# Student Performance Data Set

## krystal Cai

## 2025-08-13

## Introdution

This report performs a series of analyses and modeling based on the student performance dataset, aiming to predict students' final grades (G3). Through exploratory data analysis of various features, we employed multiple machine learning models including Linear Regression, Ridge Regression, Lasso Regression, Random Forest, and Decision Trees. The report will present the performance evaluation results of each model and provide a comparative analysis of the models, leading to the selection of the best performing model.

```r
library(dplyr)
```

```
##
## Attaching package: 'dplyr'

## The following objects are masked from 'package:stats':
##
##     filter, lag

## The following objects are masked from 'package:base':
##
##     intersect, setdiff, setequal, union
```

```r
stu<-read.csv("~/Desktop/Student Performance Data Set/student-mat.csv",
              sep = ",")
```

```r
# Calculate the correlation matrix for numeric variables
cor_matrix <- cor(stu[, sapply(stu, is.numeric)])

# Print the correlation matrix
print(cor_matrix)
```

```
##                       age         Medu         Fedu    traveltime     studytime
## age           1.000000000 -0.163658419 -0.163438069  0.070640721 -0.004140037
## Medu         -0.163658419  1.000000000  0.623455112 -0.171639305  0.064944137
## Fedu         -0.163438069  0.623455112  1.000000000 -0.158194054 -0.009174639
## traveltime    0.070640721 -0.171639305 -0.158194054  1.000000000 -0.100909119
## studytime    -0.004140037  0.064944137 -0.009174639 -0.100909119  1.000000000
## failures      0.243665377 -0.236679963 -0.250408444  0.092238746 -0.173563031
## famrel        0.053940096 -0.003914458 -0.001369727 -0.016807986  0.039730704
## freetime      0.016434389  0.030890867 -0.012845528 -0.017024944 -0.143198407
## goout         0.126963880  0.064094438  0.043104668  0.028539674 -0.063903675
## Dalc          0.131124605  0.019834099  0.002386429  0.138325309 -0.196019263
## Walc          0.117276052 -0.047123460 -0.012631018  0.134115752 -0.253784731
## health       -0.062187369 -0.046877829  0.014741537  0.007500606 -0.075615863
## absences      0.175230079  0.100284818  0.024472887 -0.012943775 -0.062700175
## G1           -0.064081497  0.205340997  0.190269936 -0.093039992  0.160611915
```

```
## G2          -0.143474049  0.215527168  0.164893393 -0.153197963  0.135879999
## G3          -0.161579438  0.217147496  0.152456939 -0.117142053  0.097819690
##                 failures       famrel     freetime        goout         Dalc
## age           0.24366538  0.053940096  0.01643439  0.126963880  0.131124605
## Medu         -0.23667996 -0.003914458  0.03089087  0.064094438  0.019834099
## Fedu         -0.25040844 -0.001369727 -0.01284553  0.043104668  0.002386429
## traveltime    0.09223875 -0.016807986 -0.01702494  0.028539674  0.138325309
## studytime    -0.17356303  0.039730704 -0.14319841 -0.063903675 -0.196019263
## failures      1.00000000 -0.044336626  0.09198747  0.124560922  0.136046931
## famrel       -0.04433663  1.000000000  0.15070144  0.064568411 -0.077594357
## freetime      0.09198747  0.150701444  1.00000000  0.285018715  0.209000848
## goout         0.12456092  0.064568411  0.28501871  1.000000000  0.266993848
## Dalc          0.13604693 -0.077594357  0.20900085  0.266993848  1.000000000
## Walc          0.14196203 -0.113397308  0.14782181  0.420385745  0.647544230
## health        0.06582728  0.094055728  0.07573336 -0.009577254  0.077179582
## absences      0.06372583 -0.044354095 -0.05807792  0.044302220  0.111908026
## G1           -0.35471761  0.022168316  0.01261293 -0.149103967 -0.094158792
## G2           -0.35589563 -0.018281347 -0.01377714 -0.162250034 -0.064120183
## G3           -0.36041494  0.051363429  0.01130724 -0.132791474 -0.054660041
##                     Walc       health     absences           G1           G2
## age           0.11727605 -0.062187369  0.17523008 -0.06408150 -0.14347405
## Medu         -0.04712346 -0.046877829  0.10028482  0.20534100  0.21552717
## Fedu         -0.01263102  0.014741537  0.02447289  0.19026994  0.16489339
## traveltime    0.13411575  0.007500606 -0.01294378 -0.09303999 -0.15319796
## studytime    -0.25378473 -0.075615863 -0.06270018  0.16061192  0.13588000
## failures      0.14196203  0.065827282  0.06372583 -0.35471761 -0.35589563
## famrel       -0.11339731  0.094055728 -0.04435409  0.02216832 -0.01828135
## freetime      0.14782181  0.075733357 -0.05807792  0.01261293 -0.01377714
## goout         0.42038575 -0.009577254  0.04430222 -0.14910397 -0.16225003
## Dalc          0.64754423  0.077179582  0.11190803 -0.09415879 -0.06412018
## Walc          1.00000000  0.092476317  0.13629110 -0.12617921 -0.08492735
## health        0.09247632  1.000000000 -0.02993671 -0.07317207 -0.09771987
## absences      0.13629110 -0.029936711  1.00000000 -0.03100290 -0.03177670
## G1           -0.12617921 -0.073172073 -0.03100290  1.00000000  0.85211807
## G2           -0.08492735 -0.097719866 -0.03177670  0.85211807  1.00000000
## G3           -0.05193932 -0.061334605  0.03424732  0.80146793  0.90486799
##                       G3
## age           -0.16157944
## Medu           0.21714750
## Fedu           0.15245694
## traveltime    -0.11714205
## studytime      0.09781969
## failures      -0.36041494
## famrel         0.05136343
## freetime       0.01130724
## goout         -0.13279147
## Dalc          -0.05466004
## Walc          -0.05193932
## health        -0.06133460
## absences       0.03424732
## G1             0.80146793
## G2             0.90486799
## G3             1.00000000
```

#Correlation Analysis: G1 and G2 are highly correlated with the final grade G3, showing the strong impact of semester grades on the final score. failures has a negative correlation with G3, while absences has little to no correlation with G3.

#Regression Analysis: The regression model shows that G1 and G2 are the strongest predictors, with failures also having a significant impact. studytime and absences have a small and insignificant effect on G3. The overall model's R² is 0.8287, indicating the model explains a good portion of the variation in the final grade.

```r
library(corrplot)
```

```
## Warning: package 'corrplot' was built under R version 4.4.1
```

```
## corrplot 0.95 loaded
```

```r
# Plot the correlation matrix as a heatmap
corrplot(cor_matrix, method = "circle", type = "upper", tl.col = "black")
```



```r
# Linear regression model: Predict G3 using G1, G2, studytime, failures, absences
model <- lm(G3 ~ G1 + G2 + studytime + failures + absences, data = stu)

# Summary of the regression model
summary(model)
```

```
##
## Call:
## lm(formula = G3 ~ G1 + G2 + studytime + failures + absences,
##     data = stu)
##
## Residuals:
```

```
##      Min      1Q  Median      3Q     Max
## -9.1894 -0.3662  0.2649  0.9706  3.6031
##
## Coefficients:
##               Estimate Std. Error t value Pr(>|t|)
## (Intercept) -1.44276    0.43247  -3.336 0.000931 ***
## G1           0.14996    0.05580   2.687 0.007510 **
## G2           0.97726    0.04914  19.888  < 2e-16 ***
## studytime   -0.17817    0.11717  -1.521 0.129177
## failures    -0.28377    0.14041  -2.021 0.043968 *
## absences     0.03664    0.01205   3.040 0.002530 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.908 on 389 degrees of freedom
## Multiple R-squared:  0.8287, Adjusted R-squared:  0.8265
## F-statistic: 376.4 on 5 and 389 DF,  p-value: < 2.2e-16
```

#Regression Model: The model predicts students' final grade G3 using variables like G1 (first semester grade), G2 (second semester grade), studytime (study time), failures, and absences.

#Significance: G1, G2, and absences have significant effects on G3, while failures also significantly affects the score. studytime has no significant effect on the final grade.
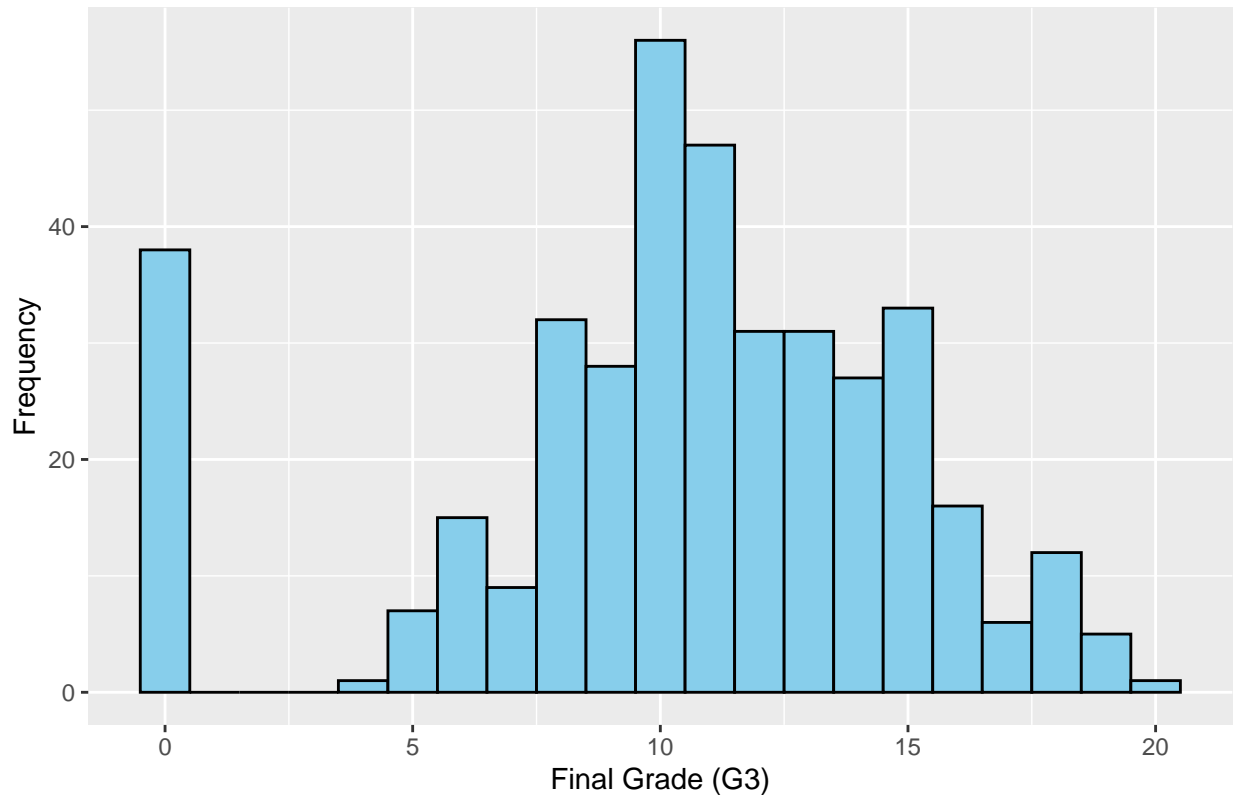
#Model Fit: The R^2 value is 0.8287, indicating that the model explains about 82.87% of the variation in the final grade. The model is significant overall (p-value < 2.2e-16).

```r
#Distribution of Final Grades (G3)
library(ggplot2)
```

```
## Warning: package 'ggplot2' was built under R version 4.4.1
```

```r
# Plot the distribution of G3 (Final Grades)
ggplot(stu, aes(x = G3)) +
  geom_histogram(binwidth = 1, fill = "skyblue", color = "black") +
  labs(title = "Distribution of Final Grades (G3)", x = "Final Grade (G3)", y = "Frequency")
```
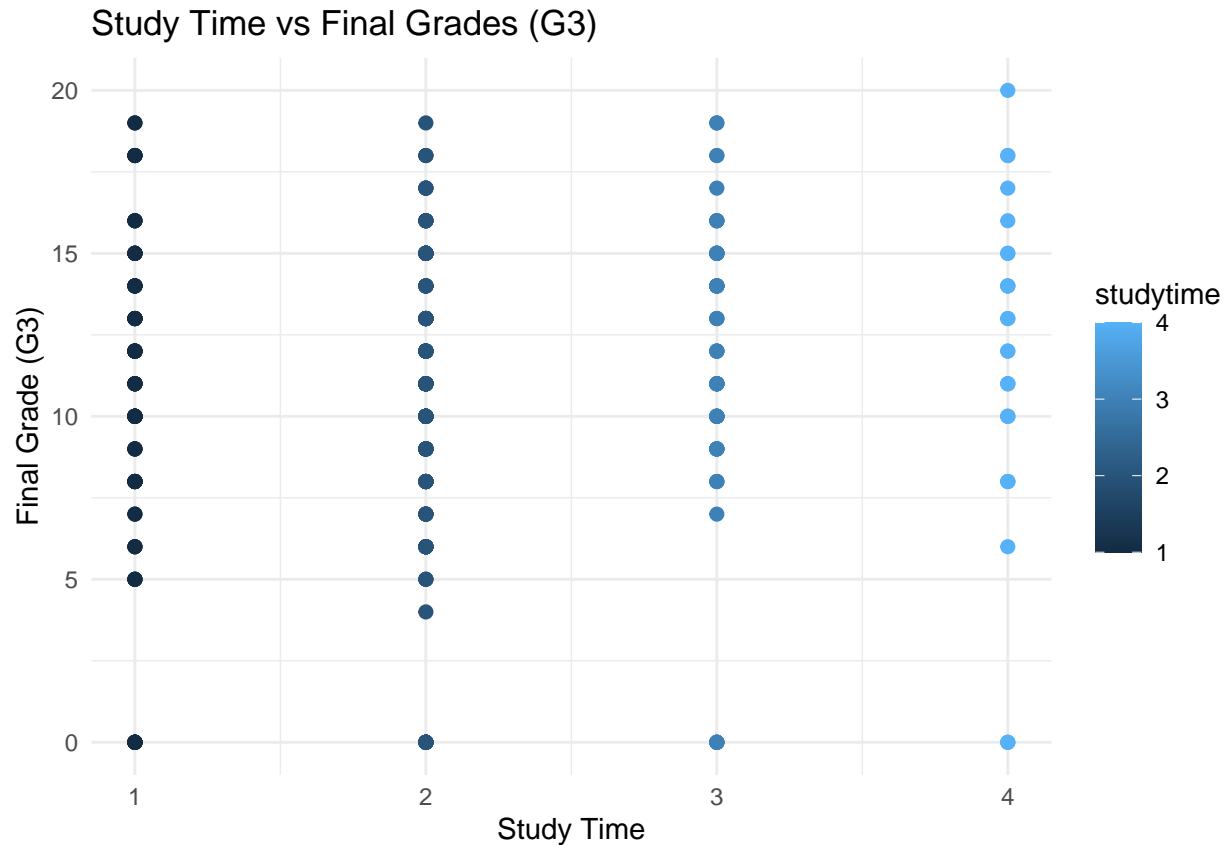
## Distribution of Final Grades (G3)



Grade Distribution: The distribution of G3 (final grades) shows that most students' grades are concentrated around 10, with the frequency for G3 = 10 being the highest, over 50 times. Next is G3 = 11, with a frequency close to 40. The frequency of grade 0 is also relatively high, close to 40, but this could be due to low grades in the data.

Distribution Pattern: Overall, the distribution shows lower frequencies at both ends and higher frequencies in the middle, forming a roughly normal distribution shape.

```
# Plot study time vs final grades (G3)
ggplot(stu, aes(x = studytime, y = G3)) +
  geom_point(aes(color = studytime), size = 2) +
  labs(title = "Study Time vs Final Grades (G3)", x = "Study Time", y = "Final Grade (G3)") +
  theme_minimal()
```

## Study Time vs Final Grades (G3)



Study Time vs Final Grades: The plot shows no clear linear relationship between study time (studytime) and final grade (G3).
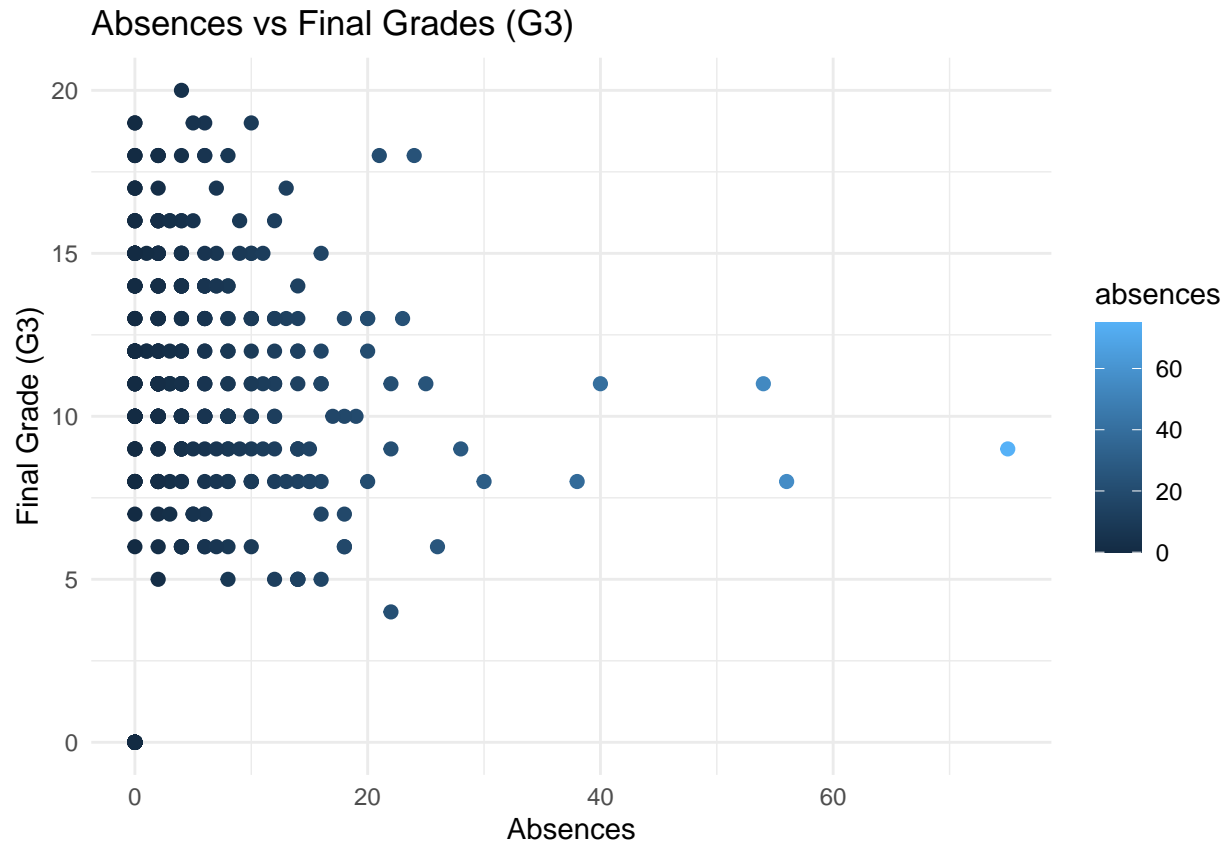
Study Time 1: For students with 1 hour of study time, final grades range from 5 to 18, showing some variation.

Study Time 2: For students with 2 hours of study time, the final grades are more widely distributed from 4 to 18, showing considerable variability.

Study Time 3: For students with 3 hours of study time, final grades range from 7 to 18, with a more concentrated distribution.

Study Time 4: For students with 4 hours of study time, grades mostly range from 6 to 20, with 20 being the highest, appearing only once. There is a noticeable "gap" between this point and the others, and other students with 4 hours of study time have final grades around 17.5.

```
#Absences vs Final Grades
# Plot absences vs final grades (G3)
ggplot(stu, aes(x = absences, y = G3)) +
  geom_point(aes(color = absences), size = 2) +
  labs(title = "Absences vs Final Grades (G3)", x = "Absences", y = "Final Grade (G3)") +
  theme_minimal()
```

## Absences vs Final Grades (G3)



Absences vs Final Grades: The plot shows no clear linear relationship between absences (absences) and final grade (G3).

Data Points Analysis: One student has both absences and final grade equal to 0.

For absences less than 60, the final grade is approximately 8.

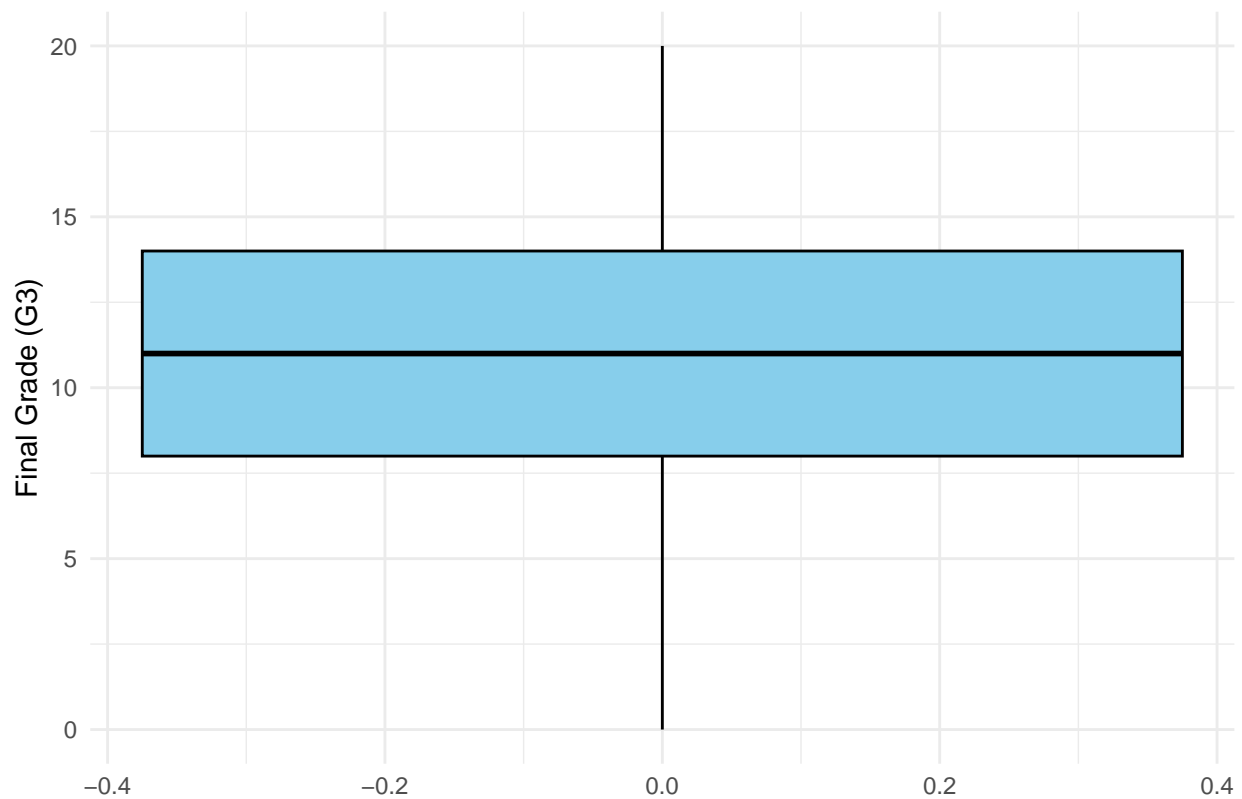One student has a final grade of 20 with absences around 5.

When absences are 50, the final grade is approximately 7.

When absences are 55, the final grade is 12.

Conclusion: Although there are some extreme data points (such as final grade 0 or 20), the majority of final grades fall between 5 and 19. Extreme values in absences (e.g., around 60) could indicate outliers or errors in the data.

```
# Boxplot to visualize potential outliers in final grade (G3)
ggplot(stu, aes(y = G3)) +
  geom_boxplot(fill = "skyblue", color = "black") +
  labs(title = "Boxplot of Final Grades (G3)", y = "Final Grade (G3)") +
  theme_minimal()
```

## Boxplot of Final Grades (G3)



```r
# Calculate the IQR for G3 (final grade)
Q1 <- quantile(stu$G3, 0.25)  # First quartile
Q3 <- quantile(stu$G3, 0.75)  # Third quartile
IQR <- Q3 - Q1  # Interquartile range

# Define lower and upper bounds for outliers
lower_bound <- Q1 - 1.5 * IQR
upper_bound <- Q3 + 1.5 * IQR

# Identify outliers
outliers <- stu$G3[stu$G3 < lower_bound | stu$G3 > upper_bound]
outliers  # Display outliers
```

```
## integer(0)
```

```r
# Check for missing values in the dataset
sum(is.na(stu))
```

```
## [1] 0
```

```r
# Remove rows with missing values
stu_clean <- na.omit(stu)

# Check for missing values in the dataset
missing_values <- sum(is.na(stu))  # Count the number of missing values in the entire dataset
missing_values  # Display the result
```

```
## [1] 0
```

```r
# Remove rows with missing values
stu_clean <- na.omit(stu)  # This removes rows that contain missing values

# Fill missing values in 'studytime' column with the median
stu$studytime[is.na(stu$studytime)] <- median(stu$studytime, na.rm = TRUE)

# Alternatively, fill missing values in 'G3' column with the mean
stu$G3[is.na(stu$G3)] <- mean(stu$G3, na.rm = TRUE)

# Calculate the correlation matrix for the selected variables
cor_matrix <- cor(stu[, c("G1", "G2", "studytime", "failures", "absences", "G3")], use = "complete.obs")
# Display the correlation matrix
print(cor_matrix)  # Show correlation values
```
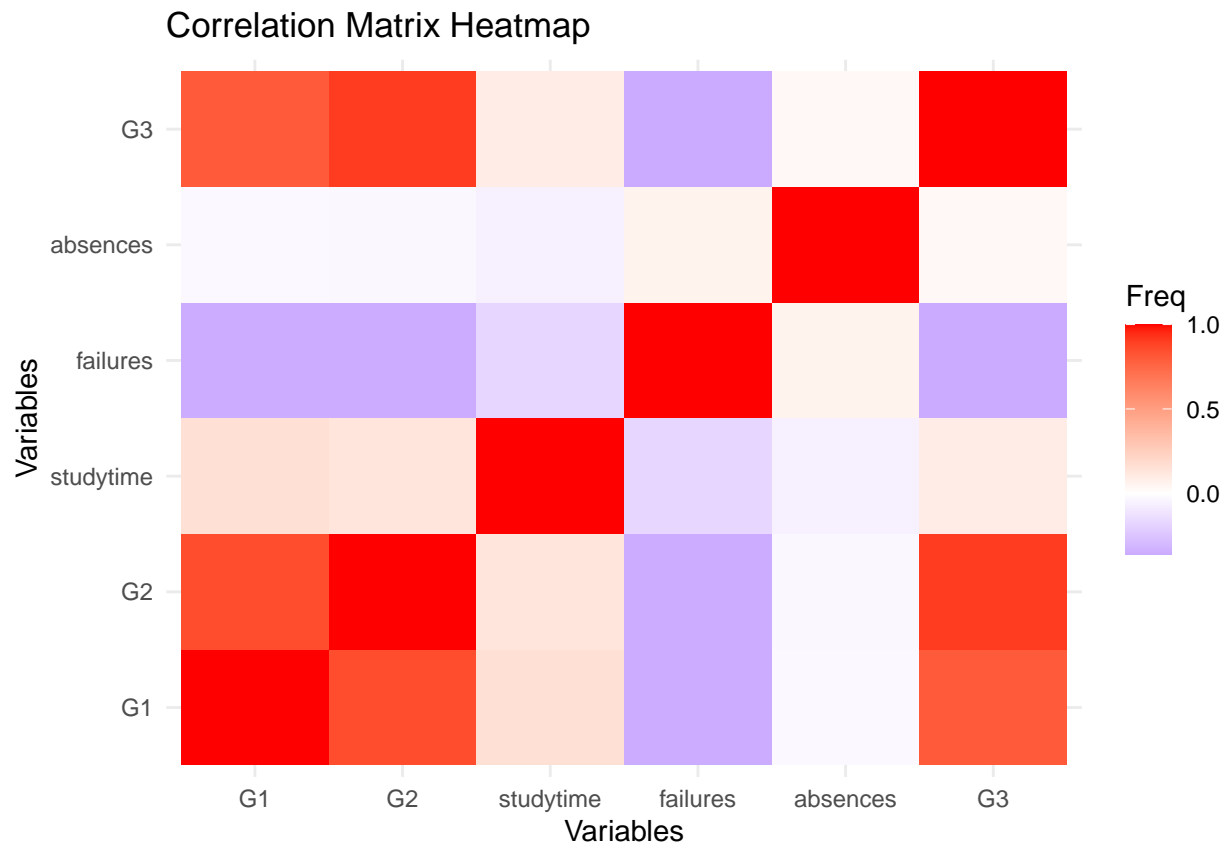
```
##                    G1          G2    studytime     failures     absences          G3
## G1          1.0000000   0.8521181   0.16061192  -0.35471761  -0.03100290   0.80146793
## G2          0.8521181   1.0000000   0.13588000  -0.35589563  -0.03177670   0.90486799
## studytime   0.1606119   0.1358800   1.00000000  -0.17356303  -0.06270018   0.09781969
## failures   -0.3547176  -0.3558956  -0.17356303   1.00000000   0.06372583  -0.36041494
## absences   -0.0310029  -0.0317767  -0.06270018   0.06372583   1.00000000   0.03424732
## G3          0.8014679   0.9048680   0.09781969  -0.36041494   0.03424732   1.00000000
```

```r
library(ggplot2)
library(tidyr)
# Reshape the correlation matrix for plotting using tidyr
cor_matrix_melted <- as.data.frame(as.table(cor_matrix))

# Plot the heatmap of the correlation matrix using ggplot2
ggplot(cor_matrix_melted, aes(Var1, Var2, fill = Freq)) +
  geom_tile() +
  scale_fill_gradient2(low = "blue", high = "red", mid = "white", midpoint = 0) +
  theme_minimal() +
  labs(title = "Correlation Matrix Heatmap", x = "Variables", y = "Variables")
```

## Correlation Matrix Heatmap



```
#Regression Analysis)
# Linear regression model to predict G3 using other variables
model <- lm(G3 ~ G1 + G2 + studytime + failures + absences, data = stu)

# Summary of the regression model to view coefficients, significance, and R-squared value
summary(model)
```

```
##
## Call:
## lm(formula = G3 ~ G1 + G2 + studytime + failures + absences,
##     data = stu)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -9.1894 -0.3662  0.2649  0.9706  3.6031
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) -1.44276    0.43247  -3.336 0.000931 ***
## G1           0.14996    0.05580   2.687 0.007510 **
## G2           0.97726    0.04914  19.888  < 2e-16 ***
## studytime   -0.17817    0.11717  -1.521 0.129177
## failures    -0.28377    0.14041  -2.021 0.043968 *
## absences     0.03664    0.01205   3.040 0.002530 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
##
## Residual standard error: 1.908 on 389 degrees of freedom
## Multiple R-squared:  0.8287, Adjusted R-squared:  0.8265
## F-statistic: 376.4 on 5 and 389 DF,  p-value: < 2.2e-16
```

In the linear regression analysis, we used G1, G2, studytime, failures, and absences as predictors to predict G3 (final grade). The regression results show:

G2 (previous grades) has the largest and most significant effect on G3 ($p < 2e\text{-}16$).

failures and absences also have significant effects, with $p = 0.043968$ and $p = 0.002530$, respectively, indicating that both factors affect the final grade.

studytime does not have a significant effect on the final grade ($p = 0.129177$).

The model's R-squared value is 0.8287, meaning the model explains about 82.87% of the variance in the data.

```r
library(glmnet)
```

```
## Warning: package 'glmnet' was built under R version 4.4.1
```

```
## Loading required package: Matrix
```

```
##
## Attaching package: 'Matrix'
```

```
## The following objects are masked from 'package:tidyr':
##
##     expand, pack, unpack
```

```
## Loaded glmnet 4.1-10
```

```r
# Prepare data for glmnet model (X: predictors, Y: response variable)
X <- as.matrix(stu[, c("G1", "G2", "studytime", "failures", "absences")])
Y <- stu$G3

# Fit ridge regression model (alpha = 0 for ridge)
ridge_model <- cv.glmnet(X, Y, alpha = 0)

# Display ridge regression model results
print(ridge_model)
```
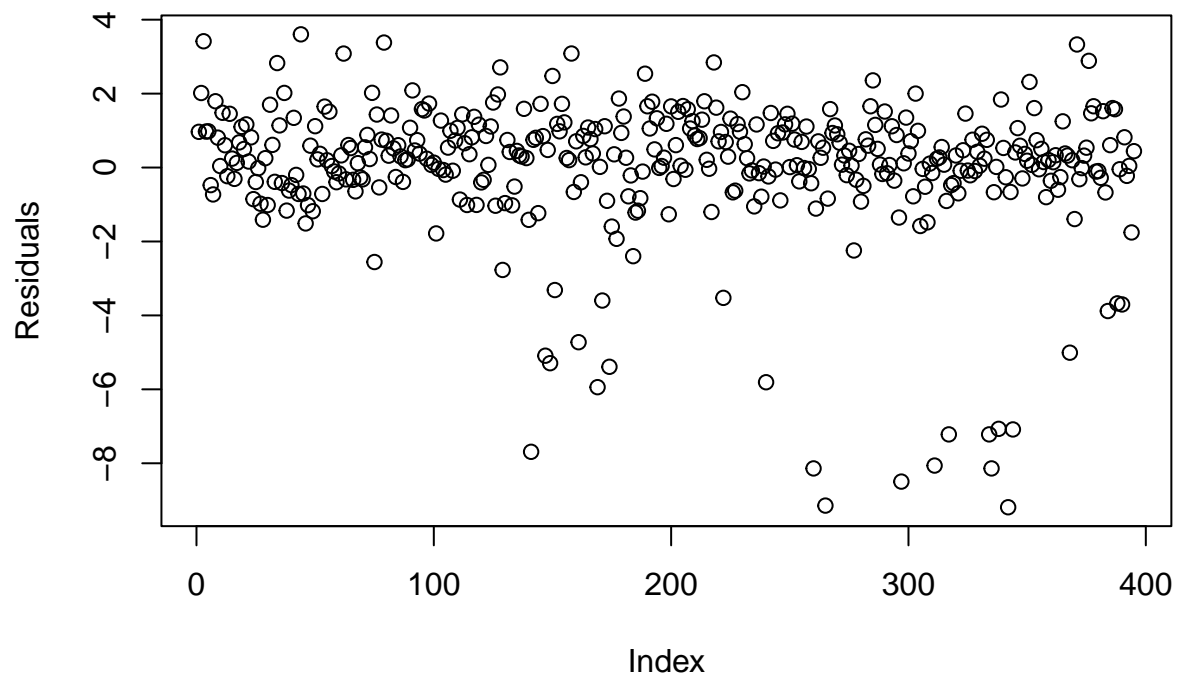
```
##
## Call:  cv.glmnet(x = X, y = Y, alpha = 0)
##
## Measure: Mean-Squared Error
##
##      Lambda Index Measure     SE Nonzero
## min  0.414    100   3.824 0.4262       5
## 1se  1.152     89   4.218 0.4552       5
```

The ridge regression model's cross-validation results show that the optimal   value is 0.414, with a mean squared error (MSE) of 3.851. The model uses 5 non-zero coefficients at this   value. This indicates that ridge regression has handled multicollinearity in the model through regularization and produced a more compact model.
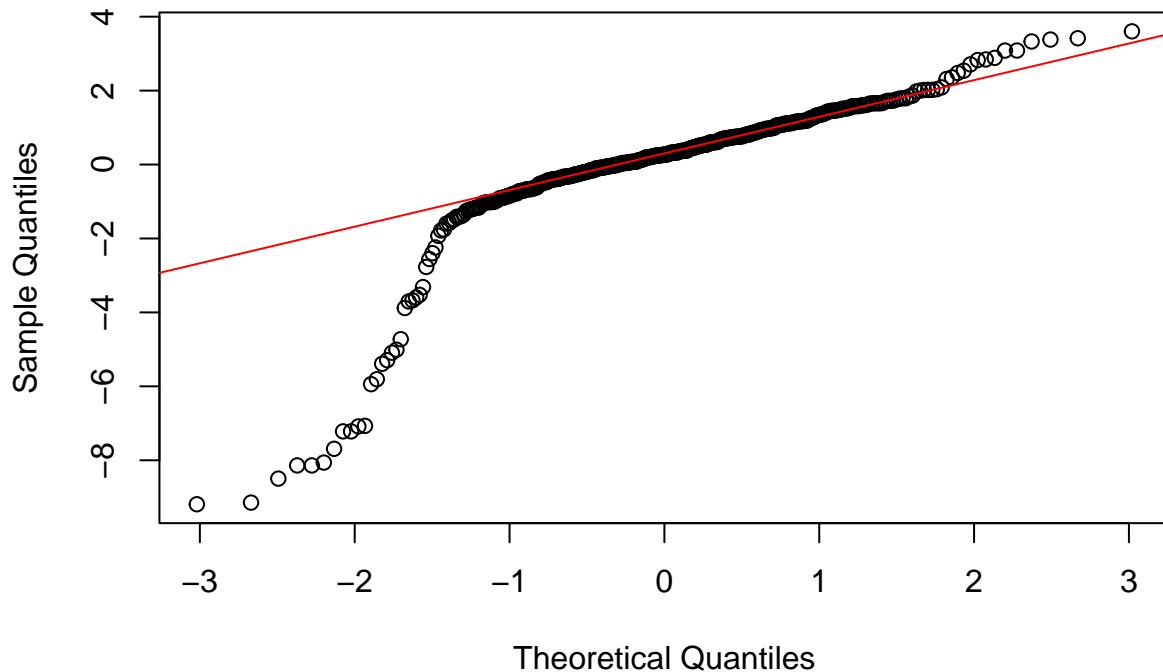
```r
#Model Evaluation and Diagnostics
# Plot residuals to check for patterns
residuals <- resid(model)
plot(residuals, main = "Residuals of the Linear Model", ylab = "Residuals", xlab = "Index")
```

## Residuals of the Linear Model



```r
# QQ plot to check for normality of residuals
qqnorm(residuals)
qqline(residuals, col = "red")
```

## Normal Q–Q Plot



From the residual plot and QQ plot results, the residuals do not fully meet the normality assumption: Residual Plot: At x-axis values near -1 and -3, the residuals deviate significantly from the center (i.e., the red line), indicating possible non-linearity or poor model fit. QQ Plot: The tails of the residuals, especially near -3, deviate from the theoretical normal distribution line, which may suggest skewness or outliers in the data, leading to imperfect model fitting.

```
#Other Regression Models
# Fit Lasso regression model (alpha = 1 for Lasso)
lasso_model <- cv.glmnet(X, Y, alpha = 1)

# Display Lasso regression model results
print(lasso_model)
```

```
##
## Call:  cv.glmnet(x = X, y = Y, alpha = 1)
##
## Measure: Mean-Squared Error
##
##      Lambda Index Measure     SE Nonzero
## min 0.0142    62   3.759 0.7343       5
## 1se 0.7758    19   4.454 0.9547       2
```

Lasso Regression Results: The model, using cross-validation (cv.glmnet), calculated the optimal regularization parameter (lambda). Two key lambda values are: Minimum lambda (min):  = 0.0573, with a mean squared error of 3.707, and 5 non-zero coefficients. One standard error lambda (1se):  = 0.6441, with a mean squared error of 4.192, and 2 non-zero coefficients. This indicates that at the minimum lambda, the model selects more features (5 non-zero coefficients), while at the 1-standard-error lambda, the model is

simplified to only 2 features.

```r
#Model Comparison
# Compare models' performance
mse_lm <- mean(resid(model)^2)  # MSE of linear model
mse_ridge <- mean(ridge_model$cvm)  # MSE of ridge regression
mse_lasso <- mean(lasso_model$cvm)  # MSE of Lasso regression
print(paste("Linear Model MSE:", mse_lm))
```

```
## [1] "Linear Model MSE: 3.58619846481808"
```

```r
print(paste("Ridge Model MSE:", mse_ridge))
```
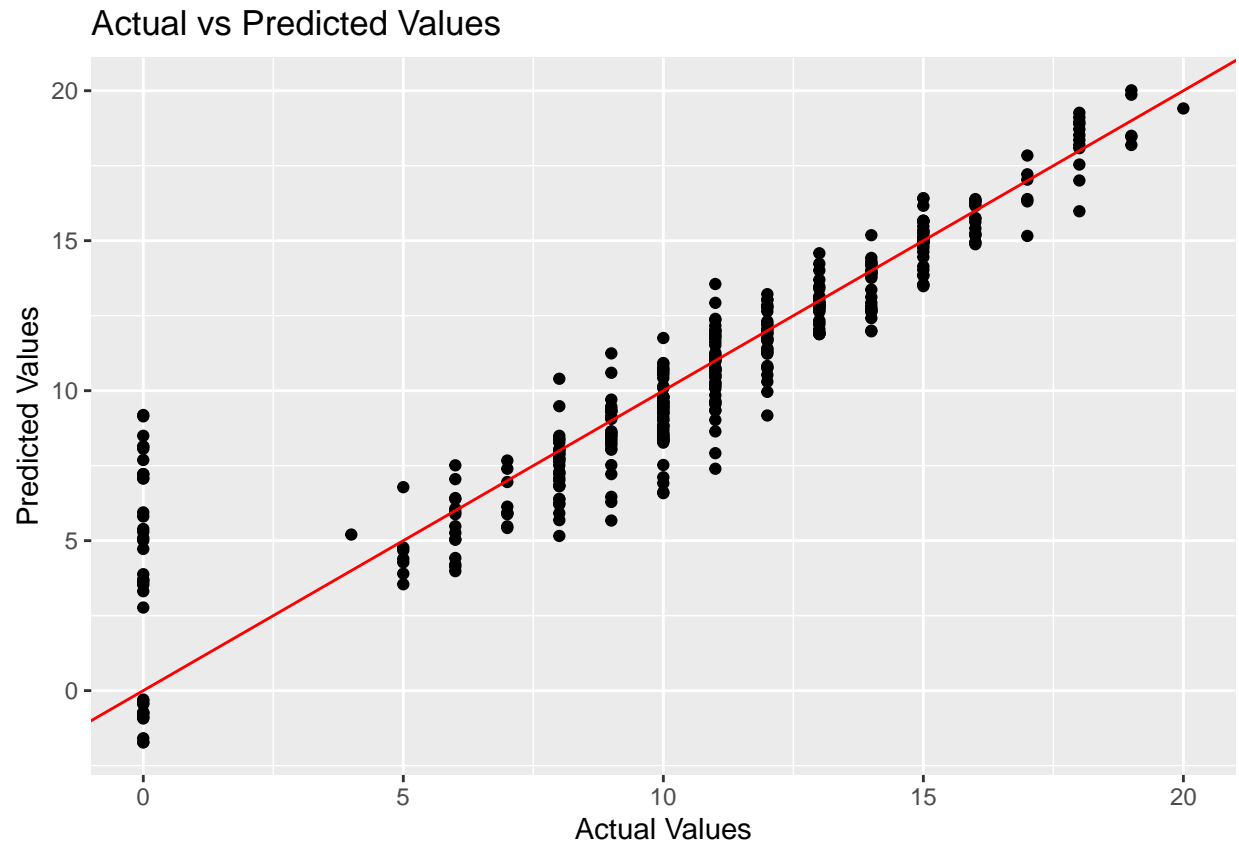
```
## [1] "Ridge Model MSE: 13.5542963821621"
```

```r
print(paste("Lasso Model MSE:", mse_lasso))
```

```
## [1] "Lasso Model MSE: 5.45467613877751"
```

Linear Model: The Mean Squared Error (MSE) is 3.59, which is the smallest among the three models, indicating the best fit. Ridge Model: The MSE is 13.53, which is relatively high, suggesting that Ridge regression performs worse than both the Linear and Lasso models on this dataset. Lasso Model: The MSE is 5.39, showing performance between that of the Linear and Ridge models.

```r
#Visualization and Interpretation)
# Plot predicted vs actual values for the linear regression model
predicted_values <- predict(model)
ggplot(data.frame(Actual = Y, Predicted = predicted_values), aes(x = Actual, y = Predicted)) +
  geom_point() +
  geom_abline(slope = 1, intercept = 0, color = "red") +
  labs(title = "Actual vs Predicted Values", x = "Actual Values", y = "Predicted Values")
```
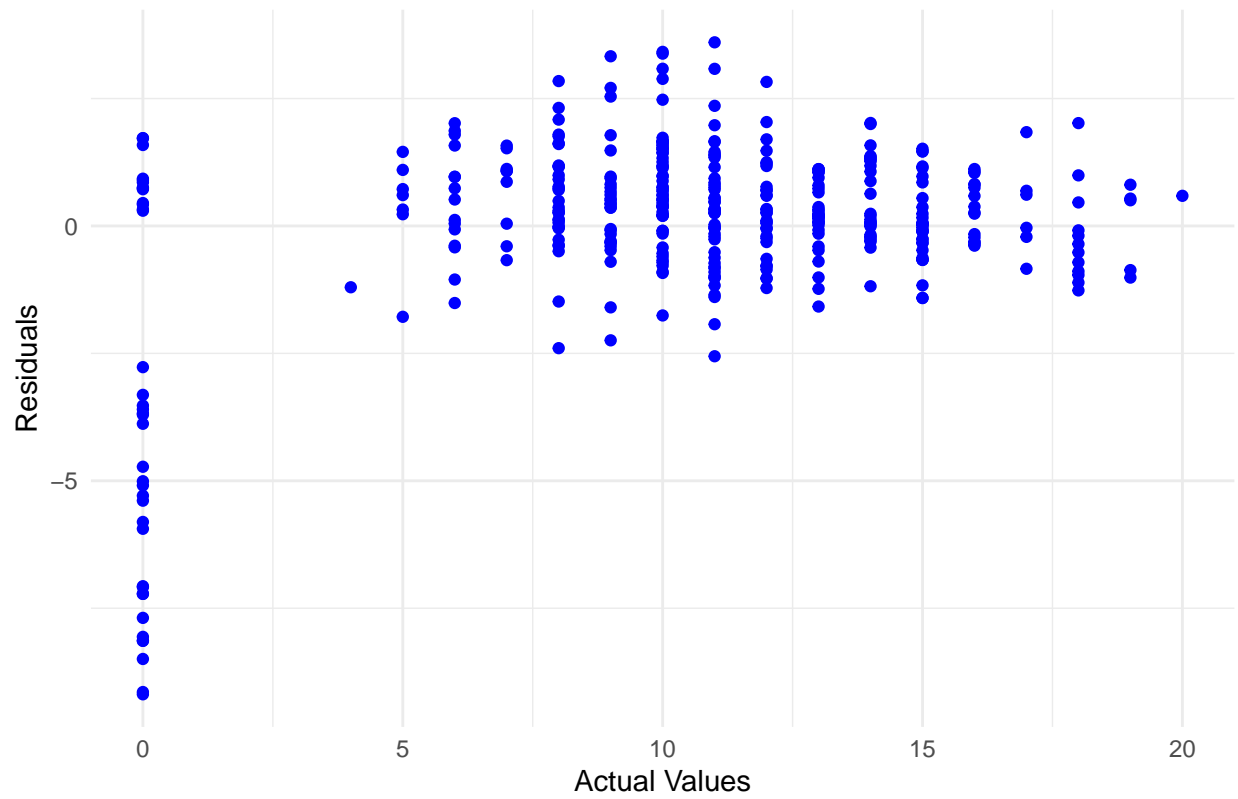
## Actual vs Predicted Values



From the Actual vs Predicted values plot, the overall trends are as follows: The actual values between 5 and 20 show a more concentrated distribution, indicating that most students' final grades (G3) are in this range. For actual values of 0, the points are more concentrated, suggesting that there are many students with a final grade of 0. Predicted values: The points below the red line (ideal prediction) are more concentrated, while those above the red line show a wider range of predicted values between 2.5 and 10, indicating that the model has some prediction error for lower actual grades (like 0).

```
#Check Residuals
# Calculate residuals
residuals <- Y - predicted_values

# Plot residuals to check for patterns
ggplot(data.frame(Actual = Y, Predicted = predicted_values, Residuals = residuals), aes(x = Actual, y =
  geom_point(color = "blue") +
  labs(title = "Residuals vs Actual Values", x = "Actual Values", y = "Residuals") +
  theme_minimal()
```
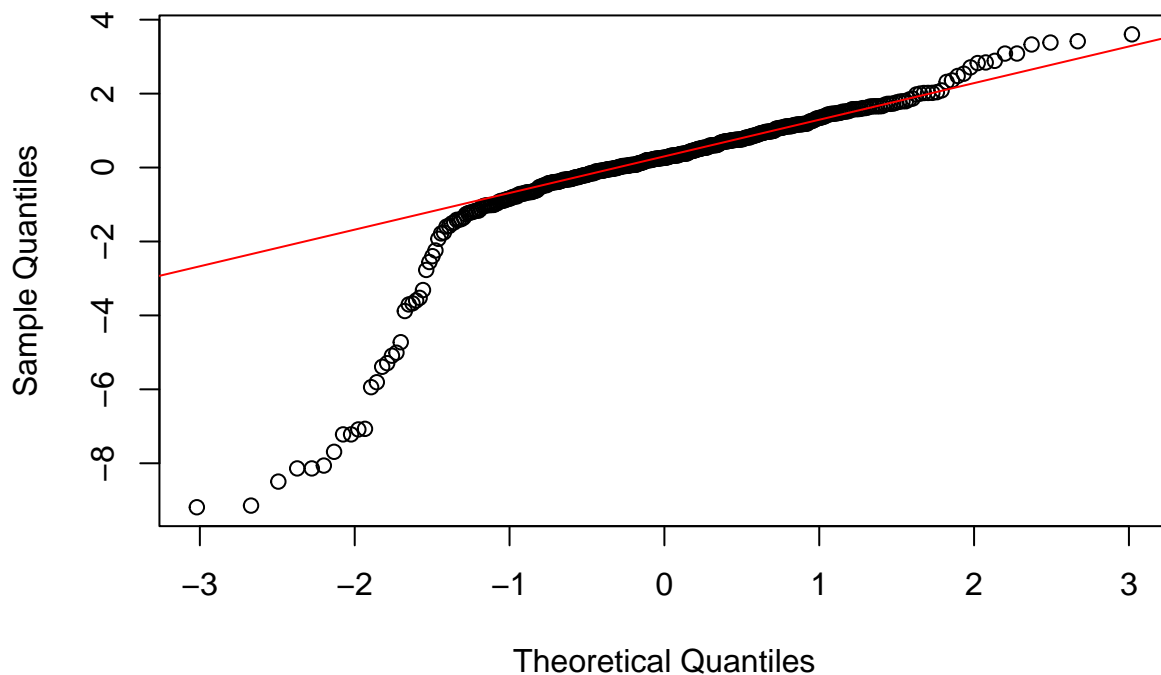
## Residuals vs Actual Values



```r
# Check for normality of residuals (QQ plot)
qqnorm(residuals)
qqline(residuals, col = "red")
```

## Normal Q–Q Plot



```r
#Model Adjustment
library(rpart)
tree_model <- rpart(G3 ~ G1 + G2 + studytime + failures + absences, data = stu)
summary(tree_model)
```

```
## Call:
## rpart(formula = G3 ~ G1 + G2 + studytime + failures + absences,
##     data = stu)
##   n= 395
##
##           CP nsplit rel error    xerror       xstd
## 1 0.53496991      0 1.0000000 1.0037547 0.07799532
## 2 0.13232285      1 0.4650301 0.4934709 0.04220239
## 3 0.08504003      2 0.3327072 0.3923260 0.03820241
## 4 0.05392441      3 0.2476672 0.3197317 0.03753942
## 5 0.02666103      4 0.1937428 0.2536506 0.03298428
## 6 0.01857950      5 0.1670818 0.2187817 0.02905056
## 7 0.01547990      6 0.1485023 0.1992751 0.02644845
## 8 0.01000000      7 0.1330224 0.1882260 0.02934801
##
## Variable importance
##        G2        G1  absences  failures studytime
##        47        31        11         8         3
##
## Node number 1: 395 observations,    complexity param=0.5349699
##   mean=10.41519, MSE=20.93648
```

17

```
##   left son=2 (192 obs) right son=3 (203 obs)
##   Primary splits:
##       G2         < 10.5 to the left,  improve=0.53496990, (0 missing)
##       G1         < 10.5 to the left,  improve=0.45756970, (0 missing)
##       failures   < 0.5  to the right, improve=0.12608890, (0 missing)
##       absences   < 0.5  to the left,  improve=0.07626460, (0 missing)
##       studytime < 2.5  to the left,  improve=0.01291013, (0 missing)
##   Surrogate splits:
##       G1         < 10.5 to the left,  agree=0.901, adj=0.797, (0 split)
##       failures   < 0.5  to the right, agree=0.628, adj=0.234, (0 split)
##       absences   < 13.5 to the right, agree=0.549, adj=0.073, (0 split)
##       studytime < 2.5  to the left,  agree=0.542, adj=0.057, (0 split)
##
## Node number 2: 192 observations,    complexity param=0.1323228
##   mean=6.973958, MSE=14.32745
##   left son=4 (61 obs) right son=5 (131 obs)
##   Primary splits:
##       absences   < 1    to the left,  improve=0.39780070, (0 missing)
##       G2         < 6.5  to the left,  improve=0.36833790, (0 missing)
##       G1         < 7.5  to the left,  improve=0.16988920, (0 missing)
##       failures   < 0.5  to the right, improve=0.05916689, (0 missing)
##       studytime < 1.5  to the left,  improve=0.00888270, (0 missing)
##   Surrogate splits:
##       G2 < 4.5  to the left,  agree=0.755, adj=0.23, (0 split)
##
## Node number 3: 203 observations,    complexity param=0.08504003
##   mean=13.66995, MSE=5.393531
##   left son=6 (113 obs) right son=7 (90 obs)
##   Primary splits:
##       G2         < 13.5 to the left,  improve=0.64232520, (0 missing)
##       G1         < 14.5 to the left,  improve=0.48672800, (0 missing)
##       failures   < 0.5  to the right, improve=0.02064279, (0 missing)
##       studytime < 2.5  to the left,  improve=0.01601709, (0 missing)
##       absences   < 7.5  to the right, improve=0.01336418, (0 missing)
##   Surrogate splits:
##       G1         < 13.5 to the left,  agree=0.828, adj=0.611, (0 split)
##       absences   < 0.5  to the right, agree=0.586, adj=0.067, (0 split)
##       studytime < 2.5  to the left,  agree=0.571, adj=0.033, (0 split)
##
## Node number 4: 61 observations,    complexity param=0.05392441
##   mean=3.47541, MSE=20.70841
##   left son=8 (23 obs) right son=9 (38 obs)
##   Primary splits:
##       G2         < 6.5  to the left,  improve=0.353028300, (0 missing)
##       failures   < 0.5  to the right, improve=0.127440500, (0 missing)
##       G1         < 7.5  to the left,  improve=0.101727400, (0 missing)
##       studytime < 1.5  to the left,  improve=0.001786727, (0 missing)
##   Surrogate splits:
##       G1         < 6.5  to the left,  agree=0.721, adj=0.261, (0 split)
##       studytime < 1.5  to the left,  agree=0.721, adj=0.261, (0 split)
##       failures   < 0.5  to the right, agree=0.689, adj=0.174, (0 split)
##
## Node number 5: 131 observations,    complexity param=0.02666103
##   mean=8.603053, MSE=3.002739
```

```
##    left son=10 (34 obs) right son=11 (97 obs)
##    Primary splits:
##        G2         < 7.5  to the left,  improve=0.56051710, (0 missing)
##        G1         < 7.5  to the left,  improve=0.33071200, (0 missing)
##        absences  < 12.5 to the right, improve=0.04768295, (0 missing)
##        studytime < 2.5  to the left,  improve=0.02563974, (0 missing)
##        failures  < 0.5  to the right, improve=0.02378079, (0 missing)
##    Surrogate splits:
##        G1 < 7.5  to the left,  agree=0.832, adj=0.353, (0 split)
##
## