# Analysis Report

This report is structured as follows.

## Contents

## Descriptive Statistics

The table below presents descriptive statistics for various variables across four different exams. The focus is on how these variables change from one exam to the next.

Exam Score: The average exam score starts high at 0.937 in Exam 1 and then shows a slight decrease in Exam 2 to 0.892. It marginally increases in Exams 3 and 4 to 0.903 and 0.913 respectively. This trend suggests a dip in performance in the second exam, followed by a slight recovery in subsequent exams.

Days Studied of 31: The mean days studied begins at 5.565 for Exam 1, decreasing significantly to 4.478 for Exam 2, and then further dropping to around 2.261 and 2.609 for Exams 3 and 4. This indicates a notable reduction in the number of days students spent studying as the exam series progressed.

Total Reviews: There is a fluctuation in the average total reviews, starting at 602.957 for Exam 1, peaking at 702.391 for Exam 2, then decreasing to 488.870 and 526.435 for Exams 3 and 4. This variable does not follow a clear trend across the exams.

Average for Days Studied: This variable also shows variation, starting at 100.174 for Exam 1, slightly decreasing in Exam 2, then increasing in Exams 3 and 4. The increase in later exams suggests a higher intensity or effectiveness of study per day.

Average Total: The mean value shows considerable variation across the exams, with the highest average total observed in Exam 2.

Cards: The variables Card_New, Card_Learning, Card_Relearning, Card_Young, Card_Mature, and Card_Suspended seem to represent different categories or types of study materials or methods. Each shows different patterns of usage across the exams. For example, Card_New usage peaks in Exam 2, while Card_Learning shows a significant increase in the same exam. Card_Relearning and Card_Mature have relatively low averages but show some increase by Exam 4.

| Variable | Exam 1 | | Exam 2 | | Exam 3 | | Exam 4 | |
|---|---|---|---|---|---|---|---|---|
| | Mean | SD | Mean | SD | Mean | SD | Mean | SD |
| Exam_Score | 0.937 | 0.056 | 0.892 | 0.075 | 0.903 | 0.046 | 0.913 | 0.046 |
| Days_studied_of_31 | 5.565 | 5.035 | 4.478 | 5.177 | 2.261 | 3.003 | 2.609 | 4.098 |
| Total_reviews | 602.957 | 598.675 | 702.391 | 1213.630 | 488.870 | 767.828 | 526.435 | 904.097 |
| Average_for_days_studied | 100.174 | 79.920 | 79.652 | 101.077 | 109.130 | 162.886 | 139.304 | 250.217 |
| Average_total | 19.478 | 19.278 | 25.652 | 44.829 | 15.739 | 24.735 | 17.043 | 29.049 |
| Card_New | 408.043 | 223.163 | 610.130 | 344.408 | 467.304 | 286.037 | 573.174 | 334.874 |
| Card_Learning | 14.870 | 26.755 | 47.652 | 177.491 | 4.217 | 9.244 | 6.696 | 10.814 |
| Card_Relearning | 0.652 | 2.288 | 0.000 | 0.000 | 0.261 | 1.251 | 1.304 | 6.041 |
| Card_Young | 235.870 | 227.481 | 174.304 | 267.239 | 189.565 | 267.099 | 218.174 | 320.214 |
| Card_Mature | 10.913 | 16.847 | 21.957 | 51.491 | 0.087 | 0.288 | 7.739 | 22.316 |
| Card_Suspended | 8.000 | 21.030 | 10.696 | 29.007 | 11.609 | 26.534 | 0.783 | 3.541 |

<u>**Descriptive Statistics – No Outliers**</u>
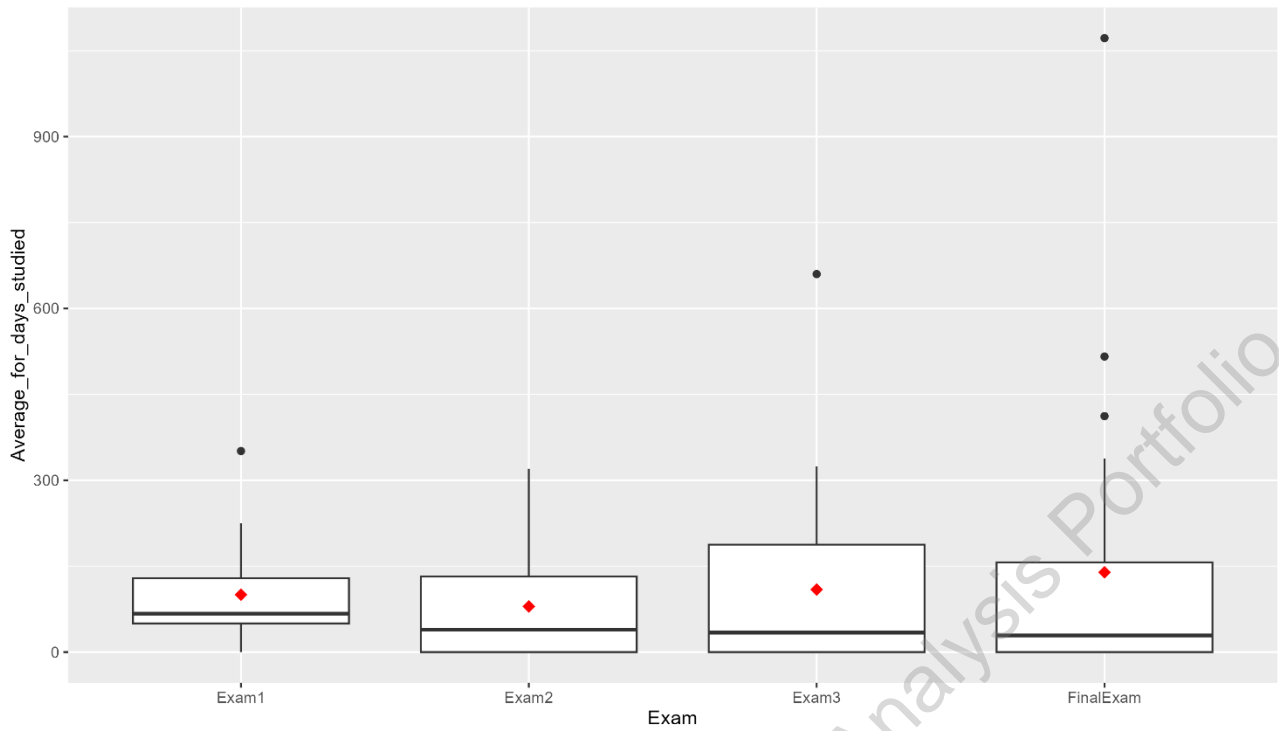
The table below contains descriptive statistics without considering Students 4 and 10.

