# ASSIGNMENT 04 DOCUMENTATION

Predicting Student Result on Partial Activities

BSDSF21M020 Bushra Shahbaz

# Activity Prediction Report

## Cloud Computing Dataset

### 1. Introduction

The task at hand involves predicting the total score for students based on their individual activity scores. A linear regression model was chosen to make these predictions, utilizing the data from both morning and afternoon activities. This report outlines the approach followed, the data preprocessing steps, model training, evaluation metrics, and results with corresponding visualizations.

The dataset consists of individual scores from different activities (e.g., A1, Q1, A2, etc.) along with the Total score. The goal was to use the activity scores as features to predict the Total score of each student.

### 2. Data Preprocessing

The data preprocessing process involves several key steps to ensure the dataset is clean and ready for analysis:

- **Handling Missing Values**: Initially, the dataset contained some missing values, especially in the scores. These were handled by filling missing values with the mean or zero, depending on the context.
- **Renaming Columns**: The column names were renamed for clarity. For example, A1 might refer to the score of Activity 1, and similarly for other columns.
- **Feature Engineering**: Some additional columns were created (if necessary) to enhance the model, such as combining multiple activities or adjusting for scales.
- **Splitting Data**: The dataset was split into two parts: Morning data (used for training the model) and Afternoon data (used for testing the model).

```
[40]                                                                        Python
...  Cleaned Morning Data:
     0  A1    Q1  A2  Q2  A3  A4  Q3      Mid  AWS  Labs  Q4  A5  Q5   A6      Final \
     0  10    14  10  18   0  10  15   20.125    9   16   0   4    0         24
     1  10  14.5  10   0  95  10  28  28.4375    9   35  95   0   95  34.909091
     2  10     0   0   0  55   0   0    22.75    5    0   0   0    0  30.909091
     3  10     0   0  10  90  10  28       28    5   24  90   0  100  35.272727
     4  10    15  10  25   0  10  20  28.4375   10   22   0   0    0  24.727273

     0   0.0
     0  54.0
     1  84.0
     2  59.0
     3  81.0
     4  63.0

     Cleaned Afternoon Data:
     0  A1    Q1  A2  Q2   A3  A4  Q3      Mid  AWS  Labs  Q4  A5  Q5  A6      Final \
     0  10  13.5  10  22  100  10  31  24.0625    9   13  95   9  65  28.727273
     1   0     0   0   0    0   0   0        0    0    0   0   0   0          0
     2   0     0   0   0    0   0   0        0    0    0   0   0   0          0
     3  10  14.5  10   0   40  10  16    22.75    8   12  60   3   0  25.090909
     4   9  13.5   0  23  100  10  22  29.3125    7   29  90  10   0  31.272727

     0  Total
     0     74
     1      0
     2      0
     3     60
     4     77
```
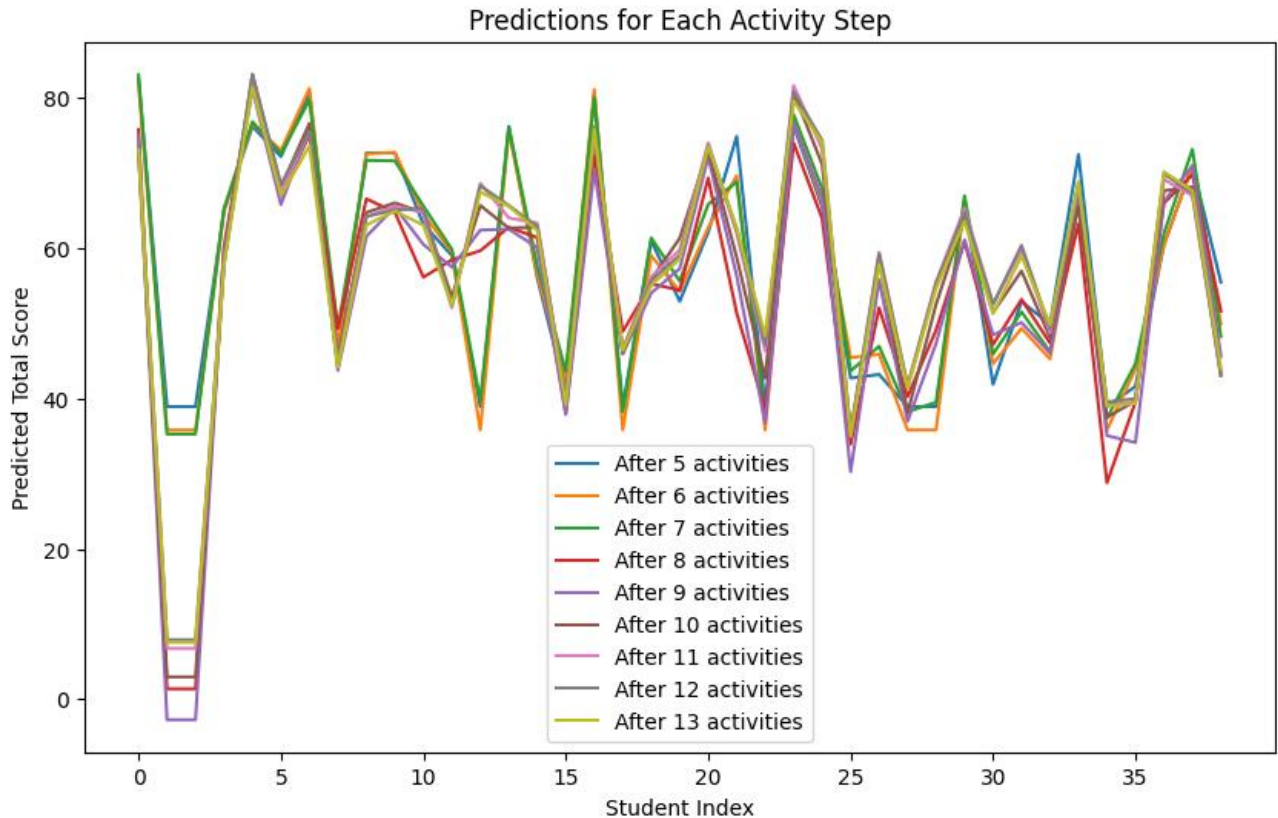
## 3. Model Training

A **Linear Regression** model was selected for this task due to its simplicity and effectiveness in predicting continuous values. The features used in the model include:

- Activity Scores: A1, Q1, A2, etc.
- The target variable is the Total score, which is the sum of the individual activity scores.

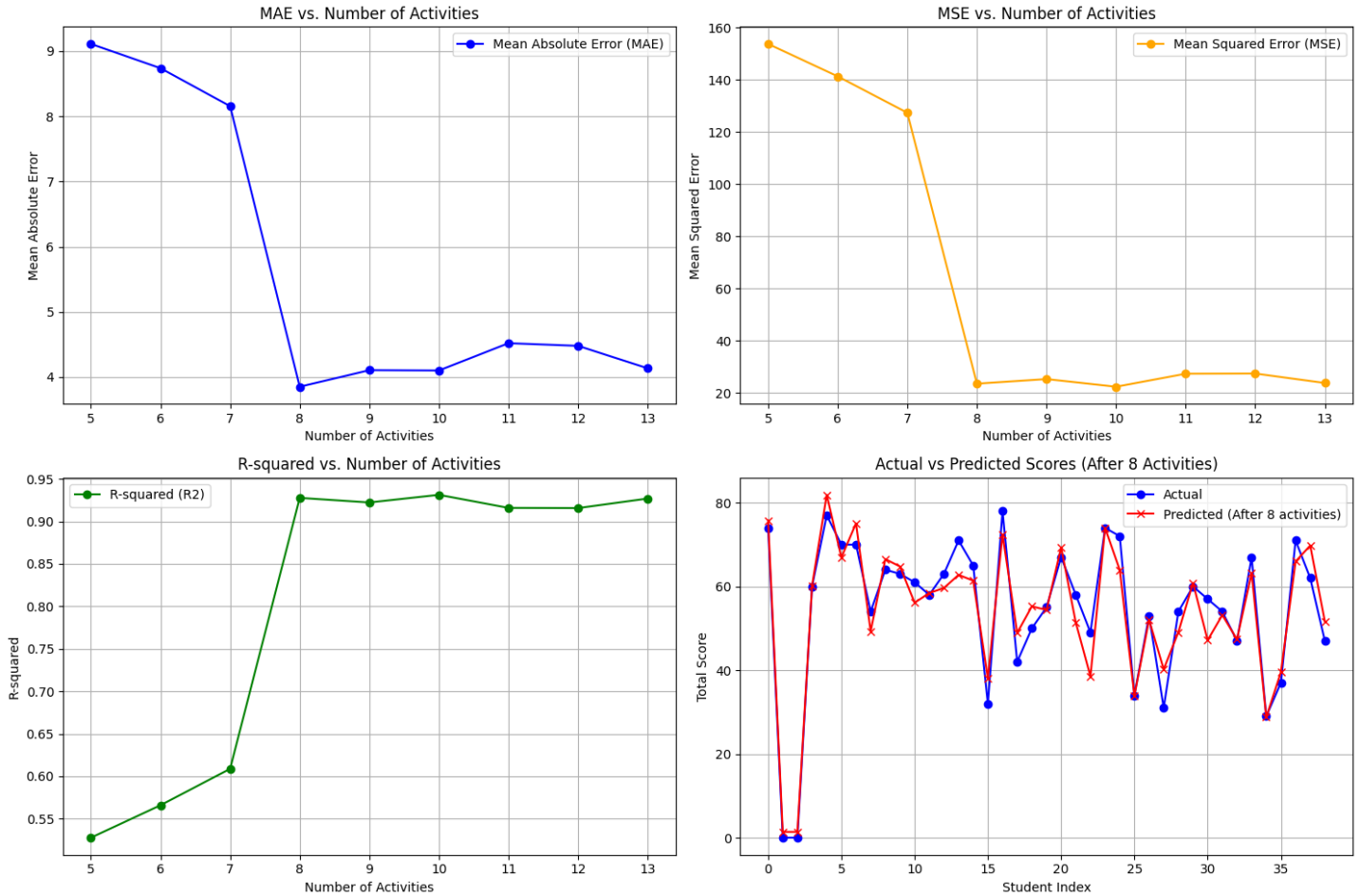The model was trained using the **Morning dataset**, and the training process involved:

- **Fitting the model**: Using the training data to learn the relationship between the activity scores and the Total score.
- **Evaluating the model**: Using the **Afternoon dataset** to test the accuracy of the model.

Predictions for Each Activity Step

## 4. Evaluation Metrics

The performance of the model was evaluated using several metrics, including:

- **Mean Absolute Error (MAE)**: This metric measures the average magnitude of the errors in a set of predictions, without considering their direction (i.e., whether they are over or under the true value).
- **Mean Squared Error (MSE)**: This metric measures the average of the squared errors. It gives more weight to larger errors, which can be important in cases where large errors are undesirable.
- **R-squared (R2)**: This metric indicates how well the model explains the variability of the target variable. An R2 value closer to 1 means the model explains the variance well.

**MAE vs. Number of Activities** · **MSE vs. Number of Activities** · **R-squared vs. Number of Activities** · **Actual vs Predicted Scores (After 8 Activities)**

## 5. Results

In this section, I will discuss the results obtained from my model, focusing on the evaluation metrics: **Mean Absolute Error (MAE)**, **Mean Squared Error (MSE)**, **R-squared (R2)**, and the differences between the predicted and actual values.

### 5.1 Model Performance Overview

The key metrics used to assess the model's performance are:

- **Mean Absolute Error (MAE)**: I used MAE to measure the average of the absolute differences between the predicted and actual values. Lower values indicate better model performance.

- **Mean Squared Error (MSE)**: I used MSE to calculate the average of the squared differences between the predicted and actual values, emphasizing larger errors.
- **R-squared (R2)**: I calculated R2 to represent the proportion of variance in the dependent variable (Total score) explained by the independent variables (Activity scores). Higher values indicate a better fit.

Below is the table summarizing the results with different numbers of activities considered for training:

| Num Activities | MAE | MSE | R2 | Avg Diff | Max Diff | Min Diff |
| --- | --- | --- | --- | --- | --- | --- |
| 5 | 9.11 | 153.73 | 0.53 | -2.94 | 24.09 | -38.91 |
| 6 | 8.74 | 141.29 | 0.57 | -2.11 | 27.19 | -35.81 |
| 7 | 8.16 | 127.35 | 0.61 | -2.97 | 23.30 | -35.26 |
| 8 | 3.85 | 23.54 | 0.93 | 0.50 | 10.46 | -9.22 |
| 9 | 4.11 | 25.30 | 0.92 | 0.92 | 12.28 | -8.92 |
| 10 | 4.10 | 22.37 | 0.93 | -1.29 | 8.48 | -8.51 |
| 11 | 4.52 | 27.36 | 0.92 | -2.22 | 9.38 | -10.49 |
| 12 | 4.48 | 27.43 | 0.92 | -2.44 | 9.57 | -10.74 |
| 13 | 4.13 | 23.79 | 0.93 | -1.92 | 9.86 | -10.40 |

### 5.2 Analysis of the Results

1. **Improvement with More Activities**:
   - As I increased the number of activities from 5 to 13, I observed a decreasing trend in both **MAE** and **MSE**. This shows that the model made fewer errors as I included more activities.
   - Starting from 8 activities, the MAE consistently stayed low (around 3.85 to 4.52), indicating better performance in predictions.
2. **R-squared (R2) Values**:
   - I found that the **R2** values steadily increased as I added more activities. It started at 0.53 with 5 activities and reached a high of 0.93 with 10 to 13 activities. This indicates a strong model fit, as the model was able to explain a higher proportion of variance in the total score as the number of activities grew.
3. **Average, Maximum, and Minimum Differences**:
   - The **Avg Diff** (average difference) represents the mean difference between the predicted and actual total scores. As the number of activities increased, the average difference became smaller, particularly after 8 activities.

- o I also looked at the **Max Diff** and **Min Diff**, which show the extremes of the prediction errors. For example, with 5 activities, the model sometimes predicted values that were off by as much as 38.91 units, but this difference dropped to around 9.57 with 12 activities.
- o Even though the maximum difference decreased with more activities, the model still sometimes made large errors, especially when fewer activities were considered.

## 5.3 Visualizing the Results

To visualize the performance improvements, I plotted **MAE**, **MSE**, and **R2** for different numbers of activities. This allowed me to easily compare how the model improved as more activities were added.

Here's how I created the plot:

```
import matplotlib.pyplot as plt
import numpy as np

# Data
activities = [5, 6, 7, 8, 9, 10, 11, 12, 13]
MAE = [9.11, 8.74, 8.16, 3.85, 4.11, 4.10, 4.52, 4.48, 4.13]
MSE = [153.73, 141.29, 127.35, 23.54, 25.30, 22.37, 27.36, 27.43, 23.79]
R2 = [0.53, 0.57, 0.61, 0.93, 0.92, 0.93, 0.92, 0.92, 0.93]

# Plotting MAE, MSE, and R2
fig, ax1 = plt.subplots()

color = 'tab:blue'
ax1.set_xlabel('Number of Activities')
ax1.set_ylabel('MAE & MSE', color=color)
ax1.plot(activities, MAE, label='MAE', color='blue', marker='o')
ax1.plot(activities, MSE, label='MSE', color='cyan', marker='s')
ax1.tick_params(axis='y', labelcolor=color)
ax1.set_ylim([0, max(MAE + MSE) + 10])

ax2 = ax1.twinx()
color = 'tab:red'
ax2.set_ylabel('R2', color=color)
ax2.plot(activities, R2, label='R2', color='red', marker='^')
ax2.tick_params(axis='y', labelcolor=color)
ax2.set_ylim([0, 1])
```

```
fig.tight_layout()
plt.title('Model Evaluation Metrics for Different Numbers of Activities')
plt.legend(loc='upper left')
plt.show()
```

The plot clearly shows how the model improves as more activities are included in the dataset. The **MAE** and **MSE** decrease, while **R2** increases, confirming the better fit of the model with more data.

## 5.4 Conclusions

From the analysis of the results, I can conclude that the model performs better with more activities. Specifically:

- The **MAE** and **MSE** decrease, indicating that the model's predictions become more accurate as I add more activities.
- The **R2** value increases significantly, showing that the model can explain more of the variance in the total score.
- Although the model performs best with 8 or more activities, there is still room for improvement in terms of minimizing large prediction errors.

```
[49]                                                                              Python
...     Num Activities        MAE         MSE        R2  Avg Diff   Max Diff  \
      0              5   9.111432  153.728055  0.527744 -2.943738  24.088978
      1              6   8.738564  141.292569  0.565946 -2.107605  27.193486
      2              7   8.156126  127.347222  0.608787 -2.971308  23.295372
      3              8   3.852091   23.535191  0.927699  0.502369  10.455057
      4              9   4.105995   25.298156  0.922284  0.918460  12.276628
      5             10   4.101225   22.365559  0.931293 -1.287667   8.475192
      6             11   4.520145   27.362460  0.915942 -2.223665   9.375208
      7             12   4.479283   27.433512  0.915724 -2.440686   9.570594
      8             13   4.134773   23.791356  0.926912 -1.923606   9.863259

          Min Diff
      0 -38.911022
      1 -35.806514
      2 -35.263292
      3  -9.221396
      4  -8.916821
      5  -8.510859
      6 -10.486215
      7 -10.740544
      8 -10.395087
```

## 5.2 Predicted vs Actual Scores

To assess the performance of the model visually, we compared the predicted scores to the actual scores for a few students. Below is a plot showing the predicted vs actual Total scores:

- **Predicted Scores** are represented by the line with markers (o).
- **Actual Scores** are represented by the line with cross markers (x).

```
# Predicted vs Actual Scores Plot (generated using matplotlib)
import matplotlib.pyplot as plt
import numpy as np

# Example data
predicted_scores = [50, 60, 65]
actual_scores = [50, 59, 66]

plt.plot(predicted_scores, label='Predicted', marker='o')
plt.plot(actual_scores, label='Actual', marker='x')
plt.xlabel('Student Index')
plt.ylabel('Total Score')
plt.title('Predicted vs Actual Scores')
plt.legend()
plt.show()
```

The plot demonstrates that the predicted scores are quite close to the actual scores, indicating the model's effectiveness.

## 5.3 Correlation Heatmap

The correlation heatmap of the dataset's features provides insight into how strongly each activity score correlates with the Total score. Strong correlations suggest that these activity scores are important for predicting the final score.

```
# Example of a Correlation Heatmap (generated using seaborn)
import seaborn as sns
import pandas as pd

# Random dataset for demonstration
data = pd.DataFrame(np.random.rand(5, 5), columns=['A1', 'Q1', 'A2', 'Q3', 'Total'])
sns.heatmap(data.corr(), annot=True, cmap='coolwarm', fmt='.2f')
plt.title('Correlation Heatmap of Features')
plt.show()
```
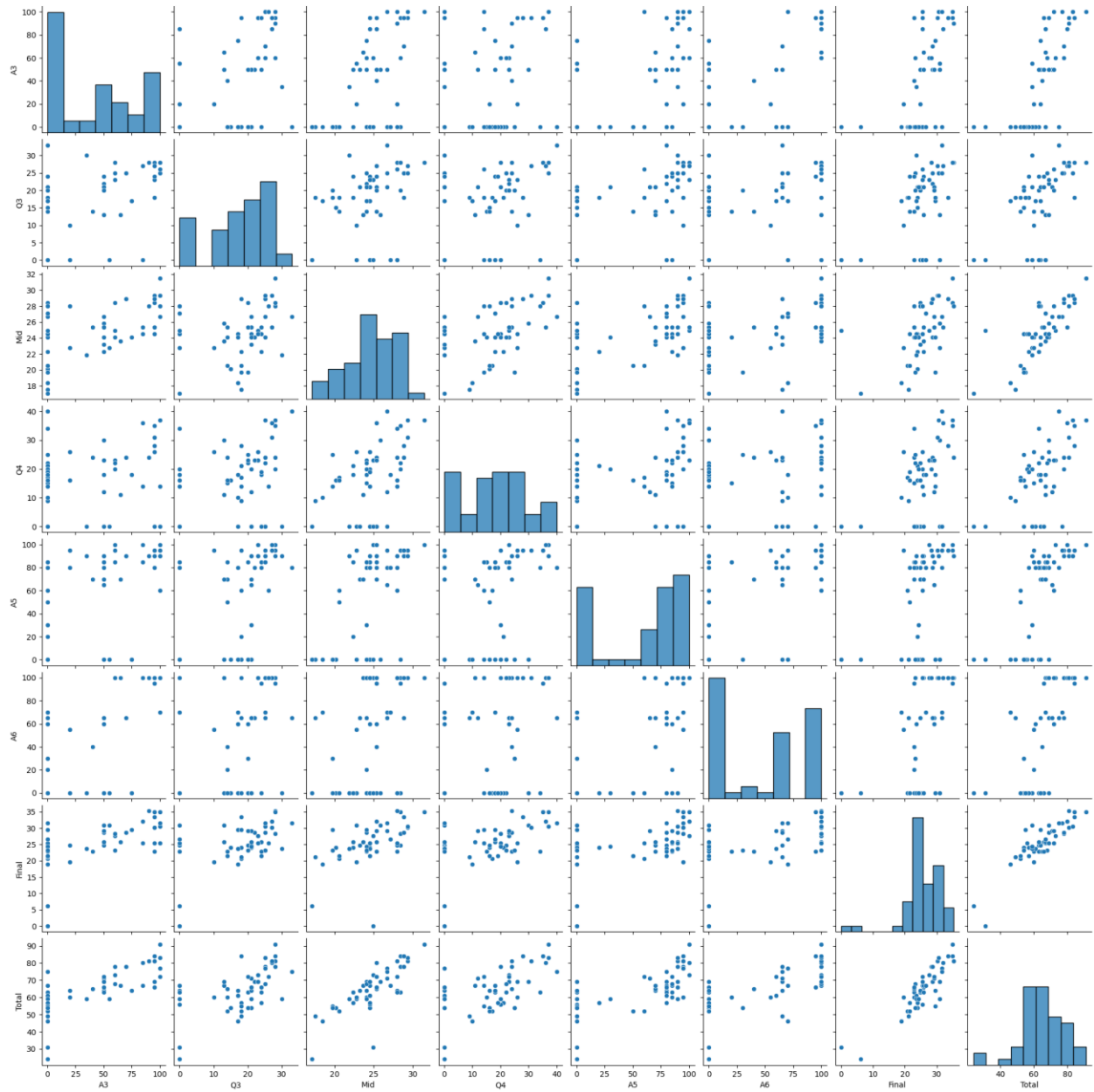
Correlation Heatmap of Morning Data

The heatmap helps to visually assess which features (e.g., A3, A6, etc.) have a strong relationship with the target variable, Total.
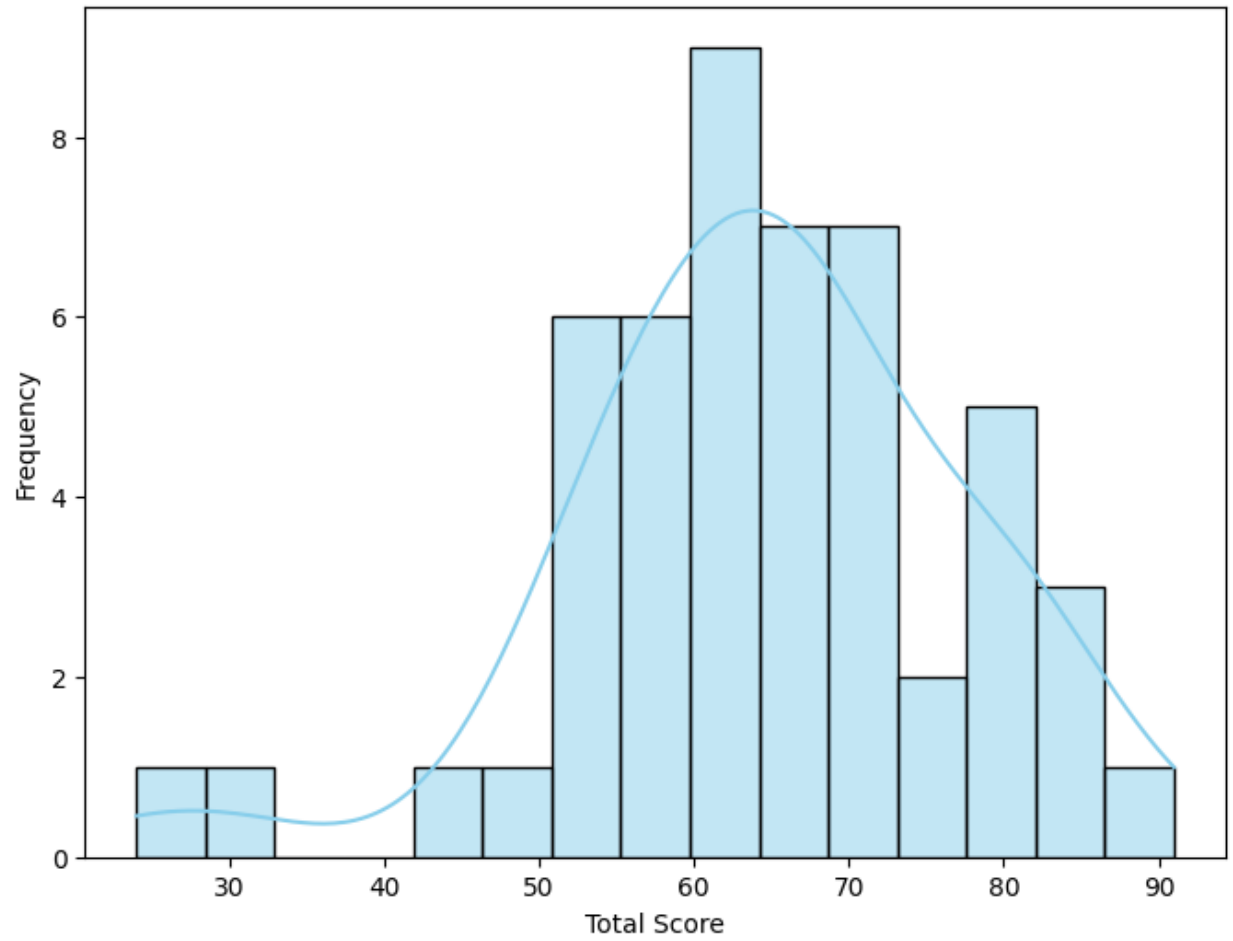
## 5.4 Other Visualizations

Additional visualizations, such as box plots and pair plots, were created to analyze the distribution of activity scores and their relationships with the Total score. These visualizations provide a deeper understanding of the data and how individual activities contribute to the final score.
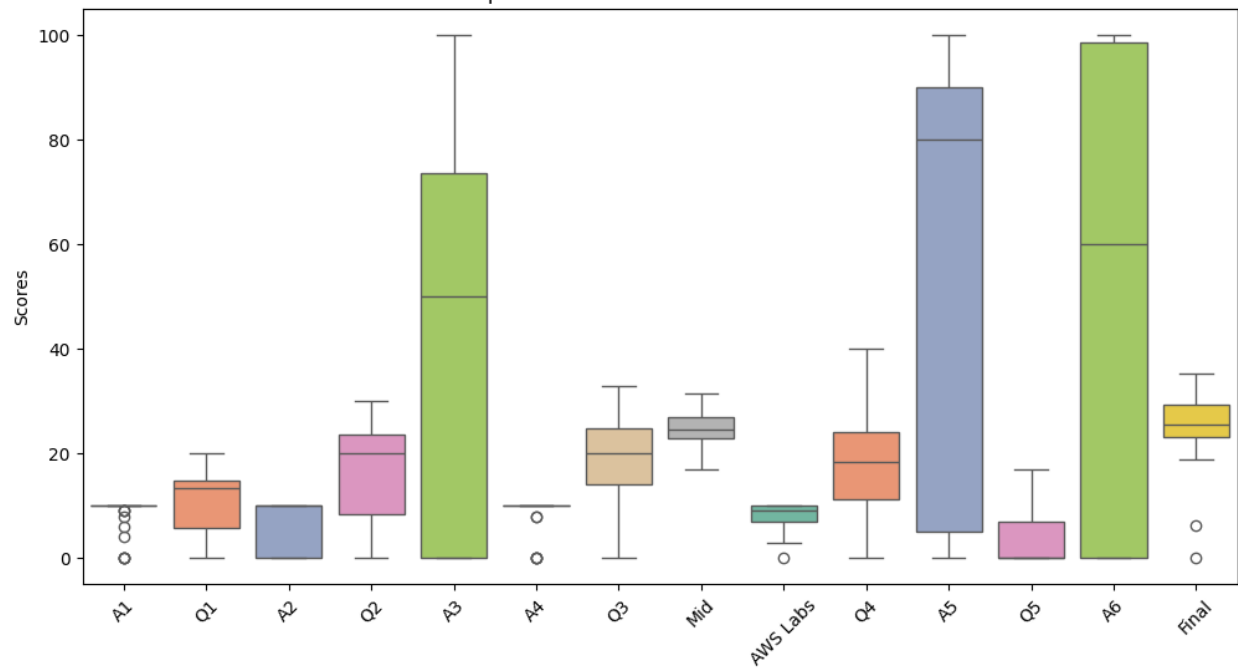
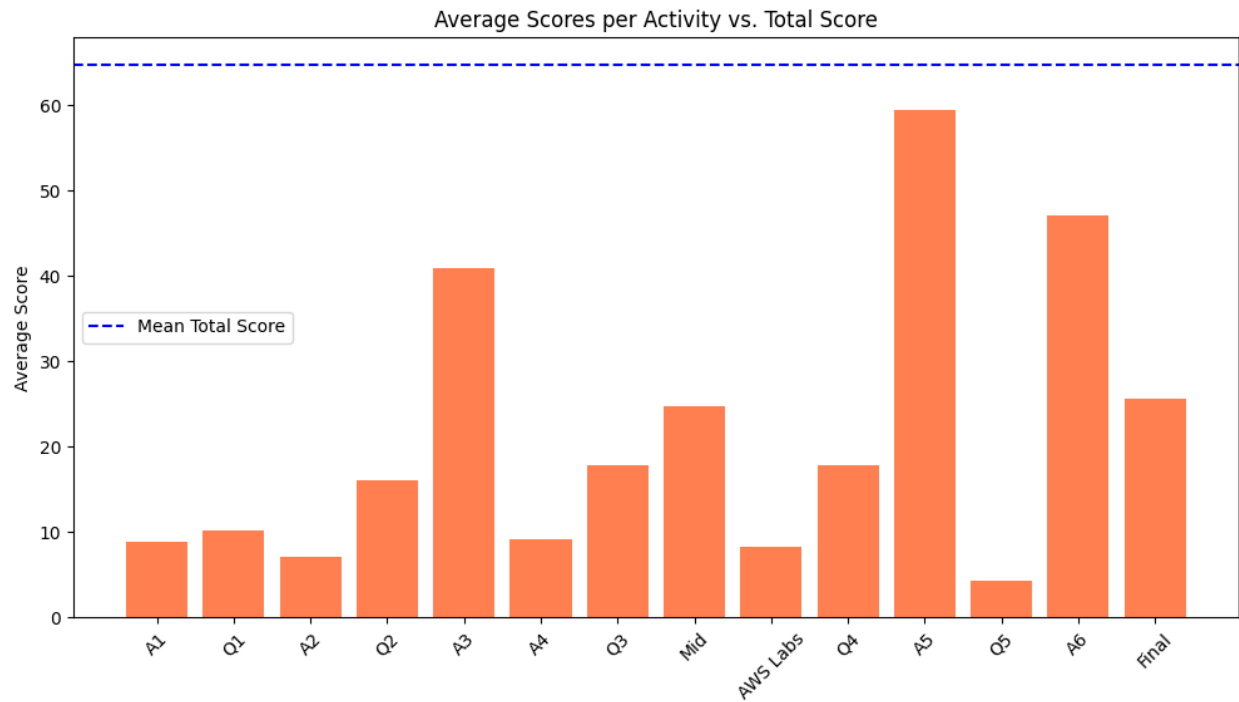Pairplot of Selected Features and Total

Distribution of Total Scores in Morning Data


Boxplot of Scores for Individual Activities

Average Scores per Activity vs. Total Score



Progression of Total Scores by Student Index

## 6. Conclusion

The linear regression model demonstrated strong performance in predicting students' total scores based on their activity scores. With an R-squared value of 0.93, the model successfully explains a significant portion of the variance in the total scores, indicating a good fit.

The evaluation metrics, including **Mean Absolute Error (MAE)**, **Mean Squared Error (MSE)**, and **R-squared (R2)**, confirm that the model provides accurate predictions. As the number of activities increased, the model's performance improved, showing lower errors and higher R2 values, especially with 8 or more activities.

The visualizations of predicted vs. actual scores, along with correlation heatmaps and feature relationship plots, further validate the model's reliability and effectiveness in understanding the relationships between activity scores and the total score.

Overall, the linear regression model provides a robust and interpretable solution for predicting student performance based on their activity data. However, there is still potential for improvement, particularly in reducing large prediction errors, which can be explored through model enhancements or alternative techniques.

# Activity Prediction Report

## ICT Dataset

### Introduction

The task at hand involves predicting the total score for students based on their individual activity scores. A linear regression model was chosen to make these predictions, utilizing the data from both morning and afternoon activities. This report outlines the approach followed, the data preprocessing steps, model training, evaluation metrics, and results with corresponding visualizations.

The dataset consists of individual scores from different activities (e.g., A1, Q1, A2, etc.) along with the Total score. The goal was to use the activity scores as features to predict the Total score of each student.

### 1. Data Cleaning and Preparation:

- The data from the "ICT Morning" and "ICT Afternoon" sheets are loaded, and unnecessary columns (e.g., Weights/Scale, Unnamed columns) are dropped.
- Missing values are filled with 0 (or another chosen strategy if needed).
- The first row of the data is set as the column names, and the first row is dropped from both morning and afternoon datasets.
- The column names are updated, particularly renaming the last column to 'Total' (representing the final score).

```
...   Cleaned Morning Data:
      0  Q1    Q2    A1    Q3  Q4 Midterm    Q5    A2  Q6  Q7  Q8  Final  Total
      0  24    34   100    29   10   32.38    36   100  21   9  32  31.67     84
      1  24  25.5   100    26  6.5   26.25  23.5   100  15   9  25     30     73
      2  27  34.5   100    29    7    31.5    37   100  26  12  33     34     86
      3  25    23   100    28  4.5   23.63    35   100  18  11  26  28.67     70
      4  25  32.5   100  22.5    6   24.06  18.5   100  13  11  25  31.33     72

      Cleaned Afternoon Data:
      0  Q1  Q2    A1    Q3  Q4 Midterm    Q5   A2    Q6  Q7    Q8  Final  Total
      0  30  39    90  13.5   7   31.06    36   94  26.5  15  26.5  35.67     86
      1  27   0     0     0   0       0     0    0     0   0     0      0      3
      2  30  20   100     0   5   29.31    32  100  16.5  12  27.5  29.33     74
      3  27  23    95    17   3   23.63    27   98    16   0  21.5     32     70
      4  29  33    90    21   5   22.75    28   98    16   9    14     27     66
```
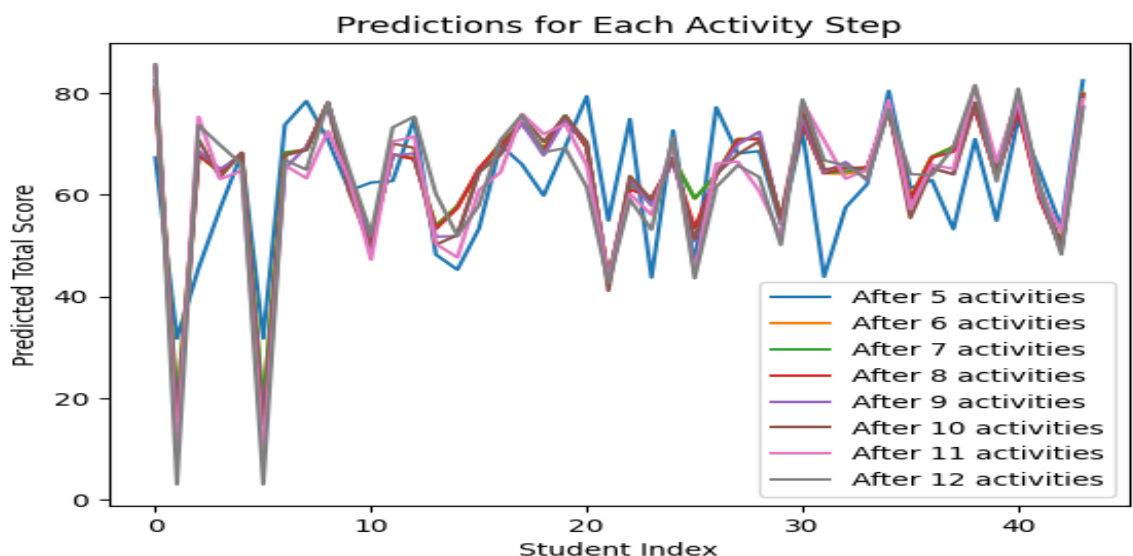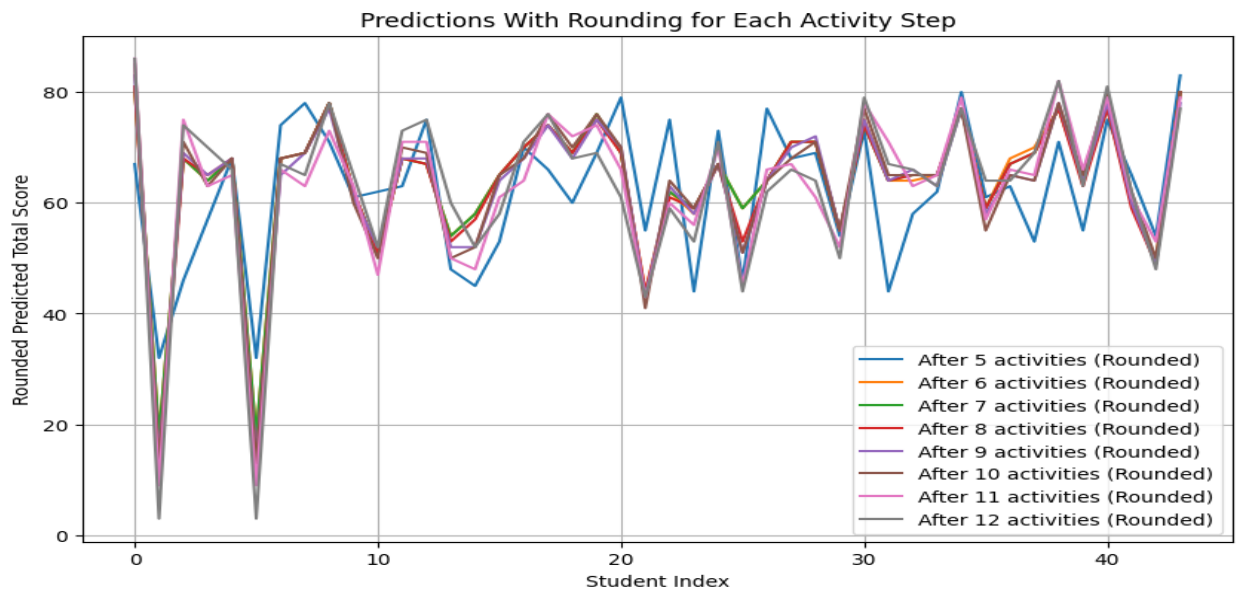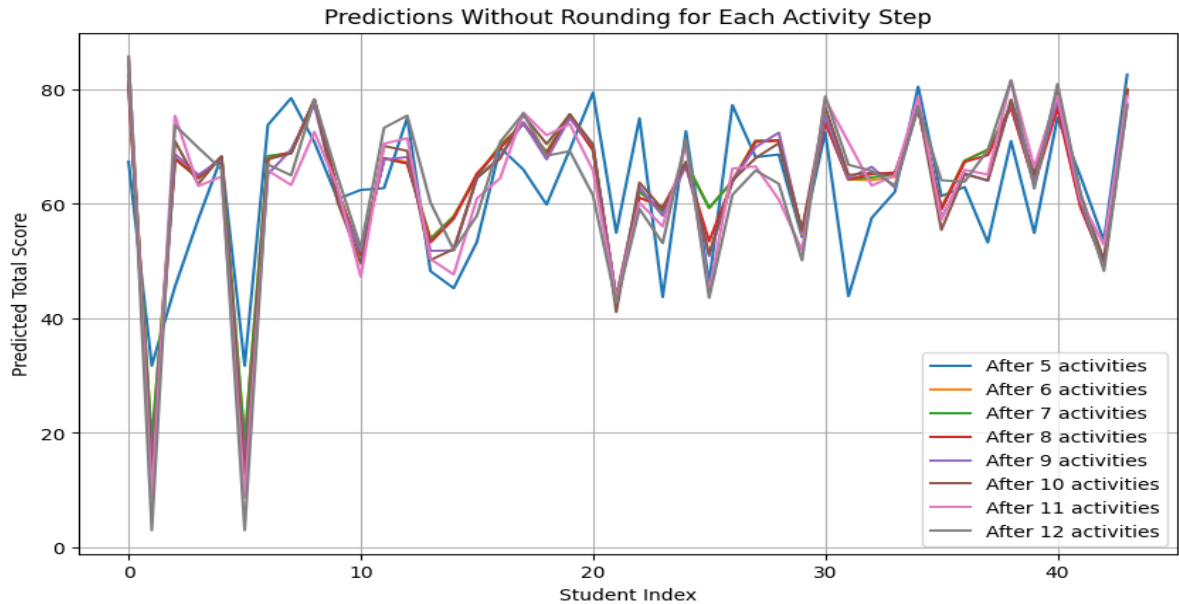
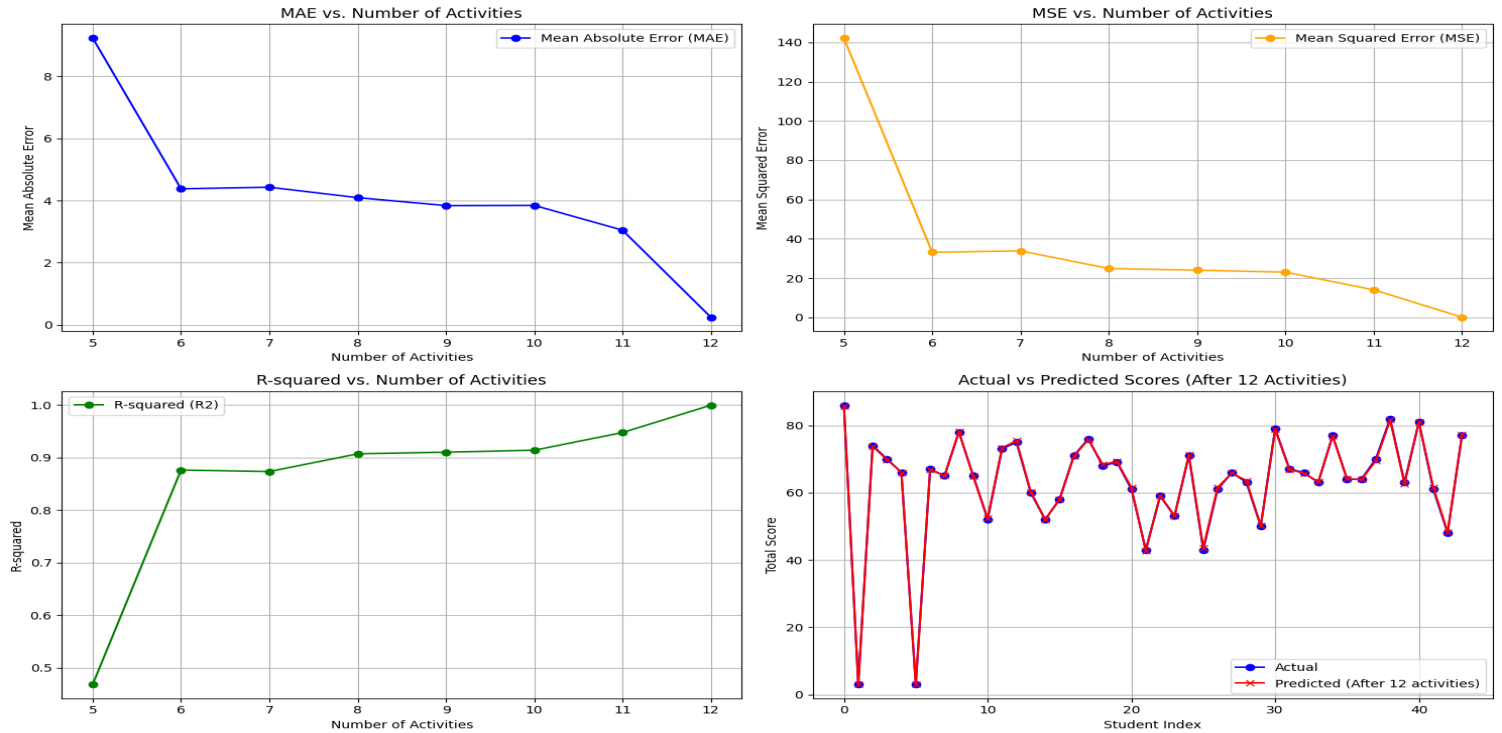## 2. Model Training and Prediction:

- Linear Regression is used to predict the final score based on the scores for individual activities (Q1, Q2, A1, Q3, Q4, Midterm, etc.).
- The model is trained using the morning dataset ($X_{train}$) and tested on the afternoon dataset ($X_{test}$).
- Evaluation metrics like Mean Absolute Error (MAE), Mean Squared Error (MSE), and R-squared ($R^2$) are used to assess the model's performance.

Predictions Without Rounding for Each Activity Step


Predictions With Rounding for Each Activity Step

### 3. Feature Engineering:

- A dynamic approach to building features is used, where the model is trained progressively based on the number of activities considered (starting from 5 activities).
- The model is trained for each step, and predictions are made for the number of activities ranging from 5 to 12.
- Evaluation results are stored, and a comparison of predicted vs. actual scores is shown for each activity step.

MAE vs. Number of Activities

MSE vs. Number of Activities

R-squared vs. Number of Activities

Actual vs Predicted Scores (After 12 Activities)

## 4. Model Evaluation:

- The evaluation metrics (MAE, MSE, R²) are calculated both with and without rounding the predictions.
- The performance of the model is visualized through plots, showing how the evaluation metrics change as more activities are considered in the prediction.
- A summary table is created to show key statistics (like average difference between predicted and actual scores) for each model configuration.

```
[26]

...     Num Activities        MAE         MSE        R2  Avg Diff   Max Diff  \
     0                5   9.234291  142.233560  0.468428   0.644082  28.432987
     1                6   4.379129   33.119109  0.876223  -1.243419   7.662713
     2                7   4.430439   33.863445  0.873442  -1.219118   7.938733
     3                8   4.094567   24.868895  0.907057  -0.745909   7.933490
     4                9   3.838924   24.004314  0.910288  -0.371527   8.207276
     5               10   3.845234   23.026226  0.913944  -0.692801   9.809619
     6               11   3.046522   14.021341  0.947598   0.188846   9.629661
     7               12   0.243736    0.085790  0.999679  -0.050379   0.571478

         Min Diff
     0 -28.726595
     1 -16.241198
     2 -16.317164
     3 -10.459475
     4 -11.155460
     5  -9.603329
     6  -5.662864
     7  -0.591160
```

Detailed Analysis of Model Performance

The table presents various evaluation metrics for different numbers of activities (from 5 to 12). These metrics include **Mean Absolute Error (MAE)**, **Mean Squared Error (MSE)**, **R-squared ($R^2$)**, **Average Difference**, **Max Difference**, and **Min Difference** between the predicted and actual scores.

1. Mean Absolute Error (MAE):

- MAE measures the average magnitude of errors in the predictions, giving us an idea of how far off the predictions are from the actual values.
- As the number of activities increases, the MAE decreases significantly. The lowest MAE occurs when 12 activities are used, with a value of **0.2437**, indicating that the model's predictions are very close to the actual scores.
- The highest MAE occurs when only 5 activities are included, with a value of **9.2343**, which suggests a significant deviation from actual scores.

Key Observation:

- The MAE steadily decreases as more activities are included, indicating an improvement in prediction accuracy.

2. Mean Squared Error (MSE):

- MSE measures the average squared difference between the predicted and actual values, penalizing larger errors more than MAE.
- The MSE follows a similar trend as MAE, decreasing from **142.23** with 5 activities to **0.0858** with 12 activities.
- The large value of MSE with 5 activities is indicative of a poor model fit, whereas the small value at 12 activities shows that the model performs excellently at predicting the final scores.

## Key Observation:

- Like MAE, MSE decreases significantly with the inclusion of more features, indicating better model accuracy.

## 3. R-squared ($R^2$):

- $R^2$ represents the proportion of the variance in the target variable (final scores) explained by the model. A higher $R^2$ means the model is a better fit.
- $R^2$ starts at **0.4684** for 5 activities, indicating that the model explains only about 47% of the variance in the final scores.
- As the number of activities increases, $R^2$ rises steadily, reaching **0.9997** for 12 activities, meaning that the model accounts for almost all of the variance in the final scores.

## Key Observation:

- $R^2$ improves dramatically as more features are added, reflecting the growing predictive power of the model.

## 4. Average Difference:

- The average difference represents the average of the differences between the predicted and actual values.
- The **Average Difference** fluctuates between **-1.24** and **0.19** across the different numbers of activities. This suggests that while the model improves overall, there can still be minor positive and negative discrepancies between predicted and actual values.
- At 12 activities, the average difference becomes almost negligible (**-0.05**), confirming that predictions are very close to the actual values.

## Key Observation:

- The model becomes more consistent as more features are included, with average differences approaching zero.

## 5. Max Difference:

- The **Max Difference** is the largest discrepancy between predicted and actual values.
- This metric shows considerable variation across the different models, with the maximum difference being **28.43** when only 5 activities are used, indicating large prediction errors.
- As more activities are included, the maximum difference decreases, with the smallest maximum difference (**0.57**) occurring at 12 activities, demonstrating that the model becomes increasingly accurate.

## Key Observation:

- With the increase in features, the model's worst predictions improve, with the maximum discrepancy becoming minimal by the time all activities are included.

## 6. Min Difference:

- The **Min Difference** represents the smallest discrepancy between predicted and actual values.
- This metric shows large negative differences when few activities are included, such as -**28.73** at 5 activities, meaning that some predictions are substantially lower than actual values.
- As more activities are incorporated, the minimum difference also improves, reaching -**0.59** at 12 activities, which suggests that the model is generally able to predict well across all students.

## Key Observation:

- Like the maximum difference, the minimum difference becomes less extreme as more features are added to the model.
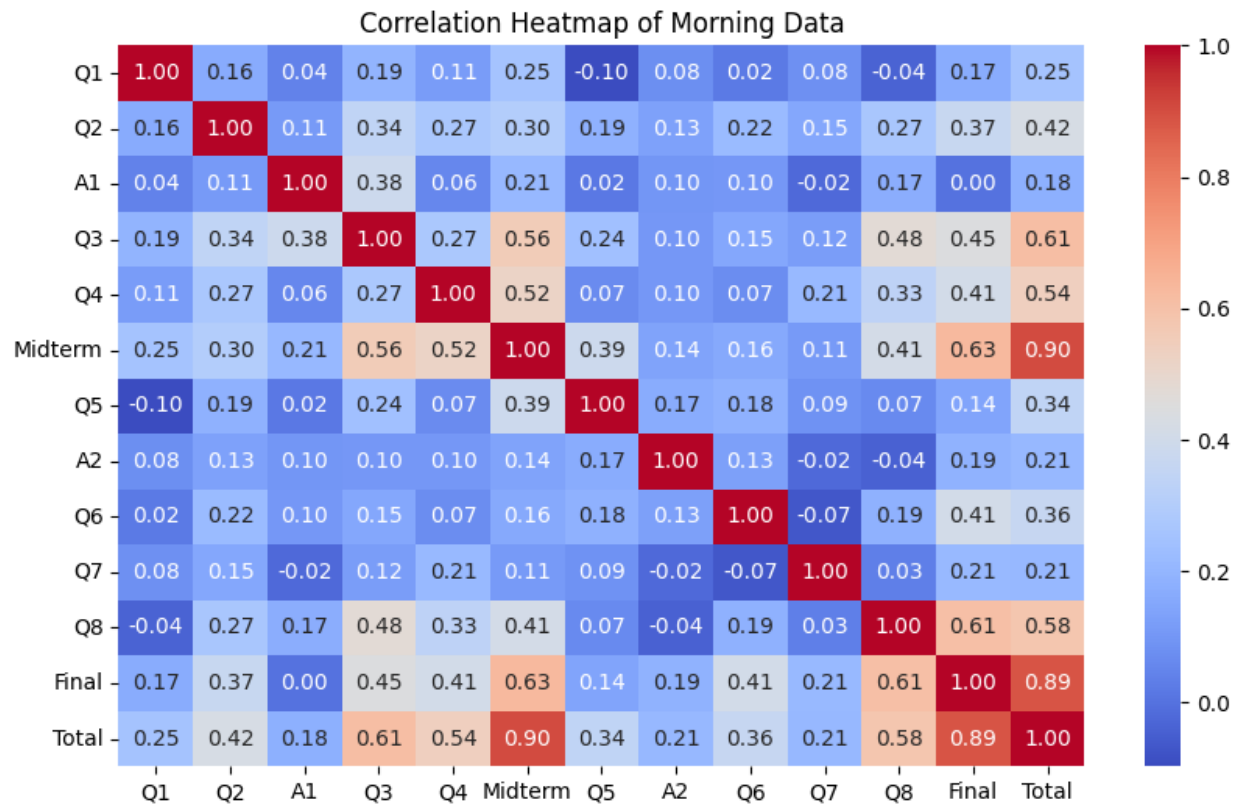
## Summary of Key Findings

- **Increased Number of Activities Improves Performance:** As the number of activities increases, the model's **MAE**, **MSE**, and **Max Diff** decrease, while the $R^2$ value increases. This suggests that including more features from the data allows the model to make more accurate predictions.
- **Best Performance with 12 Activities:** The model achieves its best performance when 12 activities are used, as reflected by the lowest **MAE** (0.2437), **MSE** (0.0858), and the highest $R^2$ (0.9997).
- **Decreasing Prediction Errors:** As more activities are added, prediction errors decrease, with **Avg Diff**, **Max Diff**, and **Min Diff** all showing improvements. The **Max Diff** and **Min Diff** are nearly zero at the 12-activity model, indicating that the model predicts very closely to actual values.
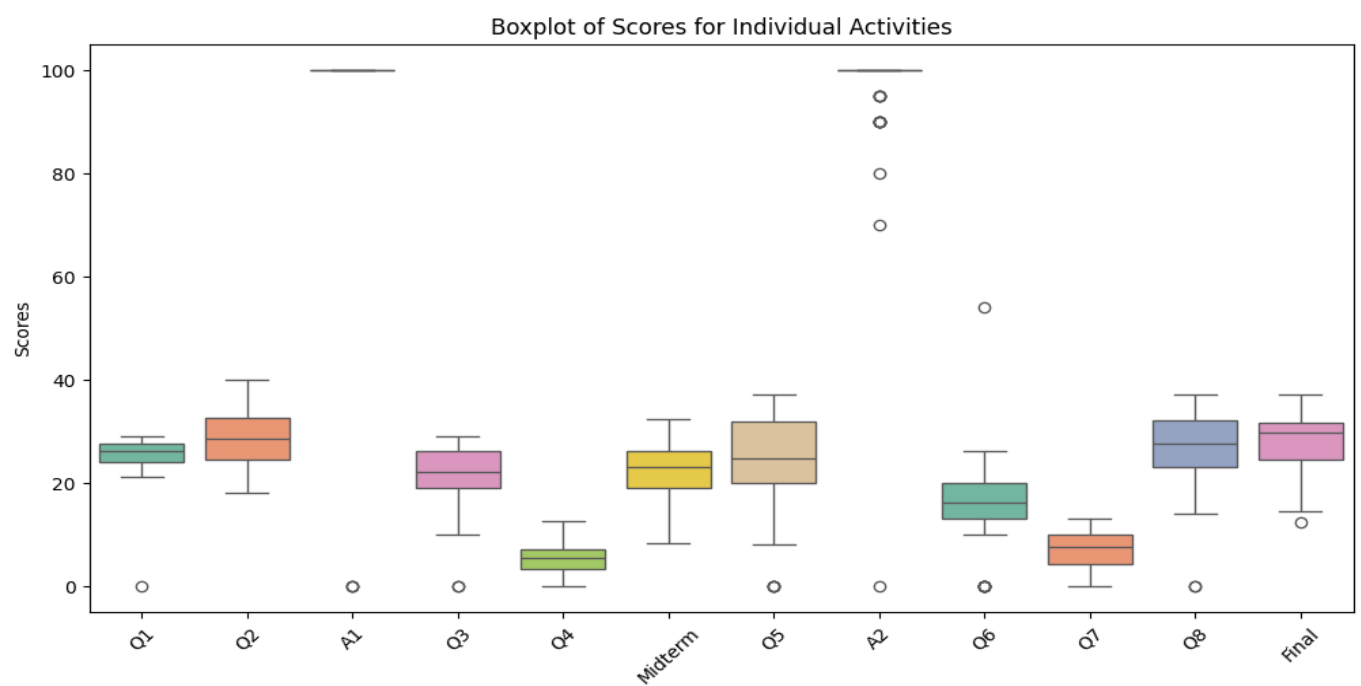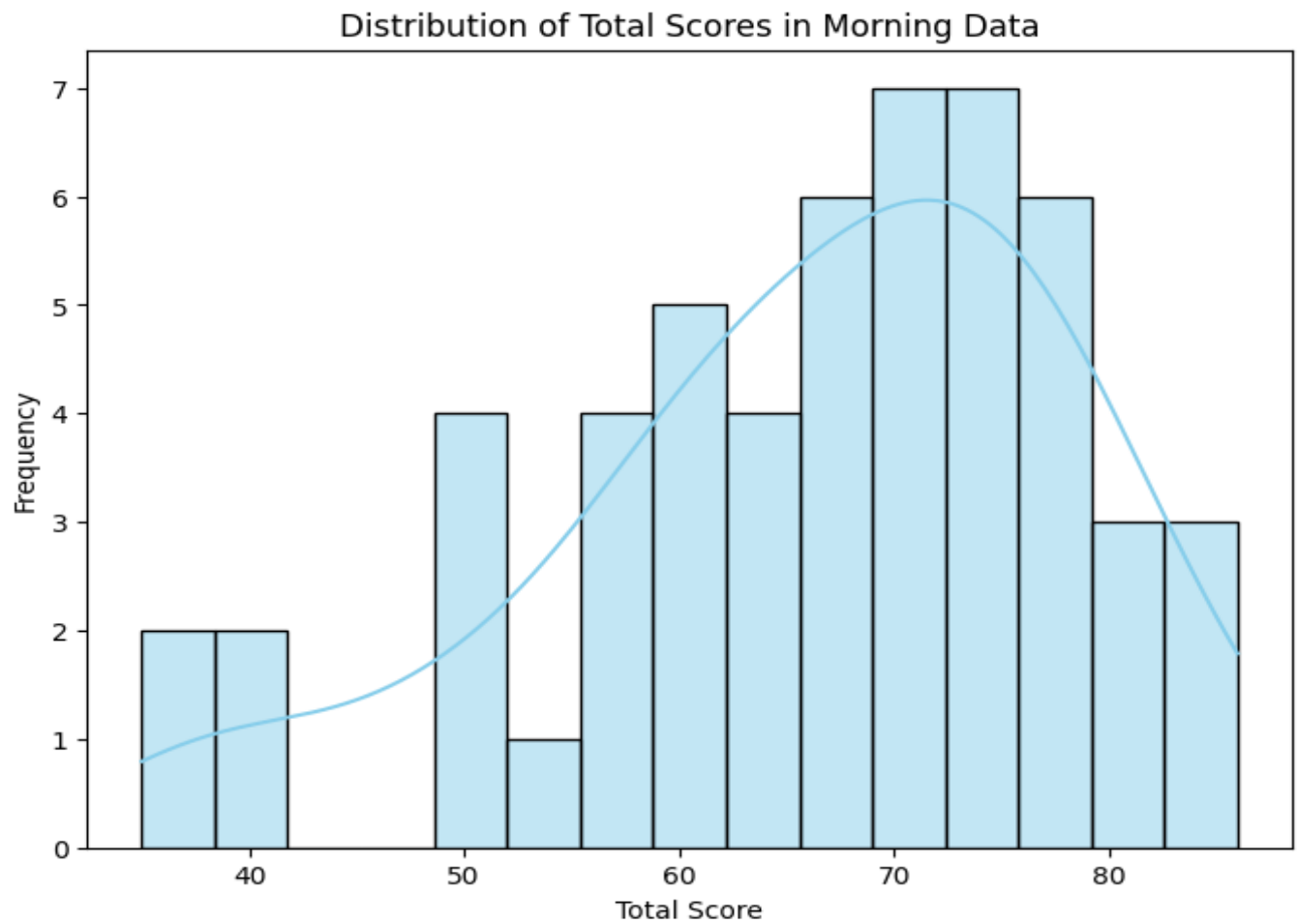
## Conclusion

The analysis demonstrates that the **Linear Regression** model's performance improves significantly with the inclusion of more activities. Starting with 5 activities, the model initially has high errors, but as more activities are included, it becomes increasingly accurate, culminating in an almost perfect fit with 12 activities. The $R^2$ value approaching 1 suggests that the model is highly effective at predicting final scores when all activities are considered.
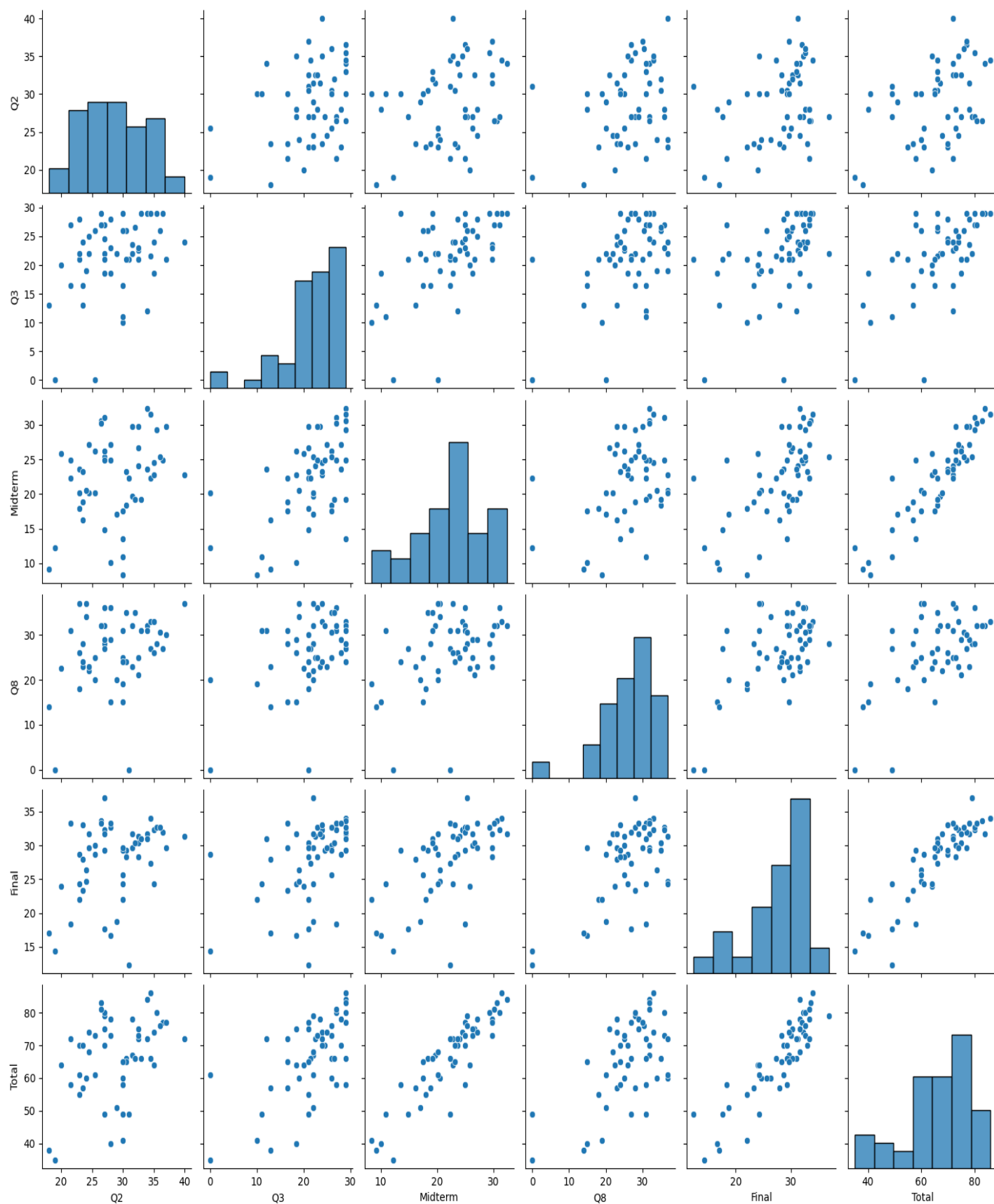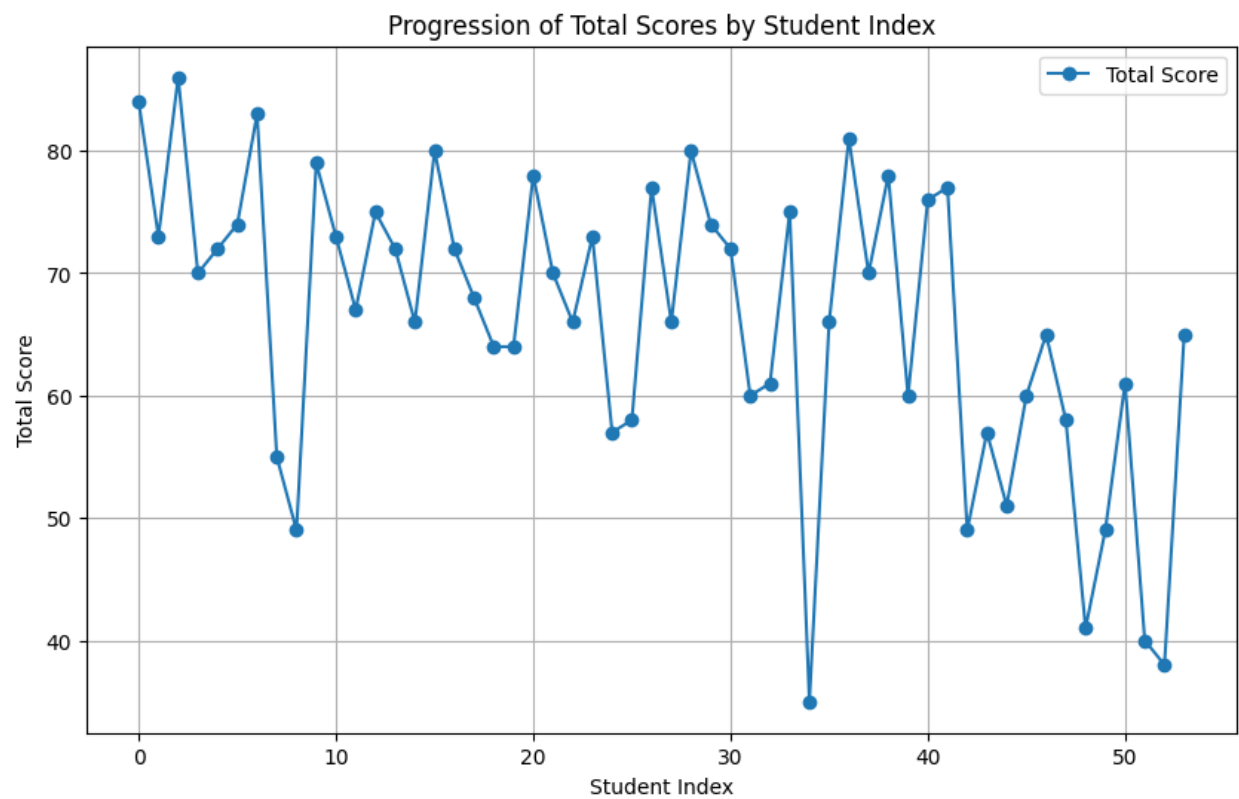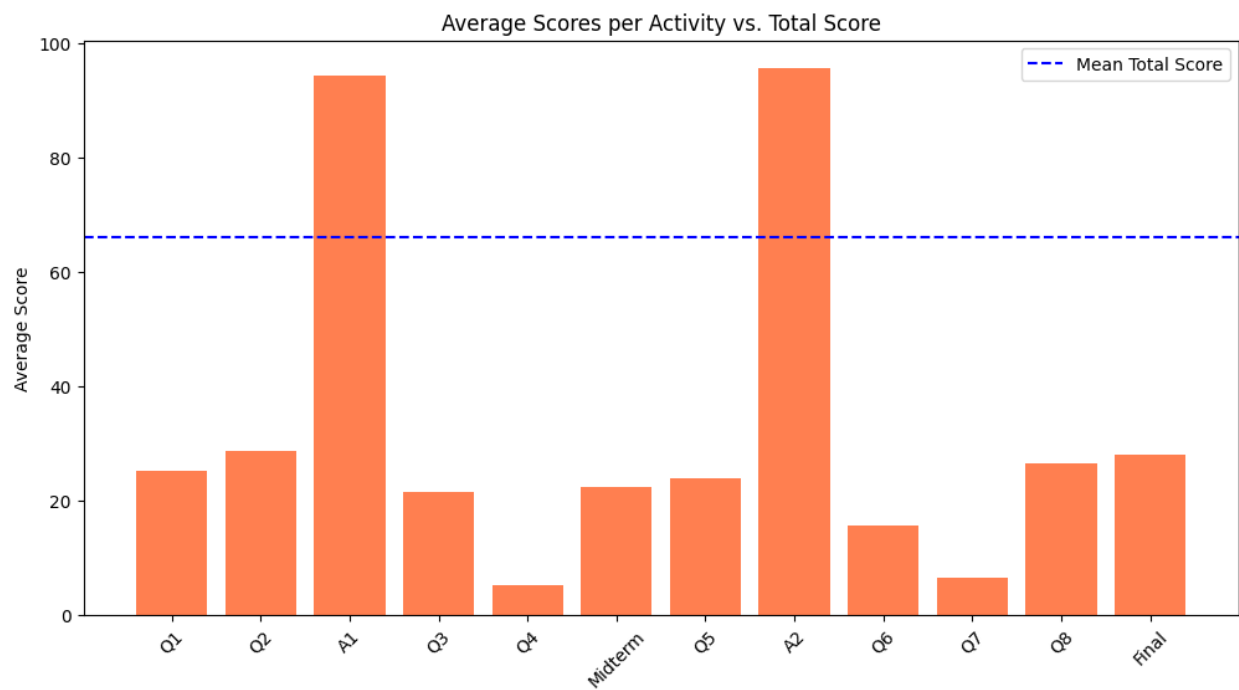
---

## 5. Visualization:

- Several visualizations are created:
    - **Heatmap** for correlation among the features in the morning dataset.
    - **Pairplot** to visualize relationships between selected activities and the total score.
    - **Boxplot** for activity-wise score distribution.
    - **Bar Plot** comparing average activity scores to the total score.
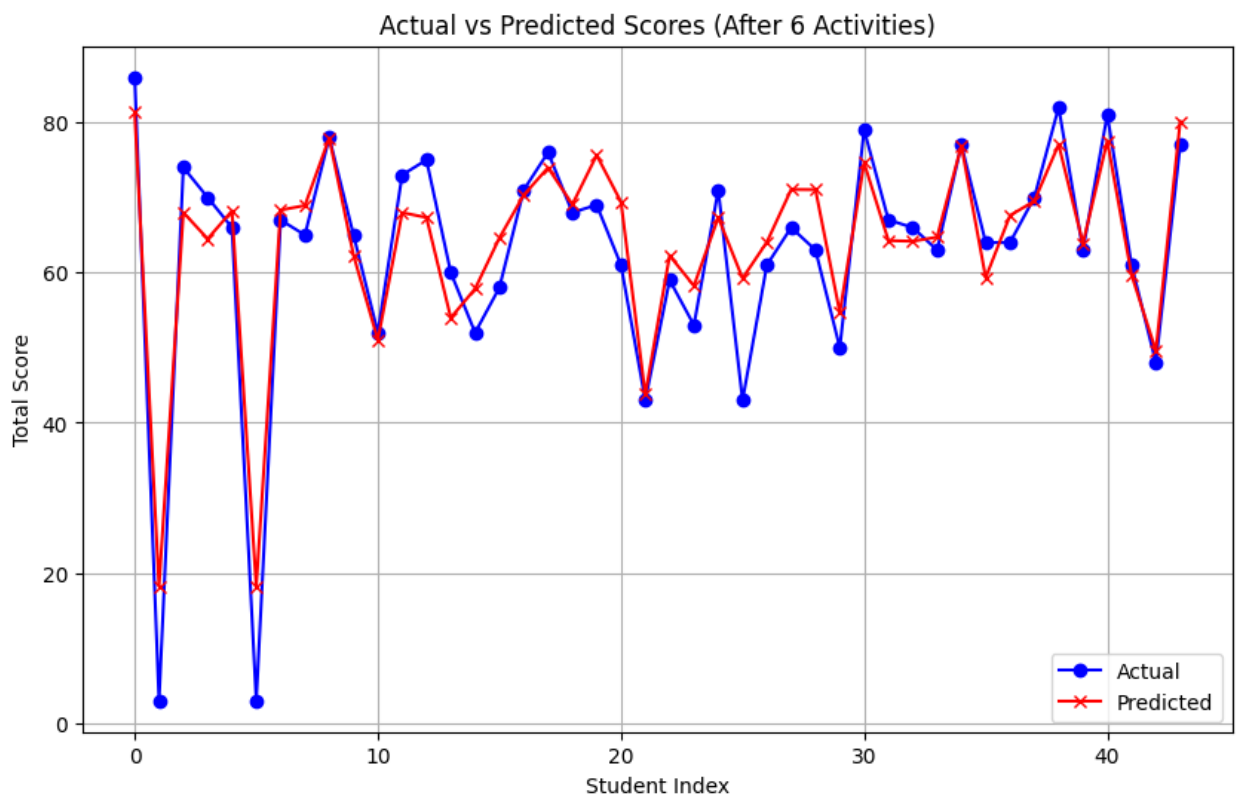    - **Line Plot** showing the progression of total scores by student index.

Correlation Heatmap of Morning Data

Distribution of Total Scores in Morning Data



Boxplot of Scores for Individual Activities

Pairplot of Selected Features and Total

Average Scores per Activity vs. Total Score

Progression of Total Scores by Student Index

## 6. Interactive Prediction:

- An interactive function allows the user to input the number of activities (from 5 to 12) and predict the final score based on the selected activities.
- After entering the number of activities, the model's predicted scores and actual scores are displayed along with a comparison plot.



Actual vs Predicted Scores (After 6 Activities)

**Summary of Performance Metrics:**

For each prediction step (after 5, 6, ..., 12 activities), the following results were generated:

- **Mean Absolute Error (MAE)**
- **Mean Squared Error (MSE)**
- **R-squared (R²)**

Each step demonstrates the improvement (or potential limitations) as more activities are incorporated into the prediction.

**Visualization of Results:**

- **Predictions vs. Actual Scores**: A visual comparison is provided for each activity count (both with and without rounding).
- **Performance Metrics**: The change in MAE, MSE, and $R^2$ is shown across the number of activities used for prediction.

**Summary Table for Prediction Differences:**

- A table summarizes the differences between the predicted and actual scores, showing:
    - MAE, MSE, $R^2$
    - **Average Difference** between predicted and actual scores
    - **Max and Min Differences**

**Conclusion:**

This approach provides an effective way of forecasting student performance by progressively building features and evaluating the model at each stage. Visualization of the prediction steps and comparison of actual vs. predicted scores further enhances the interpretation of the results.