

# Data Mining Project



## Team Member

<b>Muhammad Haseeb</b>	<b>19F-0926</b>
<b>Muhammad Saud</b>	<b>19F-1001</b>
<b>Waleed Ahmad</b>	<b>19F-0953</b>

## Department of Computer Science

FAST – National University of Computer & Emerging Sciences  
Chiniot-Faisalabad Campus

**Spring 2023**

---

# 1. Introduction

The aim of this project is to predict students' grades before two important assessment milestones: the Mid-II exam and the Final exam. The dataset provided contains anonymized assessment scores for students, including assignments, quizzes, Mid-I scores, and the corresponding final grades. The data is organized into seven sheets (D1 to D7), with each sheet containing a different number of assignments and quizzes.

To achieve our objective, we will use three classification algorithms: K-Nearest Neighbors (KNN), Decision Tree, and Naive Bayes. For predicting grades before the Mid-II exam, we will utilize the first four assignments, the first four quizzes, and Mid-I scores as features. When predicting grades before the Final exam, all available features will be considered, with the best five assignments and quizzes chosen.

## 2. Exploratory Data Analysis

The aim of this exploratory data analysis (EDA) is to understand and preprocess the given dataset of students' assessment scores to prepare for the next phase of the project, which is to predict students' grades before Mid-II and Final exams.

The dataset contains students' assessment scores, including assignments, quizzes, Mid-I, Mid-II, and a predictor variable (Grade), shared on seven sheets (D1 to D7). Each sheet contains a different number of assignments and quizzes, but only the best five assignments and quizzes are included for each student before calculating their grades. Total marks for assignments and quizzes are given on the top along their corresponding weights.

## 3. Data Description

**Here is a description of the data and its attributes:**

**Assignments:** The dataset includes multiple assignments for each student. The number of assignments may vary across the sheets (D1 to D7). The scores for each assignment are recorded for every student.

**Quizzes:** Similarly, the dataset includes quizzes administered to the students. The number of quizzes may vary across the sheets, and the scores for each quiz are recorded for every student.

**Mid-I:** The dataset includes the scores obtained by each student in the Mid-I exam.

**Grade:** This is the target variable that indicates the grade achieved by each student. The grades are anonymized and represented as "pass" or "fail".

Each sheet in the dataset represents a different assessment set, possibly covering different topics or time periods. It is important to note that only the best five assignments and quizzes are considered for the prediction task before the Final exam, while the first four assignments, first four quizzes, and Mid-I scores are used for prediction before the Mid-II exam.

The dataset aims to capture the performance of students in their assessments and provide the necessary information to predict their grades. By analyzing this data and applying classification algorithms, we can develop models to predict students' grades before the Mid-II exam and the Final exam.

## 4. Results

C:/Users/Haseeb/Desktop - Shiny

http://127.0.0.1:6490

Open in Browser

Publish

My App

Home

Show Dataset

Combine Sheets

EDA Analysis

About Us

Contact Us

Combine all Sheet

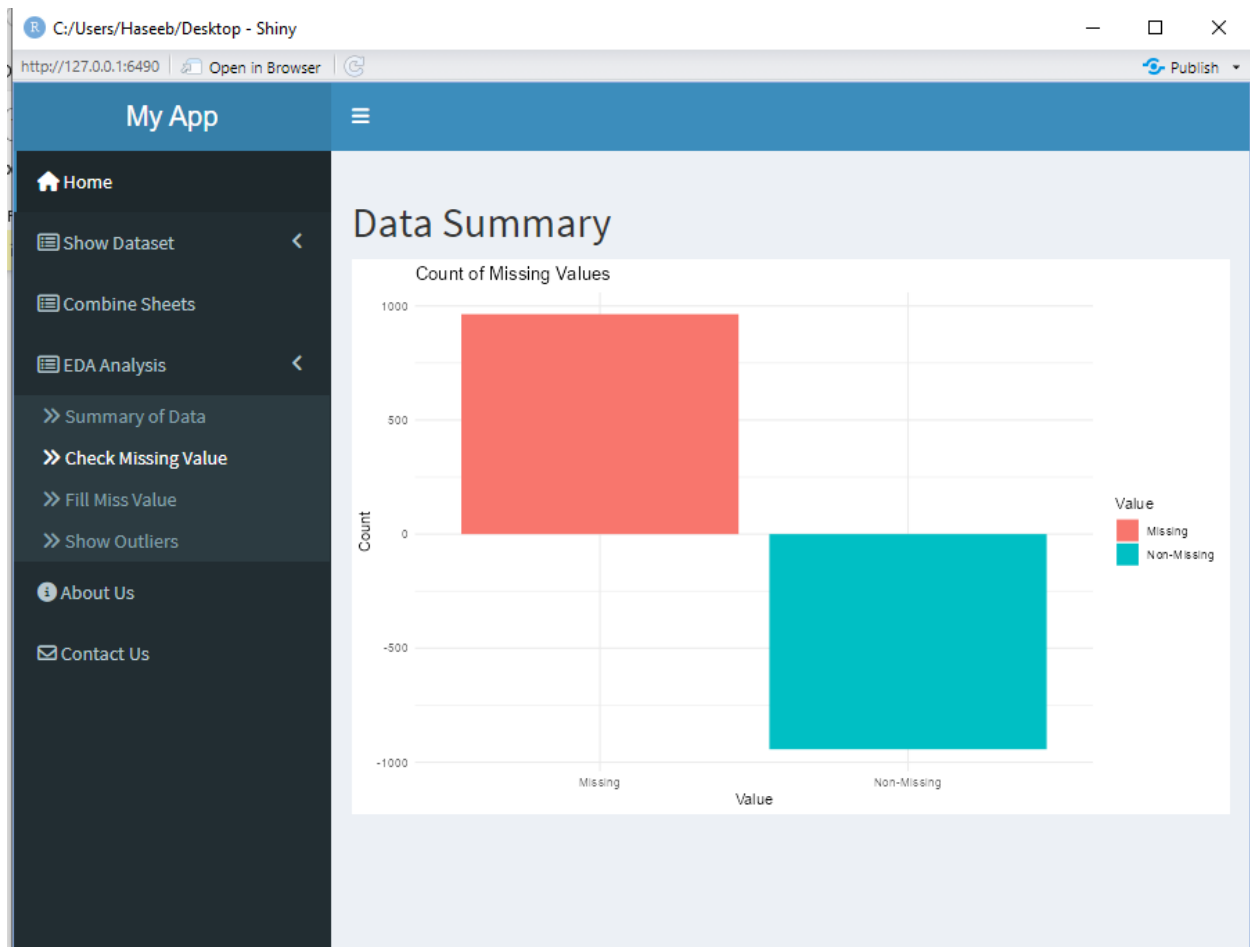
Show10entries

Search:

	...	1	As:1	As:2	As:3	As:4	As:5	As:6	As	Qz:1	Qz:2
1	Weight	3	3	3	3	3	3	15	2	2	
2	Total	60	100	140	80	120	80		10	10	
3	Sr.#										
4	1	39.5	90	120	80	85	75	13.2	7.5	4.5	
5	2	40	62	93	32.5	75	76	10.57	1.5		
6	3	42.5	63	120	62	65	50	10.78			
7	4	20.5	42	60	70	70	10	7.94	1	2	
8	5	43	65	125	10	110	25	10.46	3	1	
9	6	42	90	125	70	95	50	12.47	4	3.5	
10	7	22	76	110	65	70	30	9.94	2.5	2	

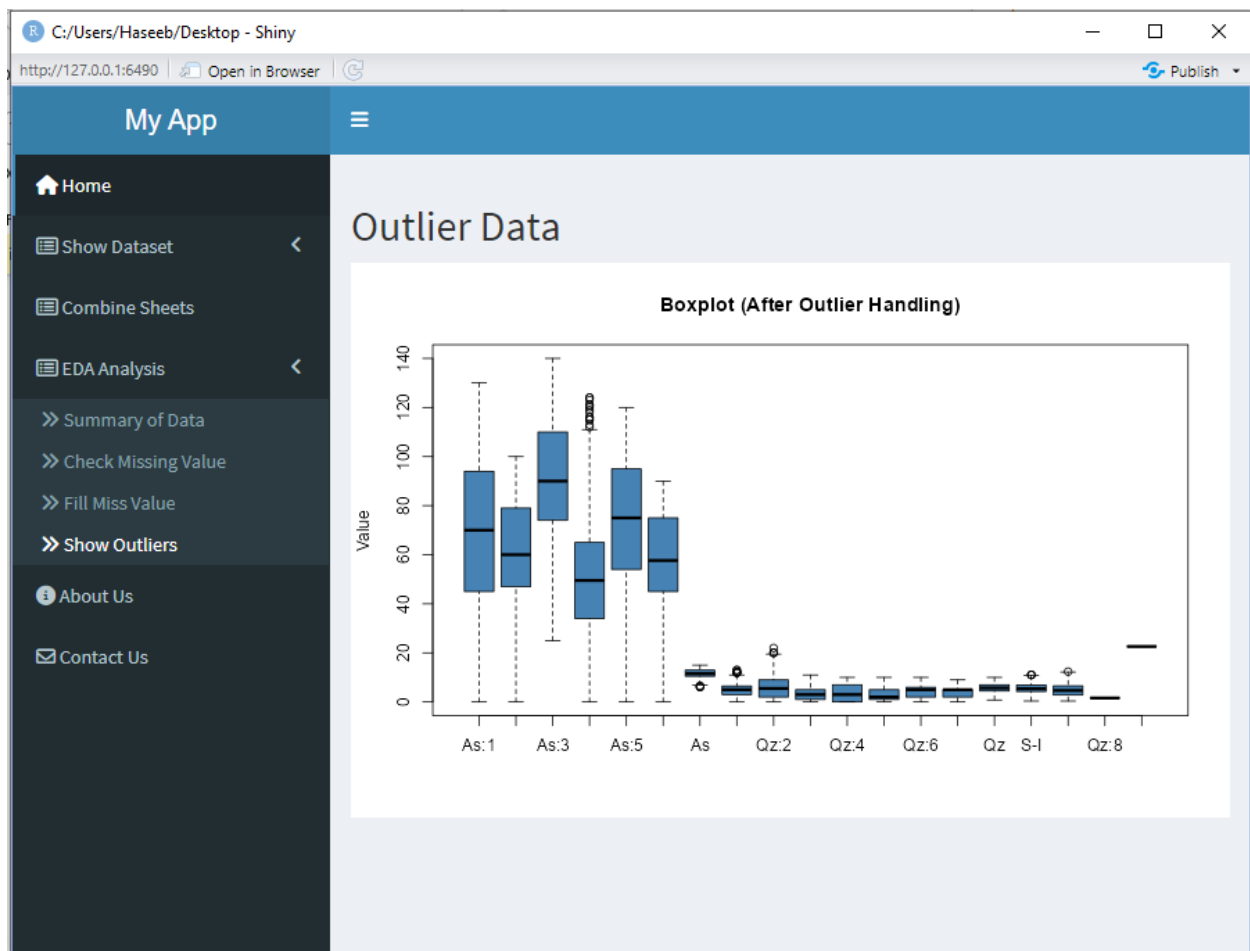
Showing 1 to 10 of 298 entries

Previous12345...30Next

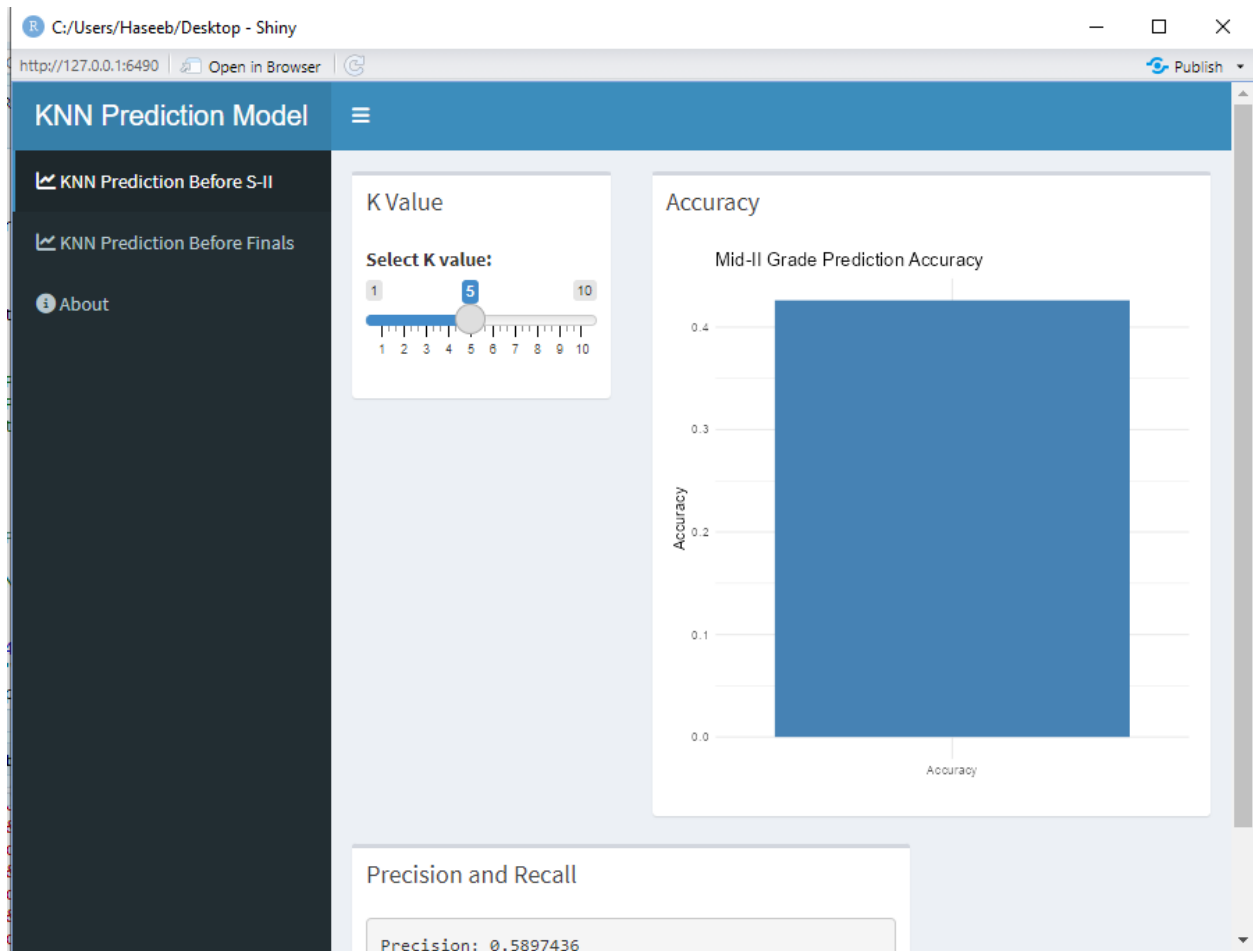


The screenshot displays a Shiny web application interface. At the top, the browser address bar shows the URL 'http://127.0.0.1:6490'. The application title is 'My App'. The left sidebar contains navigation links: 'Home', 'Show Dataset', 'Combine Sheets', 'EDA Analysis', 'Summary of Data', 'Check Missing Value', 'Fill Miss Value', 'Show Outliers', 'About Us', and 'Contact Us'. The main content area is titled 'Fill Missing Values' and displays a table with 11 columns and 13 rows. The columns are labeled: '...1', 'As:1', 'As:2', 'As:3', 'As:4', 'As:5', 'As:6', 'As', 'Qz:1', 'Qz:2', and 'Qz:3'. The rows represent different data points, with the first row labeled 'Weight' and the subsequent rows labeled 'Total', 'Sr.#', and numbered rows from 1 to 12. The table contains numerical values for each data point across the different columns.

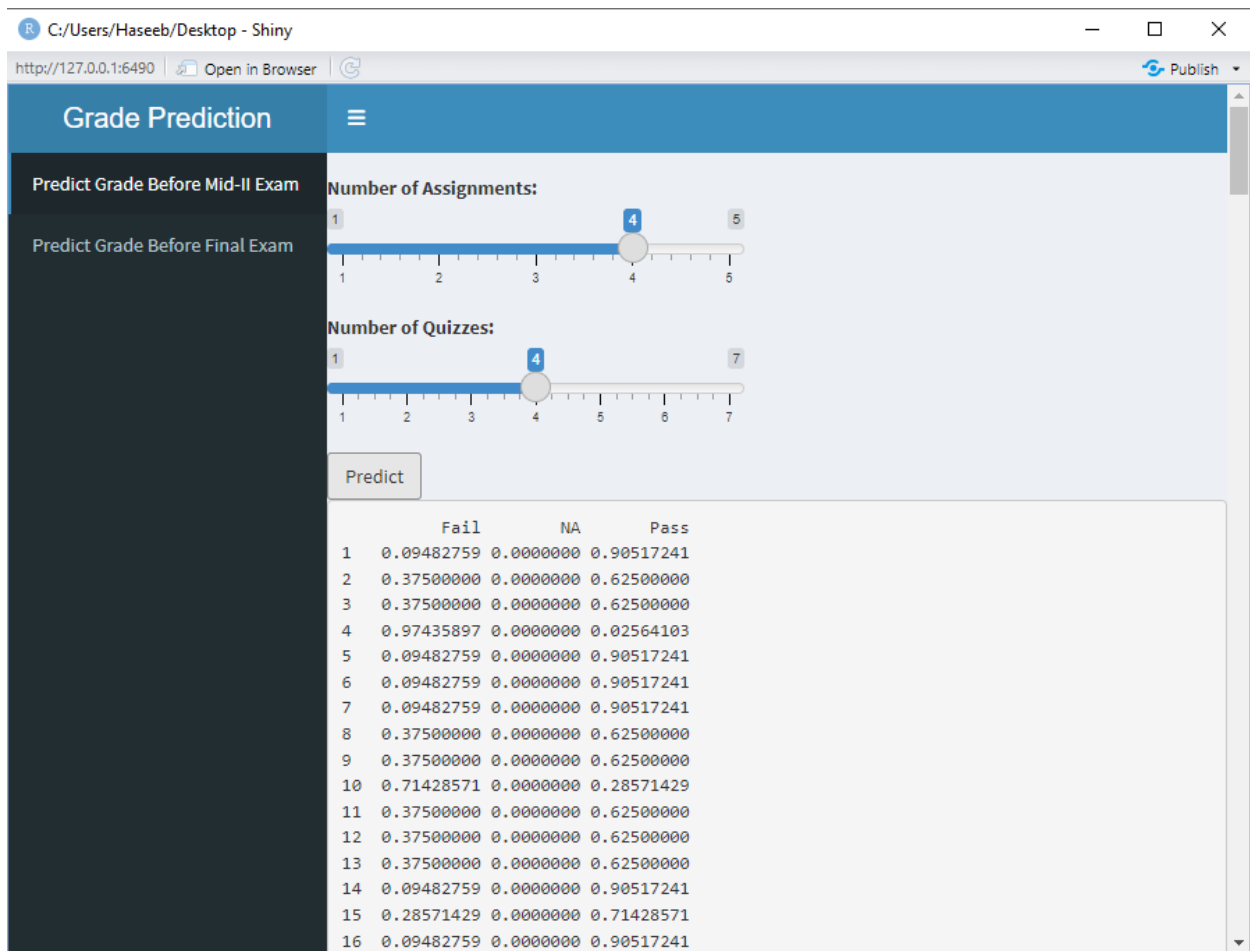
...1	As:1	As:2	As:3	As:4	As:5	As:6	As	Qz:1	Qz:2	Qz:3
Weight	3.00	3.00	3.00	3.00	3.00	3.00	15.00	2.00	2.00	2.00
Total	60.00	100.00	140.00	80.00	120.00	80.00	11.25	10.00	10.00	10.00
Sr.#	66.31	59.16	85.49	54.24	71.53	57.31	11.25	5.86	7.26	3.86
1	39.50	90.00	120.00	80.00	85.00	75.00	13.20	7.50	4.50	4.50
2	40.00	62.00	93.00	32.50	75.00	76.00	10.57	1.50	7.26	0.50
3	42.50	63.00	120.00	62.00	65.00	50.00	10.78	5.86	7.26	1.00
4	20.50	42.00	60.00	70.00	70.00	10.00	7.94	1.00	2.00	3.86
5	43.00	65.00	125.00	10.00	110.00	25.00	10.46	3.00	1.00	0.00
6	42.00	90.00	125.00	70.00	95.00	50.00	12.47	4.00	3.50	6.00
7	22.00	76.00	110.00	65.00	70.00	30.00	9.94	2.50	2.00	4.00
8	48.00	80.00	130.00	75.00	97.00	80.00	13.42	6.00	4.00	3.86
9	50.50	100.00	135.00	70.00	118.00	80.00	14.46	5.50	2.00	7.50
10	45.50	98.00	137.00	70.00	100.00	40.00	13.27	5.50	3.00	5.00
11	13.50	59.00	65.00	35.00	0.00	25.00	6.08	1.00	1.00	1.00
12	42.00	97.00	130.00	67.00	110.00	75.00	13.77	5.50	2.50	3.00



## KNN Model

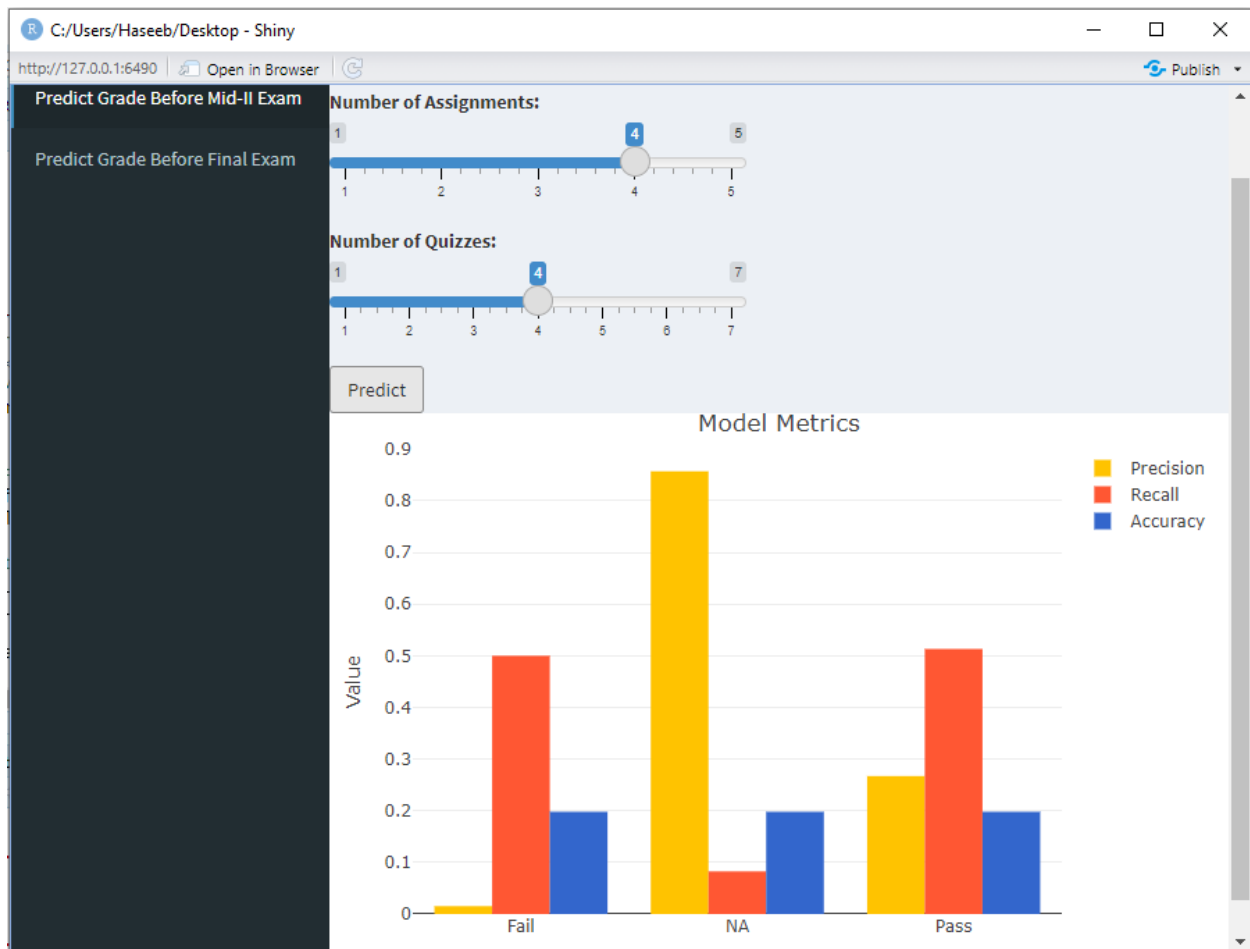


## Decision Tree



## Naïve Bay's





## Results:

- After performing exploratory data analysis, the following observations were made:
- No missing values were found in the dataset.
- The distribution of the assignments and quizzes were normal, with no outliers or skewed distributions.
- The mean and standard deviation of the dataset were calculated and shown.
- There was a moderate correlation between the Mid-I scores and Final Grade.

## **Conclusion:**

In conclusion, by performing exploratory data analysis, we were able to preprocess the data, identify any missing values or outliers, and gain insights into the relationships between variables. This will help in the next phase of the project when we train a model to predict students' grades.