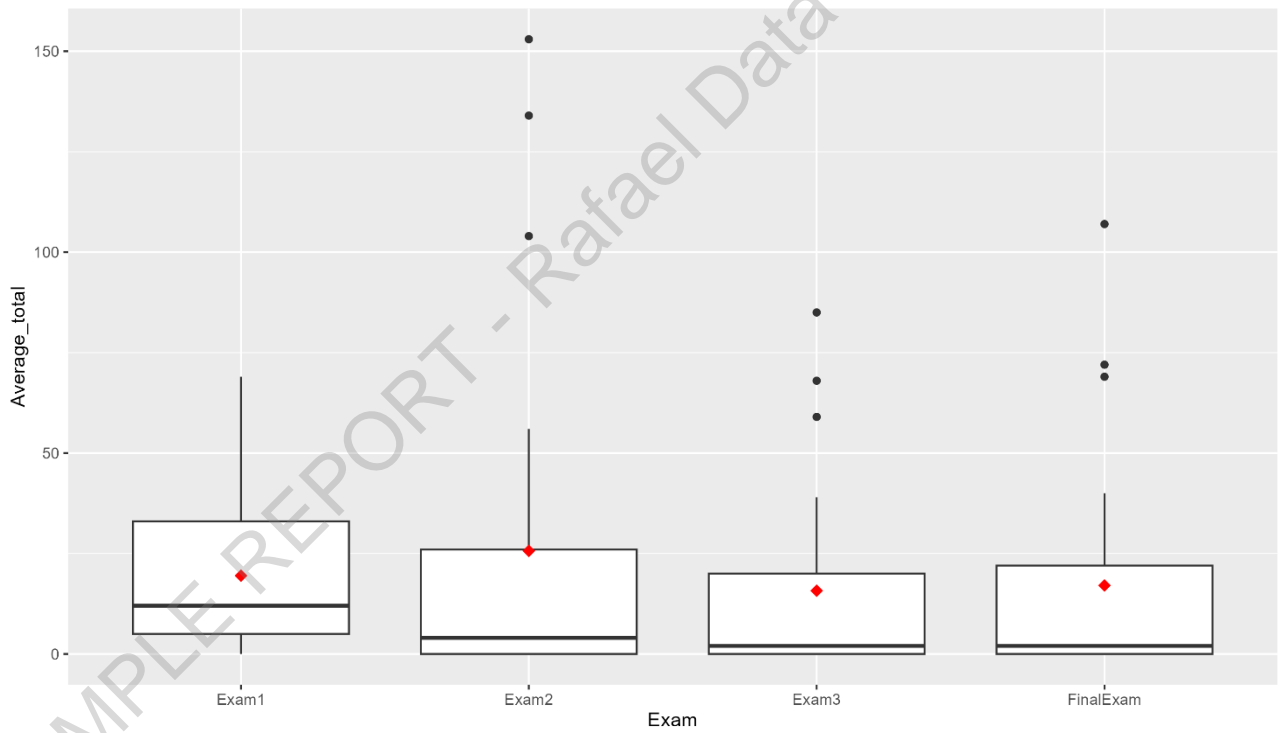| Variable | Exam 1 | | Exam 2 | | Exam 3 | | Exam 4 | |
|---|---|---|---|---|---|---|---|---|
| | **Mean** | **SD** | **Mean** | **SD** | **Mean** | **SD** | **Mean** | **SD** |
| Exam_Score | 0.933 | 0.057 | 0.897 | 0.066 | 0.904 | 0.047 | 0.916 | 0.038 |
| Days_studied_of_31 | 5.333 | 4.564 | 4.762 | 5.319 | 2.143 | 2.920 | 2.238 | 3.548 |
| Total_reviews | 558.762 | 506.774 | 761.000 | 1256.167 | 448.571 | 738.738 | 469.762 | 857.071 |
| Average_for_days_studied | 103.381 | 80.371 | 84.476 | 104.274 | 107.095 | 165.626 | 144.333 | 260.415 |
| Average_total | 18.048 | 16.280 | 27.810 | 46.405 | 14.429 | 23.775 | 15.238 | 27.555 |
| Card_New | 397.238 | 229.000 | 589.619 | 353.958 | 479.714 | 277.097 | 589.286 | 322.332 |
| Card_Learning | 16.190 | 27.681 | 50.714 | 185.785 | 4.619 | 9.594 | 7.333 | 11.124 |
| Card_Relearning | 0.714 | 2.390 | 0.000 | 0.000 | 0.286 | 1.309 | 1.429 | 6.321 |
| Card_Young | 227.333 | 213.184 | 189.190 | 275.409 | 176.667 | 256.719 | 204.762 | 312.619 |
| Card_Mature | 10.333 | 16.710 | 24.048 | 53.513 | 0.048 | 0.218 | 4.190 | 13.902 |
| Card_Suspended | 8.619 | 21.946 | 11.714 | 30.216 | 11.619 | 27.591 | 0.857 | 3.705 |

