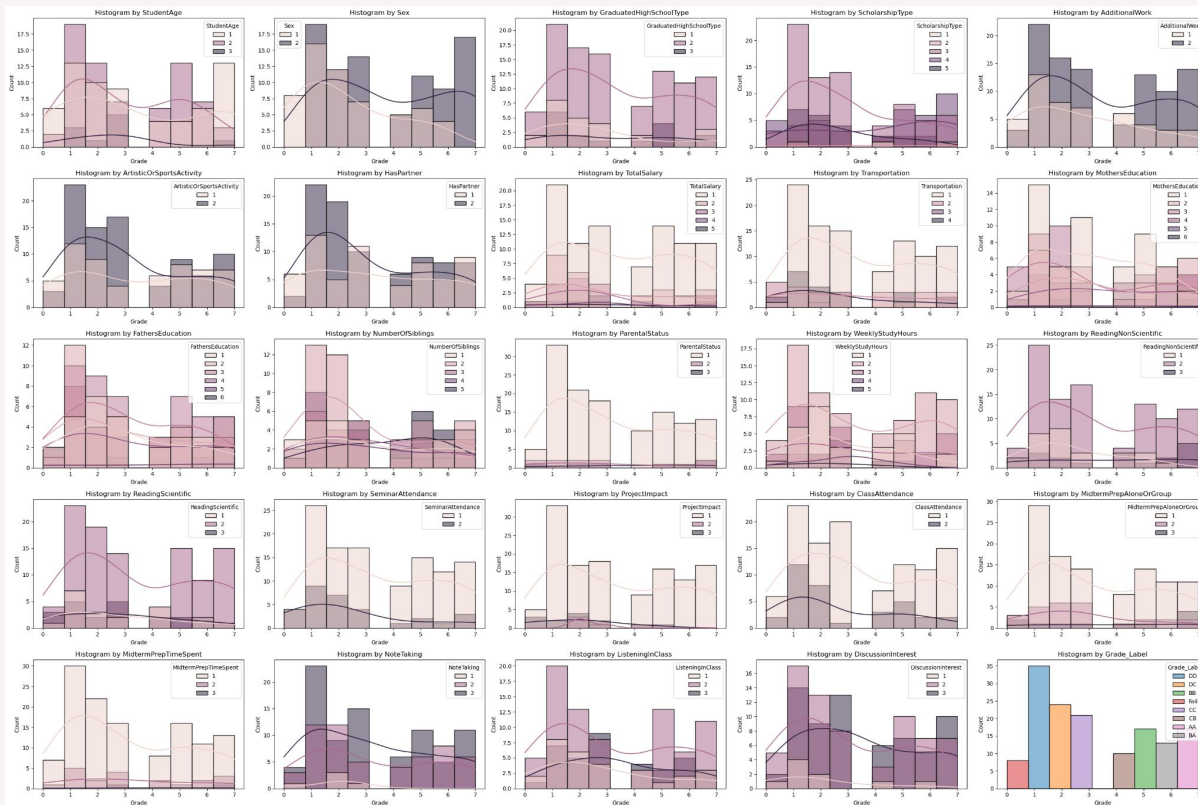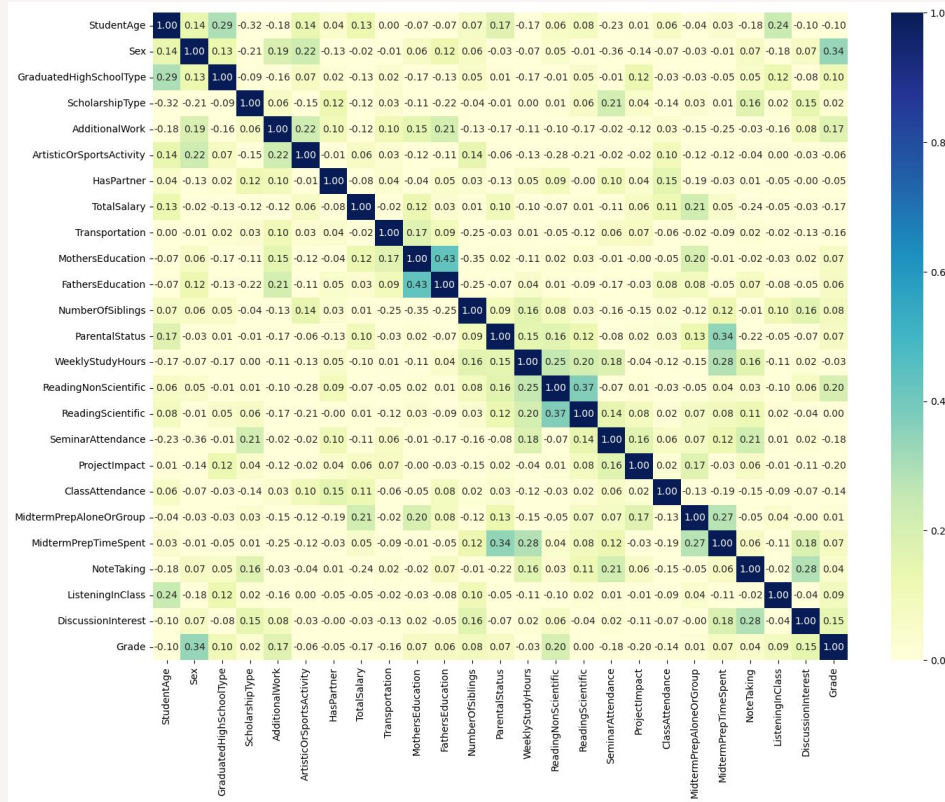
## Step 8
Data is fully clean!

# Exploratory Analysis (1/3)

- **Target Variable: The final grade of the students**
- **Interactive Dashboard:** https://andreawcy.github.io/Project_4_Students-performance-analysis/



Distribution of Grade

# Exploratory Analysis (2/3)

# Exploratory Analysis (3/3)

# Database Design



**Students**
| | |
|---|---|
| **StudentID** | varchar(50) |
| StudentAge | int |
| Sex | int |
| GraduatedHighSchoolType | int |
| ScholarshipType | int |

**StudentActivites**
| | |
|---|---|
| **ActivityID** | int |
| **StudentID** | varchar(50) |
| AdditionalWork | int |
| ArtisticOrSportsActivity | int |
| HasPartner | int |
| Transportation | int |
| WeeklyStudyHours | int |

**AcademicPerformance**
| | |
|---|---|
| **PerformanceID** | int |
| **StudentID** | varchar(50) |
| SeminarAttendance | int |
| ProjectImpact | int |
| ClassAttendance | int |
| MidtermPrepAloneOrGroup | int |
| MidtermPrepTimeSpent | int |
| NoteTaking | int |
| ListeningInClass | int |
| DiscussionInterest | int |
| Grade | int |

**FamilyBackground**
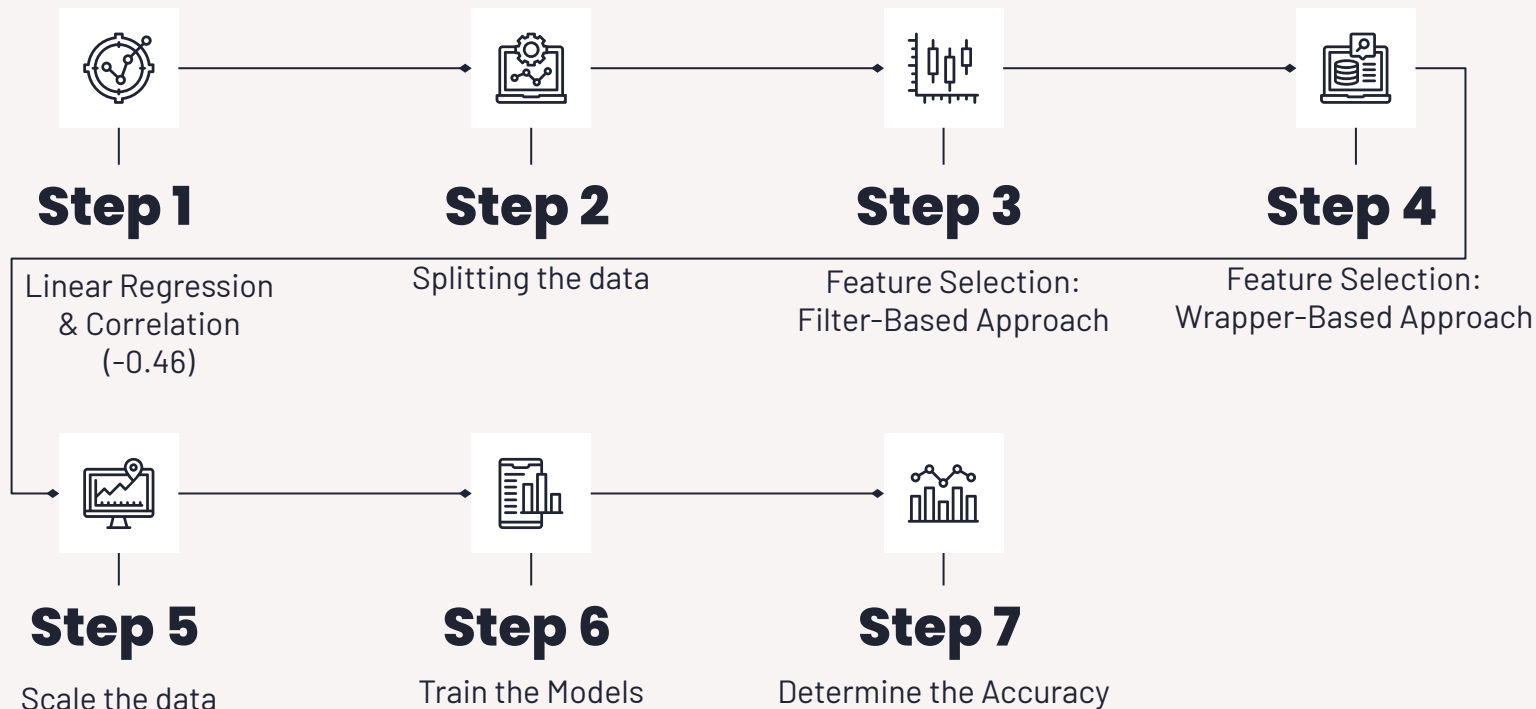| | |
|---|---|
| **BackgroundID** | int |
| **StudentID** | varchar(50) |
| MothersEducation | int |
| FathersEducation | int |
| NumberOfSiblings | int |
| ParentalStatus | int |
| TotalSalary | int |

The database was designed with four tables in mind, connected through primary and foreign keys. Each table stores related information, ensuring data integrity and allowing for more efficient queries.

1. **Students**
2. **Student Activities**
3. **Academic Performance**
4. **Family Background**

# Modelling Experiment Design

**Step 1**

Linear Regression
& Correlation
(-0.46)

**Step 2**

Splitting the data

**Step 3**

Feature Selection:
Filter-Based Approach

**Step 4**

Feature Selection:
Wrapper-Based Approach

**Step 5**

Scale the data

**Step 6**

Train the Models

**Step 7**

Determine the Accuracy

# Filter-Based Approach of Feature Selection



Basically, according to the mutual_info_clasif result above, in these particular train and test buckets, whether you prepped alone or in a group is the most "important" feature to consider.

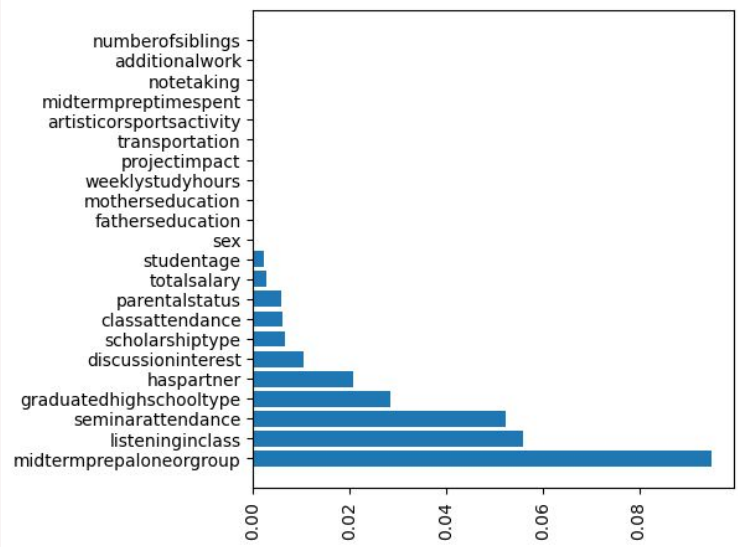## Mutual Information Classification

Select subsets of features based on their relationship with the target.

**Advantages**:

- Faster than using a wrapper method
- More generalized and easier to interpret

**Disadvantage**:

- Each feature is looked at in isolation, i.e. it's prone to discarding useful features that are weak predictors on their own, but useful when combined with others

# Wrapper-Based Approach of Feature Selection

## RandomForestClassifier

A supervised machine learning algorithm based on ensemble learning, a technique that creates multiple models using decision trees and then combines them to produce improved results.

**Advantages**:

- **High accuracy**— because it uses multiple decision trees, then merges them together, it lessens the variation associated with individual trees.
- **Adaptable**— can be applied to many datasets and problem areas.

**Disadvantages**:

- **Computationally expensive**— it uses a lot of computing power and memory.
- **High prediction time**— it can take longer to make predictions than other algorithms, making it ill-fit for real-time applications.

# Wrapper-Based Approach of Feature Selection

## SequentialFeatureSelector

An algorithm that iteratively adds or removes features from a dataset in order to improve the performance of a predictive model.

**1)** Initialize with a predictive model, number of features to select, the scoring metric, and a tolerance for improvement → **2)** Fit the model on full set of features → **3)** Evaluate the model on the training set using the scoring metric → **4)** The feature that most improves the model's cross-validation score is added to the set, or the feature that least reduces the cross-validation score is removed from the set → **5)** Repeat steps 2-4 until the desired number of features has been selected.

**Advantages**:

- **Simple**— it's a simple and efficient algorithm, especially for a wrapper approach.
- **Flexible**— it can be used with any type of predictive model.

**Disadvantages**:

- **Bias**— it can be biased towards features that are highly correlated with the target feature.
- **Computationally Expensive**— it uses a lot of computing power and memory.

# Modeling the Data

1.  **Choose the Feature Selection**

2.  **Scale the data**
    Using StandardScaler()

3.  **Train the Models**

    **Models used:**
    - LinearSVC()
    - SGDClassifier()
    - DecisionTreeClassifier()
    - KNeighborsClassifier()
    - RandomForestClassifier()

\*  More models, such as a Logistic Regression model, can be found in the project notebook.

# Model Evaluation

| | Model | Accuracy | Confusion Matrix |
|---|---|---|---|
| 0 | RandomForestClassifier() | 0.272727 | [[0, 3, 1, 0, 0, 0, 0, 1], [0, 6, 2, 1, 0, 1, ... |
| 1 | KNeighborsClassifier() | 0.250000 | [[1, 2, 1, 0, 0, 0, 0, 1], [0, 5, 2, 0, 0, 1, ... |
| 2 | LinearSVC() | 0.204545 | [[0, 2, 1, 1, 0, 0, 1, 0], [0, 2, 4, 0, 0, 3, ... |
| 3 | SGDClassifier() | 0.204545 | [[1, 2, 1, 0, 0, 0, 0, 1], [0, 2, 4, 0, 0, 3, ... |
| 4 | DecisionTreeClassifier() | 0.159091 | [[0, 2, 2, 0, 0, 1, 0, 0], [0, 1, 2, 1, 3, 2, ... |

After training each model in an iterative loop, I stored the accuracy result and the confusion matrix for each model in a Pandas dataframe. As you can see, my models weren't so accurate...

# Key Findings

Based on the dataset we gathered  we couldn't predict a student's academic performance (e.g., grade point average or success in courses) based on their personal background, study habits, and extracurricular activities because the accuracy was low in every model.

Moreover, there were no key factors that influence the success of a student.

# Limitations

## Data Structure

The dataset is primarily categorical, which may not provide sufficient detail for models to identify complex patterns or make accurate predictions.

## Multicollinearity

In other words, it's difficult to determine the individual effects of each variable on the student's performance, because many of them are dependent on each other.

## Limited Data

While the dataset had a wide breadth of categories, there were only a little under 150 entries in the dataset. This helps to explain the poor accuracy of our models.

# Considerations for Future

**Expand Data Sources:**

Integrate additional data such as student surveys and extracurricular participation records.

**Enhance Data Accuracy:**

Implement automated validation tools and advanced imputation techniques to address data gaps and ensure accuracy.

# THANKS!

**Do you have any questions?**