<u>**Boxplots**</u>

The boxplots on this section provide a visualization of the data distribution of different metrics on different exams. Some variables contain several outliers, which in some cases make the box barely visible in the graph.
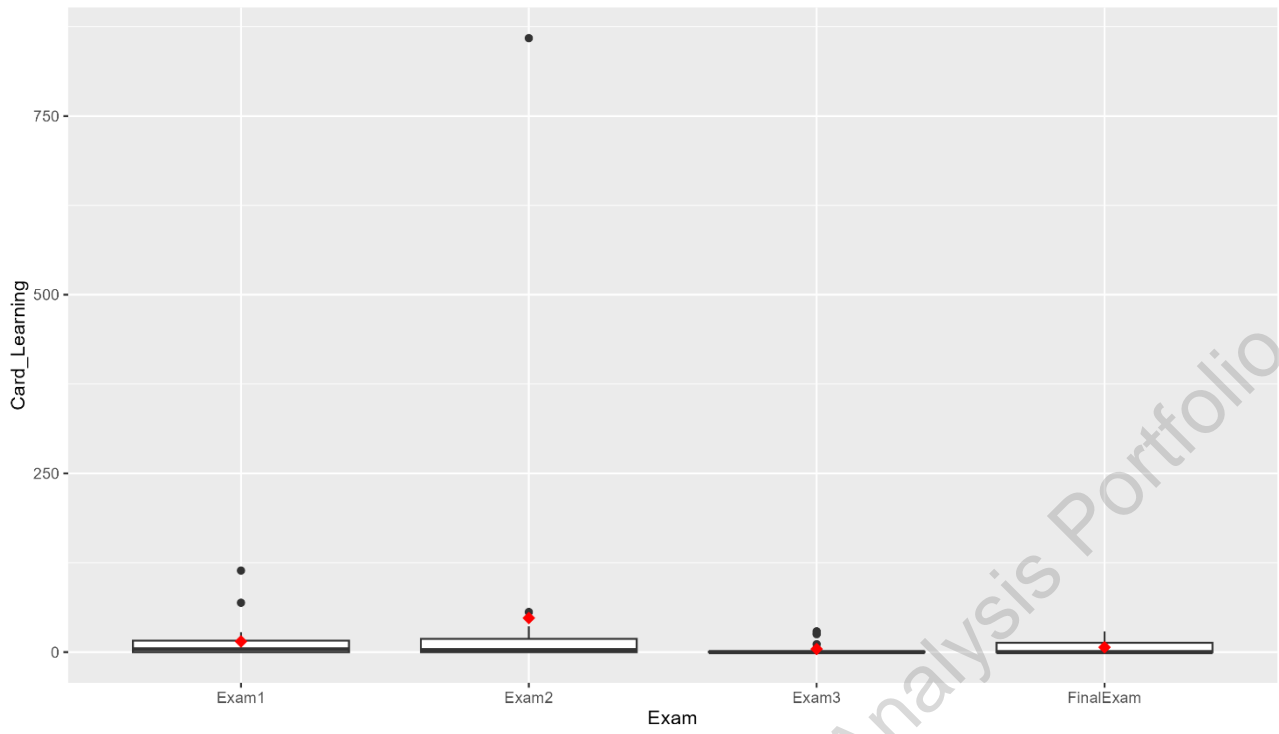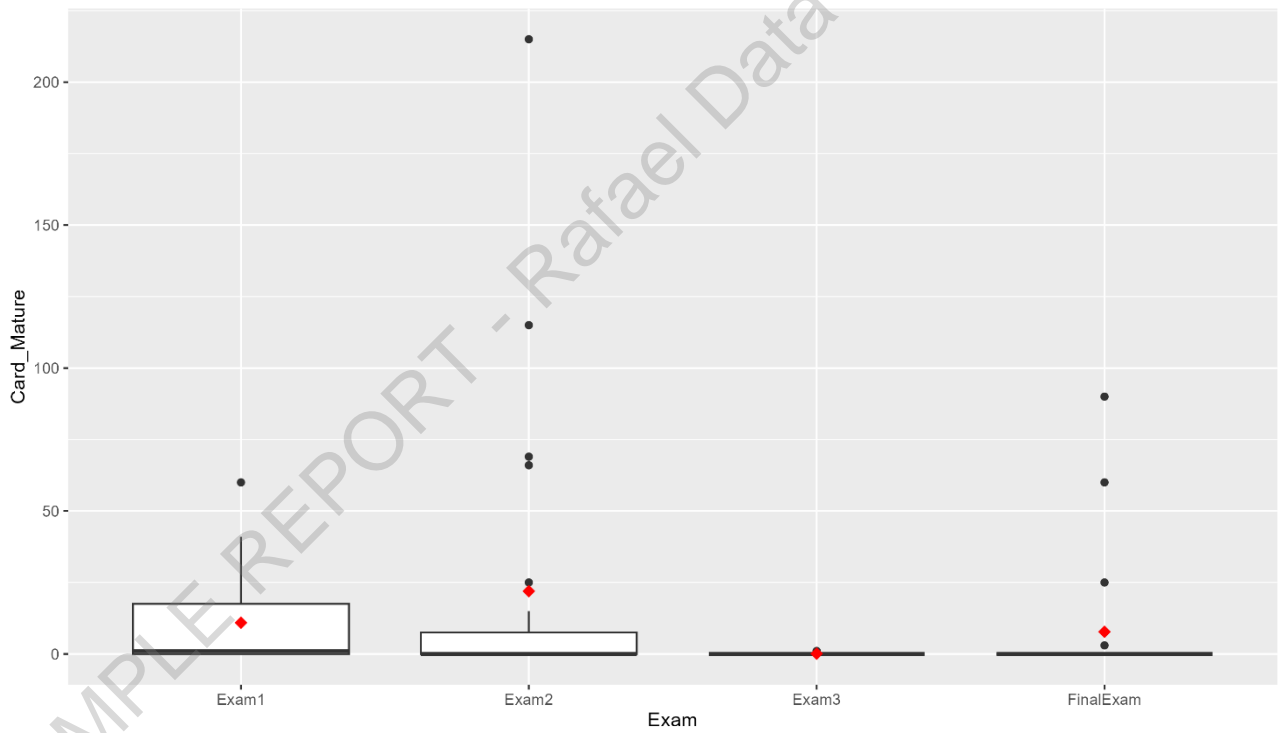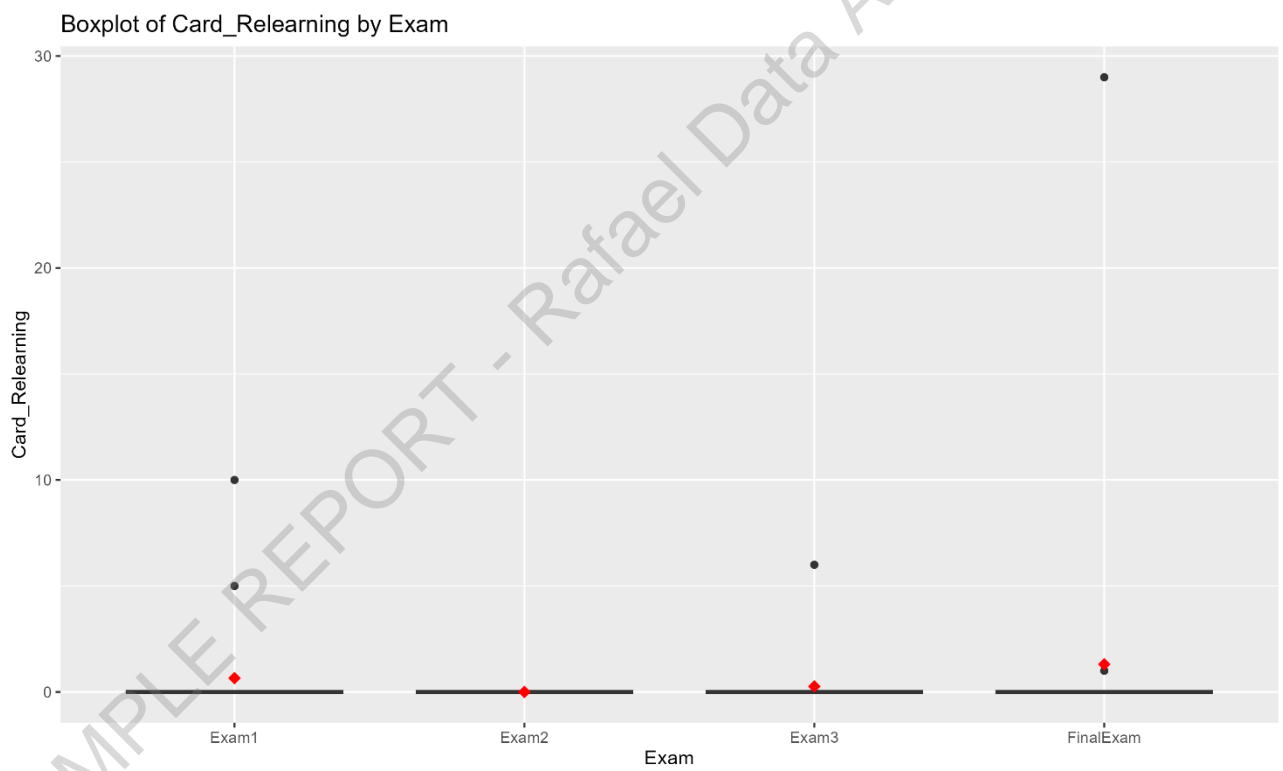
## Boxplot of Average_for_days_studied by Exam
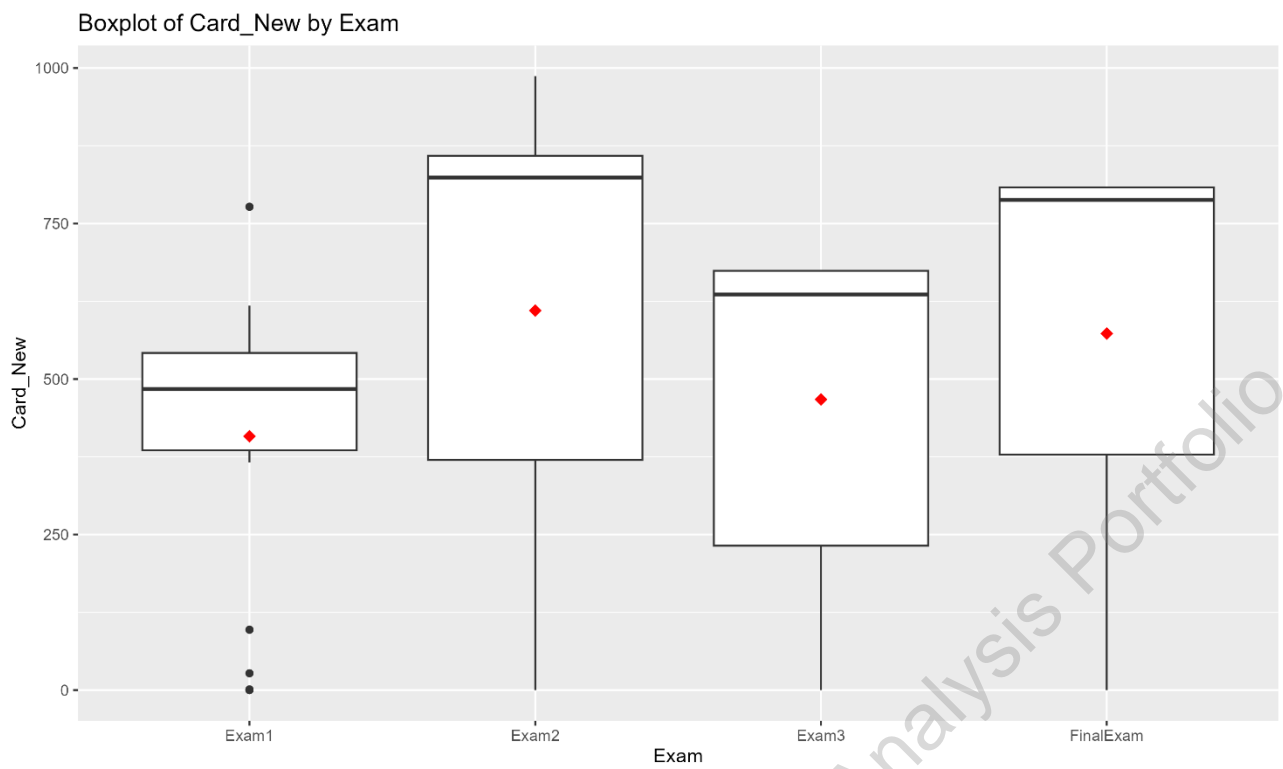


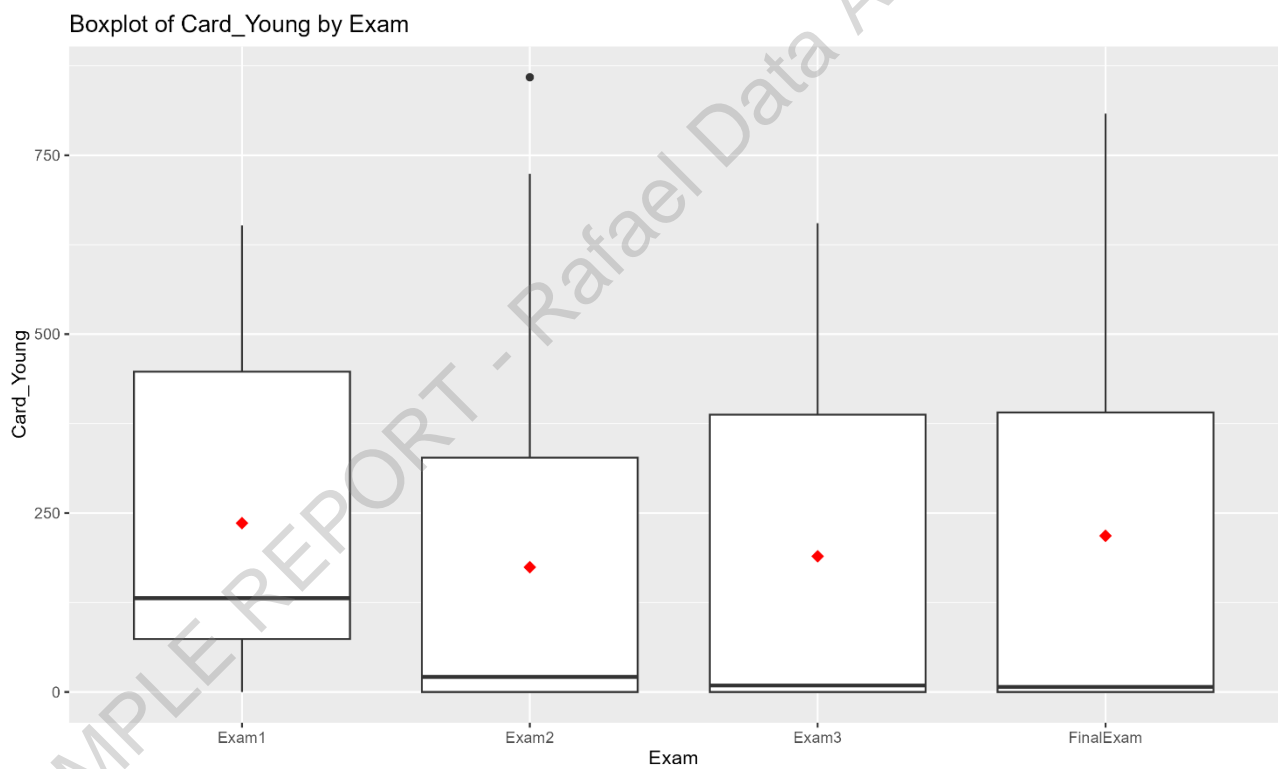## Boxplot of Average_total by Exam

Boxplot of Card_Learning by Exam



Boxplot of Card_Mature by Exam

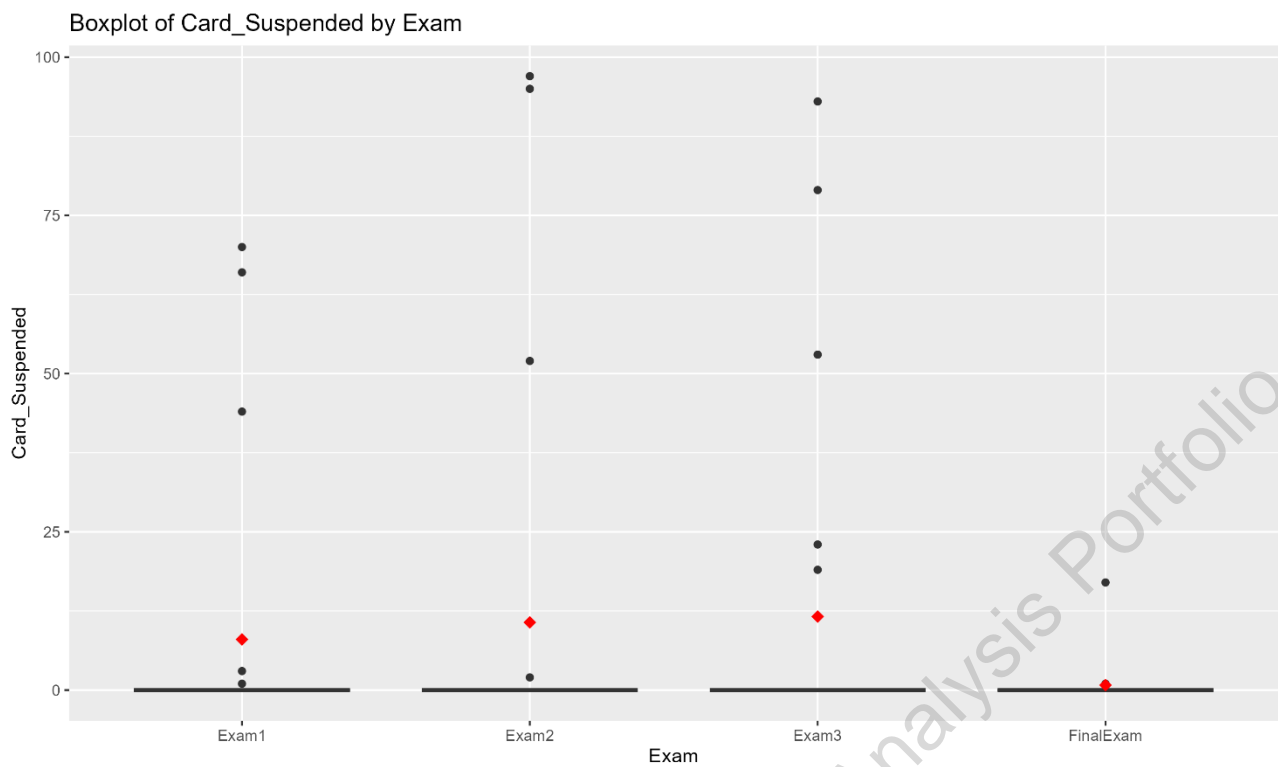## Boxplot of Card_New by Exam



## Boxplot of Card_Relearning by Exam

Boxplot of Card_Suspended by Exam
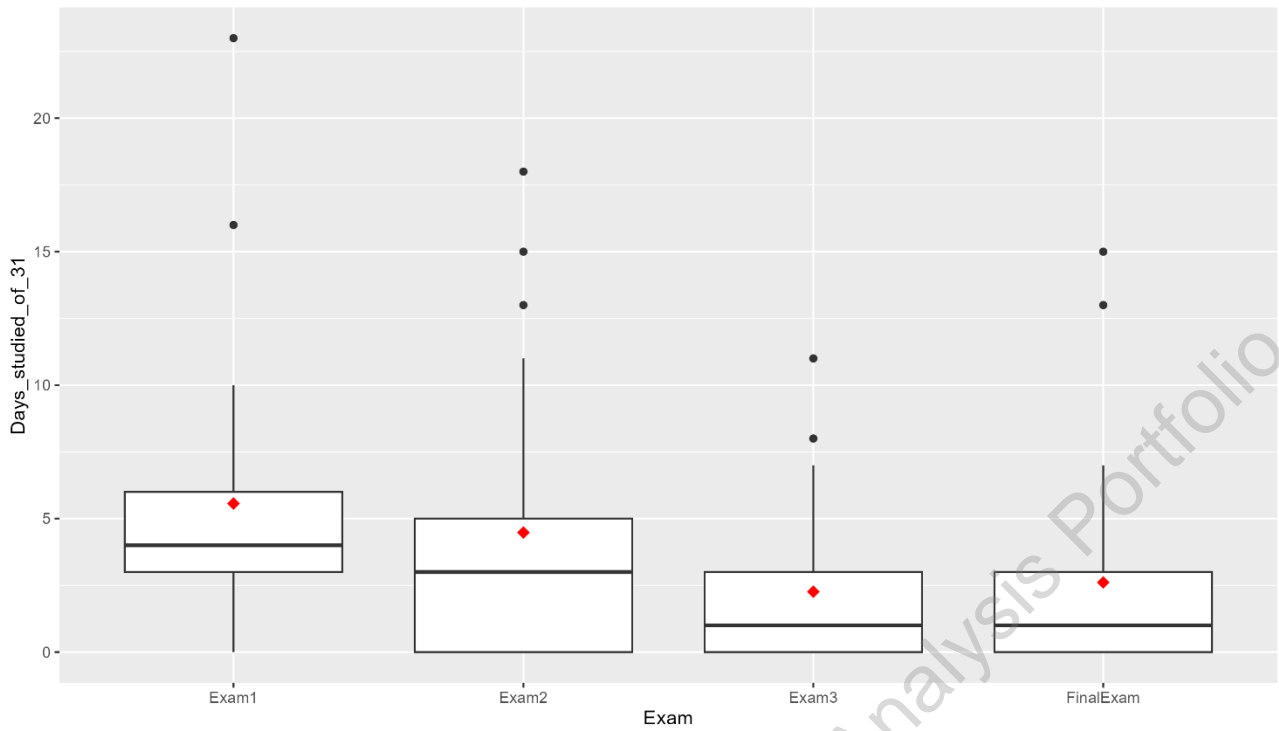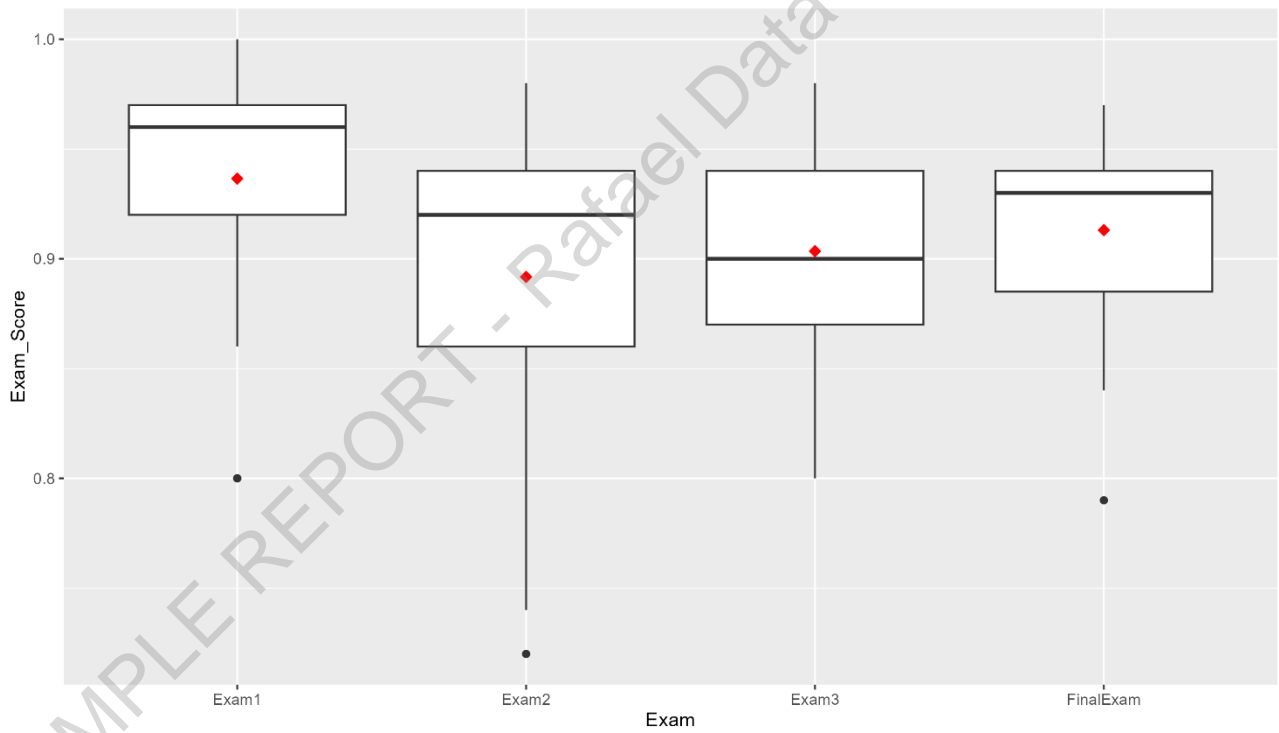

Boxplot of Card_Young by Exam

Boxplot of Days_studied_of_31 by Exam
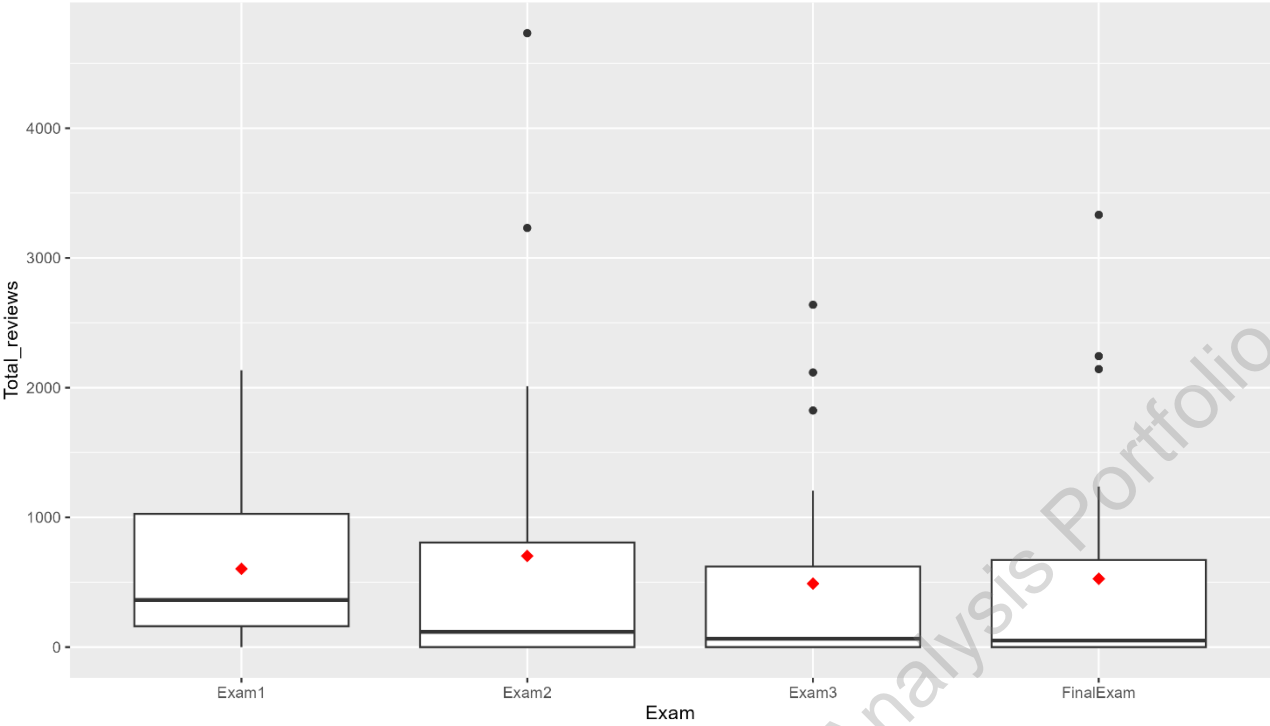


Boxplot of Exam_Score by Exam

Boxplot of Total_reviews by Exam

<u>**Correlation Analysis**</u>

The correlation table reveals several interesting relationships between "Exam_Score" and other variables, though it's important to note that these correlations, while statistically significant in some cases, are generally weak to moderate in strength. "Exam_Score" shows a positive correlation with "Days_studied_of_31" ($r = 0.173$, $p < 0.1$), indicating that students who spend more days studying tend to have slightly higher exam scores, although the correlation is weak. A similar relationship is observed with "Total_reviews" ($r = 0.185$, $p < 0.1$), suggesting that students who review more also tend to score slightly higher on exams.

However, for "Average_for_days_studied", the correlation with "Exam_Score" is positive ($r = 0.132$) but not statistically significant, implying that the average score per day studied does not have a clear relationship with overall exam scores. On the other hand, "Average_total" exhibits a positive correlation ($r = 0.185$, $p < 0.1$) with "Exam_Score", again indicating a weak association where higher averages are slightly related to better exam performance.

Interestingly, "Exam_Score" has negative correlations with "Card_New" ($r = -0.229$, $p < 0.05$) and "Card_Learning" ($r = -0.252$, $p < 0.05$), both significant at the 0.05 level. These correlations suggest that students who are more engaged with new cards or learning cards tend to score lower, although these relationships are not strong. Conversely, "Card_Young" shows a moderately positive correlation ($r = 0.246$, $p < 0.05$) with "Exam_Score", indicating that engagement with this category of cards is associated with higher exam scores.

For "Card_Mature" and "Card_Suspended", the correlations with "Exam_Score" are positive but relatively weak ($r = 0.090$ and $r = 0.221$, $p < 0.05$, respectively). These correlations suggest a slight association between these card types and exam scores, with higher engagement correlating with marginally higher exam scores.

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **Exam_Score (1)** | 1.000 | 0.173* | 0.185* | 0.132 | 0.185* | -0.229** | -0.252** | 0.080 | 0.246** | 0.090 | 0.221** |
| **Days_studied_of_31 (2)** | 0.173* | 1.000 | 0.800*** | 0.249** | 0.770*** | -0.596*** | -0.058 | 0.257** | 0.647*** | 0.578*** | 0.461*** |
| **Total_reviews (3)** | 0.185* | 0.800*** | 1.000 | 0.615*** | 0.973*** | -0.757*** | -0.067 | 0.313*** | 0.859*** | 0.393*** | 0.473*** |
| **Average_for_days_studied (4)** | 0.132 | 0.249** | 0.615*** | 1.000 | 0.597*** | -0.661*** | -0.063 | 0.099 | 0.750*** | 0.068 | 0.247** |
| **Average_total (5)** | 0.185* | 0.770*** | 0.973*** | 0.597*** | 1.000 | -0.748*** | -0.064 | 0.287*** | 0.849*** | 0.379*** | 0.537*** |
| **Card_New (6)** | -0.229** | -0.596*** | -0.757*** | -0.661*** | -0.748*** | 1.000 | -0.159 | -0.224** | -0.889*** | -0.305*** | -0.526*** |
| **Card_Learning (7)** | -0.252** | -0.058 | -0.067 | -0.063 | -0.064 | -0.159 | 1.000 | -0.031 | -0.090 | -0.045 | -0.036 |
| **Card_Relearning (8)** | 0.080 | 0.257** | 0.313*** | 0.099 | 0.287*** | -0.224** | -0.031 | 1.000 | 0.230** | 0.050 | 0.081 |
| **Card_Young (9)** | 0.246** | 0.647*** | 0.859*** | 0.750*** | 0.849*** | -0.889*** | -0.090 | 0.230** | 1.000 | 0.293*** | 0.466*** |
| **Card_Mature (10)** | 0.090 | 0.578*** | 0.393*** | 0.068 | 0.379*** | -0.305*** | -0.045 | 0.050 | 0.293*** | 1.000 | 0.237** |
| **Card_Suspended (11)** | 0.221** | 0.461*** | 0.473*** | 0.247** | 0.537*** | -0.526*** | -0.036 | 0.081 | 0.466*** | 0.237** | 1.000 |

*: $p < 0.1$ **: $p < 0.05$ *** $p < 0.01$

## Correlation Analysis – No Outliers

The matrix below presents correlation coefficients without Students 4 and 10. Now 'Average for Days Studied' is significantly positively correlated with Exam Scores ($r = 0.192$, $p < 0.10$). The magnitude of the correlations between exam scores and both total reviews and days studied have increased without the outliers.

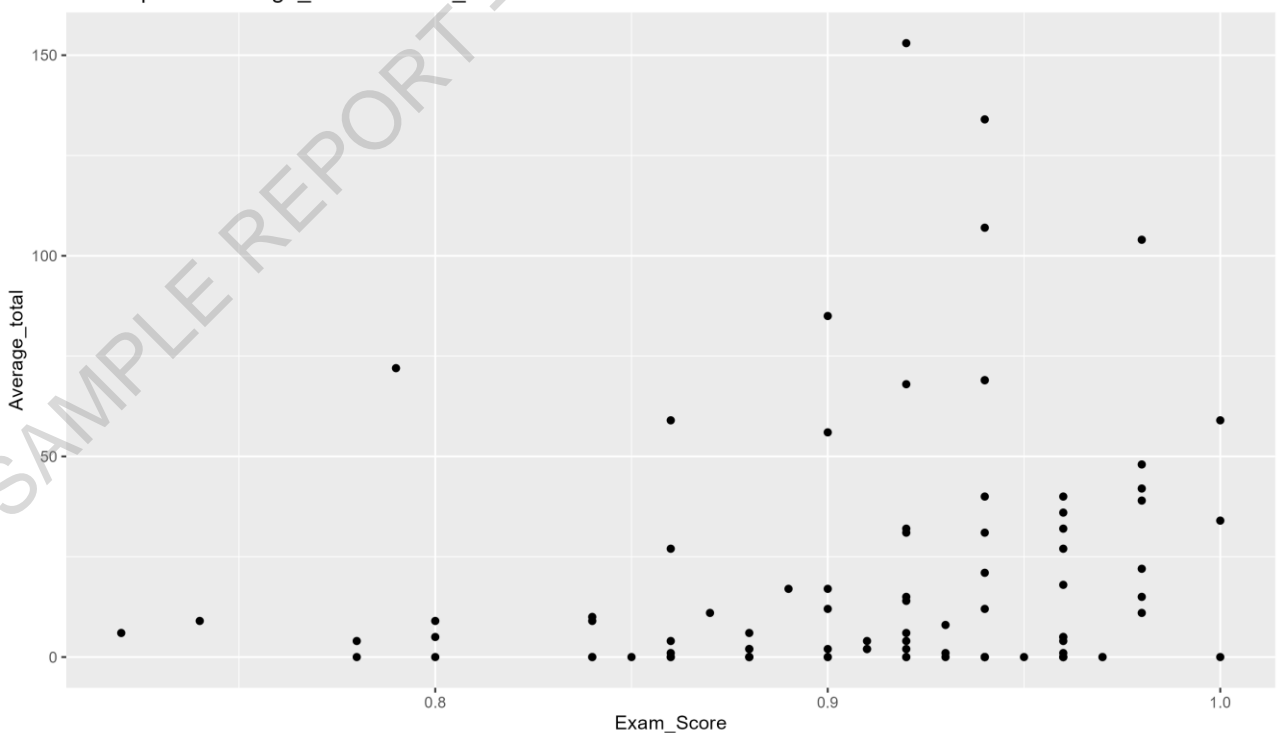|  | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **Exam_Score (1)** | 1.000 | 0.298*** | 0.293*** | 0.192* | 0.285*** | -0.333*** | -0.278** | 0.086 | 0.373*** | 0.174 | 0.255** |
| **Days_studied_of_31 (2)** | 0.298*** | 1.000 | 0.778*** | 0.226** | 0.747*** | -0.586*** | -0.058 | 0.286*** | 0.601*** | 0.555*** | 0.503*** |
| **Total_reviews (3)** | 0.293*** | 0.778*** | 1.000 | 0.615*** | 0.970*** | -0.746*** | -0.065 | 0.338*** | 0.841*** | 0.353*** | 0.495*** |
| **Average_for_days_studied (4)** | 0.192* | 0.226** | 0.615*** | 1.000 | 0.592*** | -0.655*** | -0.066 | 0.098 | 0.760*** | 0.051 | 0.236** |
| **Average_total (5)** | 0.285*** | 0.747*** | 0.970*** | 0.592*** | 1.000 | -0.736*** | -0.062 | 0.306*** | 0.832*** | 0.340*** | 0.559*** |
| **Card_New (6)** | -0.333*** | -0.586*** | -0.746*** | -0.655*** | -0.736*** | 1.000 | -0.171 | -0.236** | -0.889*** | -0.271** | -0.537*** |
| **Card_Learning (7)** | -0.278** | -0.058 | -0.065 | -0.066 | -0.062 | -0.171 | 1.000 | -0.034 | -0.091 | -0.042 | -0.039 |
| **Card_Relearning (8)** | 0.086 | 0.286*** | 0.338*** | 0.098 | 0.306*** | -0.236** | -0.034 | 1.000 | 0.250** | 0.055 | 0.078 |
| **Card_Young (9)** | 0.373*** | 0.601*** | 0.841*** | 0.760*** | 0.832*** | -0.889*** | -0.091 | 0.250** | 1.000 | 0.241** | 0.488*** |
| **Card_Mature (10)** | 0.174 | 0.555*** | 0.353*** | 0.051 | 0.340*** | -0.271** | -0.042 | 0.055 | 0.241** | 1.000 | 0.259** |
| **Card_Suspended (11)** | 0.255** | 0.503*** | 0.495*** | 0.236** | 0.559*** | -0.537*** | -0.039 | 0.078 | 0.488*** | 0.259** | 1.000 |

*: $p < 0.1$ **: $p < 0.05$ *** $p < 0.01$

## Scatterplots

This section shows bivariate scatterplots between Exam Scores and other factors, which illustrate how they are correlated. In some cases such as the two figures below, it is noticeable that higher values of 'Exam_Score' tend to be associated with higher values of 'Days_Studied' or 'Average_Total'.
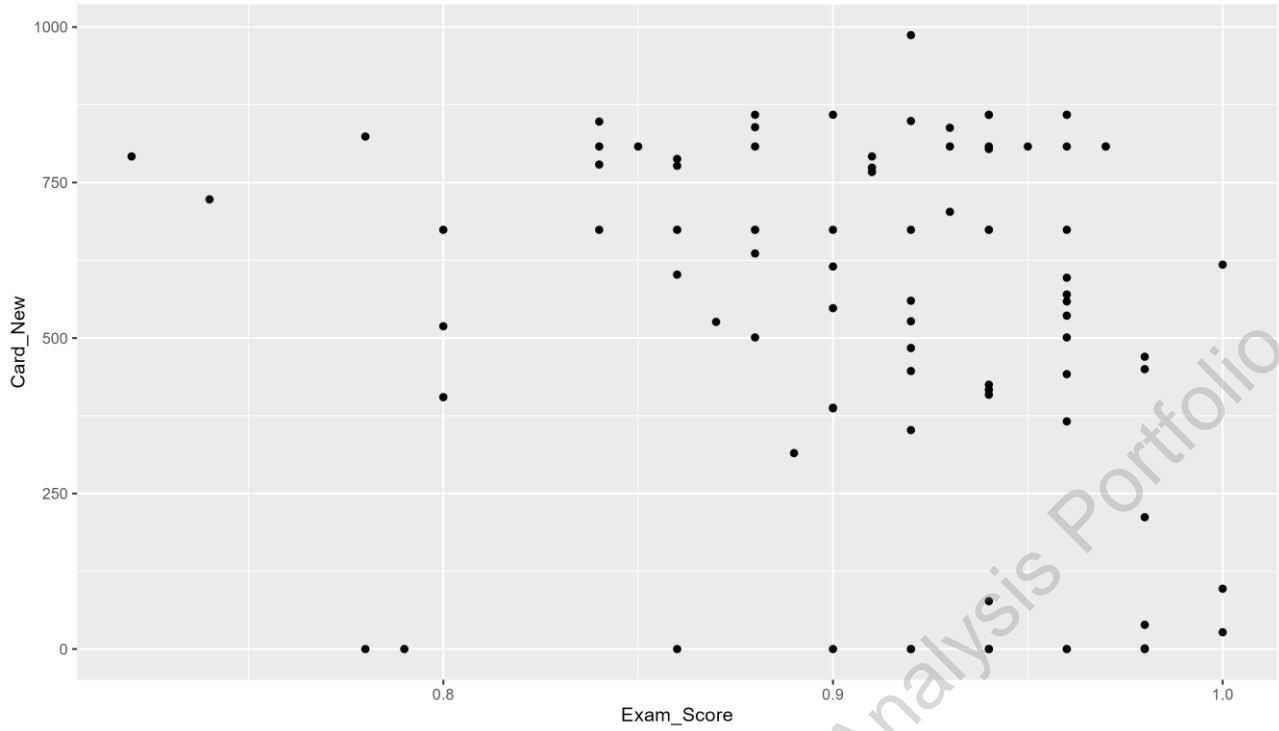


Scatterplot of Average_for_days_studied vs Exam_Score



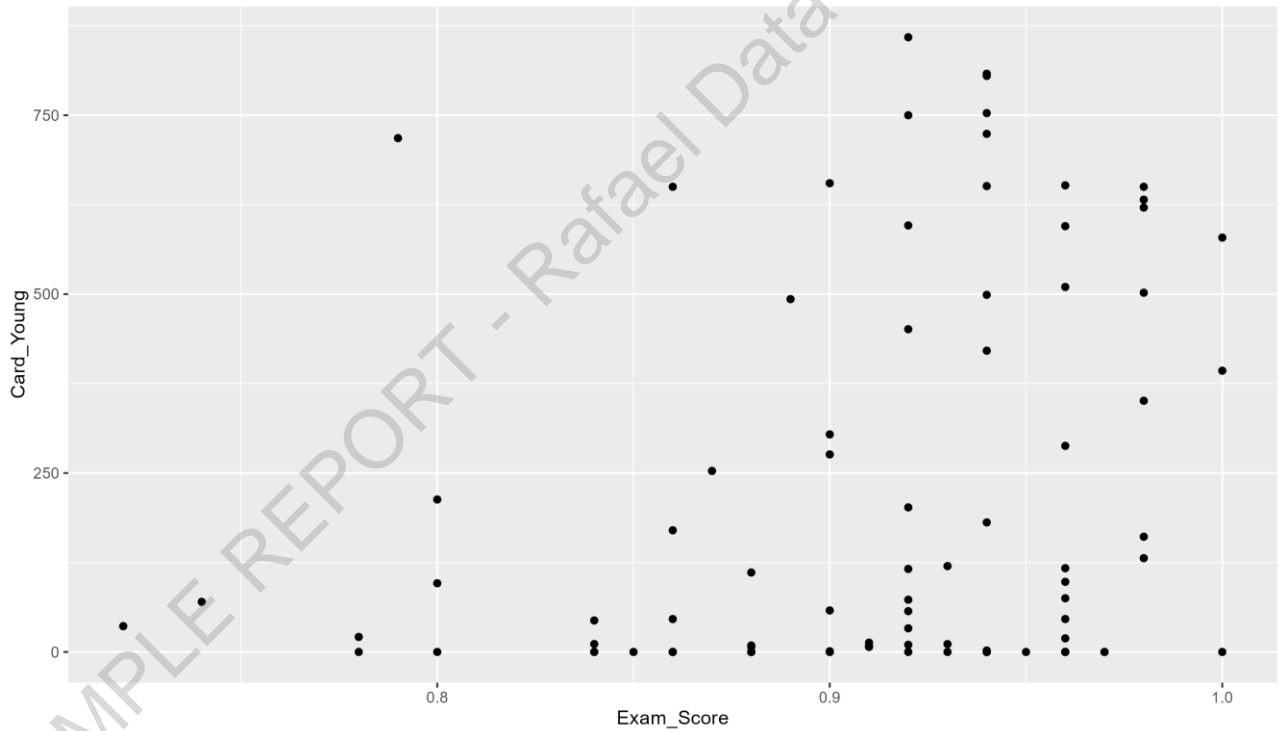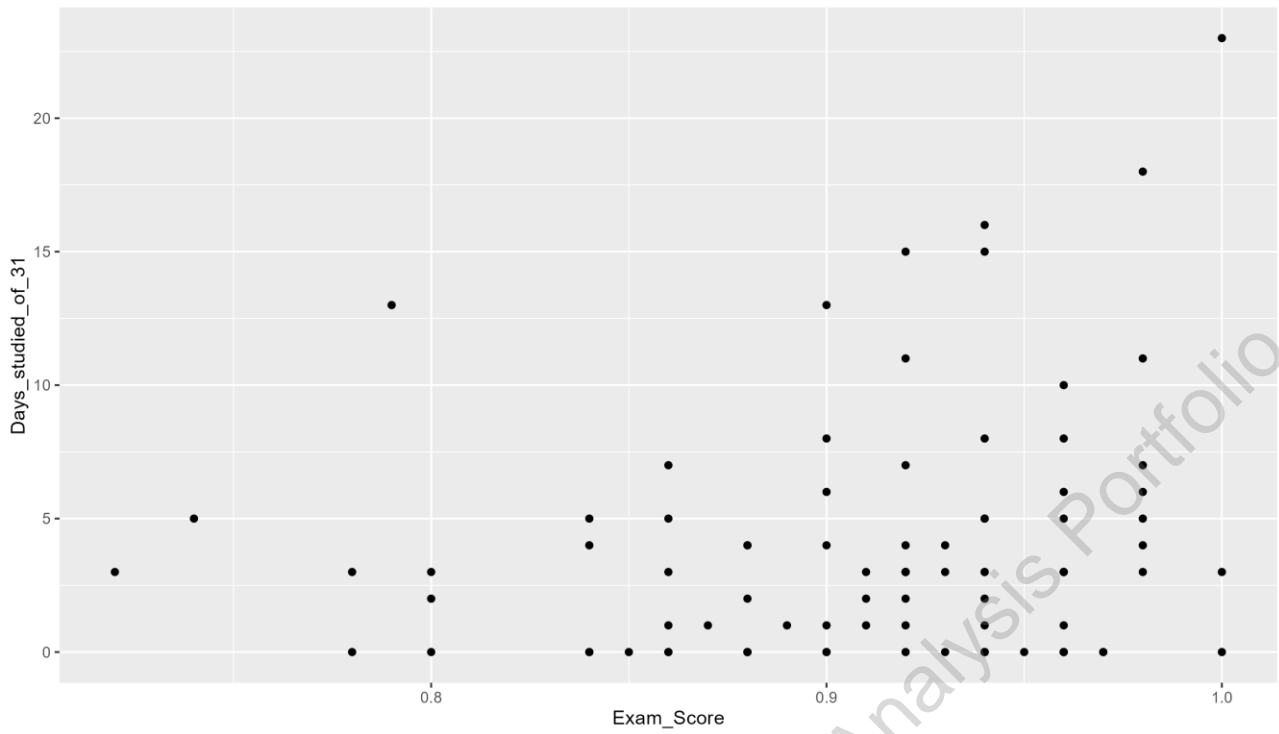Scatterplot of Average_total vs Exam_Score

## Scatterplot of Card_New vs Exam_Score



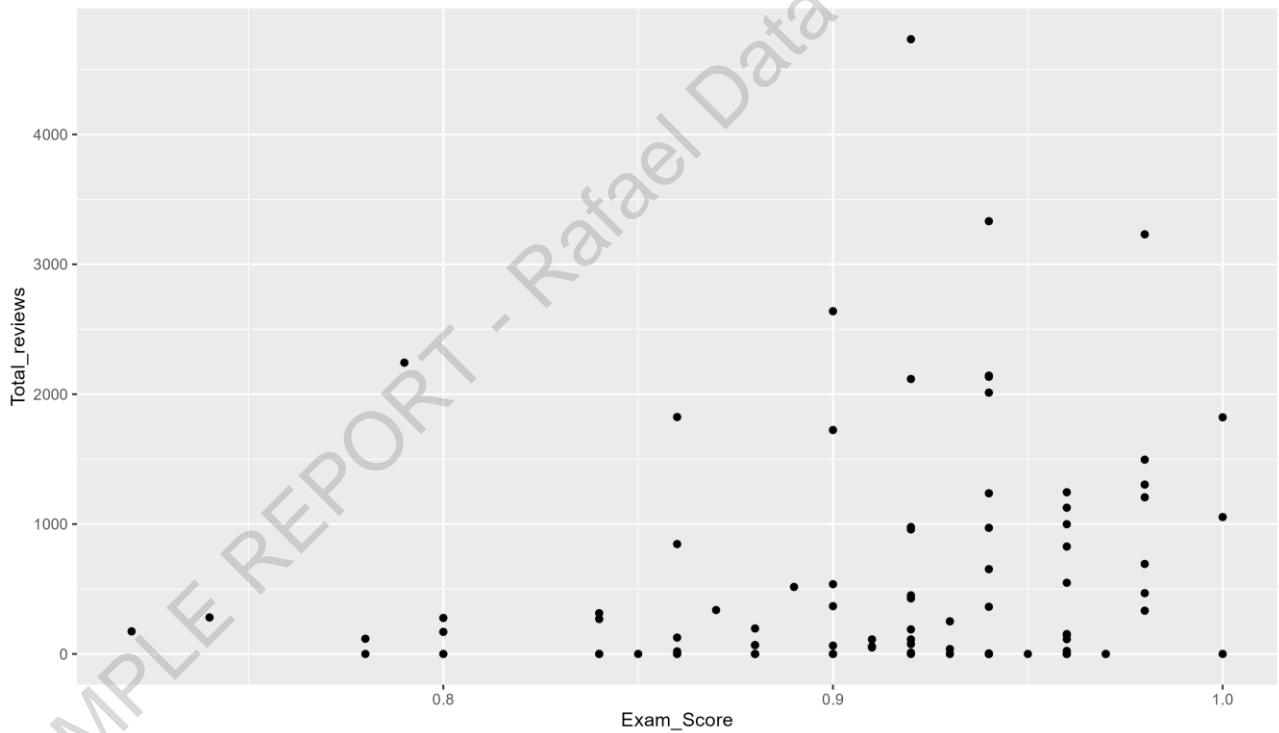## Scatterplot of Card_Young vs Exam_Score

Scatterplot of Days_studied_of_31 vs Exam_Score



Scatterplot of Total_reviews vs Exam_Score

## Linear Mixed Model

Lastly, a Linear Mixed Model was fit to the data. Not all variables were inserted here due to their correlational nature. We chose the 'Total_Reviews' variable due to having the highest correlation with exam score. Nevertheless, several models were attempted in order to explore for significant results, which were not achieved.

A Linear Mixed Model (LMM) is a statistical model designed to handle complex data structures, particularly those involving hierarchical or nested data. It's commonly used when data points are not independent of each other, which is often the case in repeated measures, longitudinal studies, or when data are collected from clusters or groups (like schools, hospitals, etc.). The key feature of an LMM is its ability to model both fixed and random effects.

Fixed effects represent the average relationship between the predictor(s) and the response variable across all groups or levels in the data. In our model, "Total_reviews" is the fixed effect.
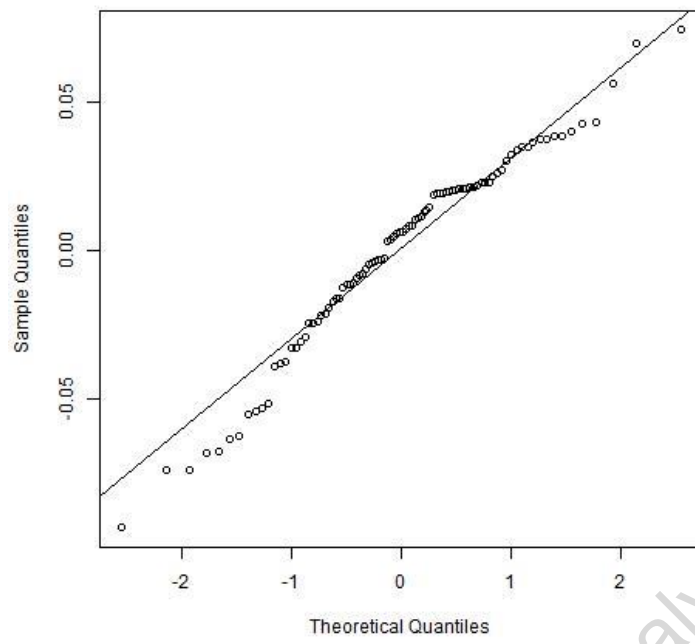
Random effects account for variations that are not captured by the fixed effects. They are used to model the impact of factors that vary across groups or levels, like individual differences or variations among clusters. In our model, "Exam" and "Student" are the random effects. This means the model accounts for the fact that exam scores might naturally vary from one exam to another and from one student to another, beyond what can be explained by the number of total reviews alone. The results are presented below.

| Variable | estimate | std.error | statistic | df | p.value |
|---|---|---|---|---|---|
| (Intercept) | 0.905 | 0.014 | 66.688 | 10.356 | 0.000 |
| Total_reviews | 9.85e-06 | 0.000 | 1.340 | 82.913 | 0.184 |

The estimate for Total_Reviews indicates the amount of change in the exam score expected for each one-unit increase in "Total_reviews". The estimate is very small (0.00000985), and with a p-value of 0.184, it's not statistically significant This implies that, after accounting for the variability in exams and students, the number of total reviews is not a strong predictor of exam scores within this particular dataset.

The figures below are used to assess the residual assumptions of the model. Both figures indicate that residuals follow a normal distribution, which makes the model valid and interpretable.

## Q-Q Plot



## Residuals vs Fitted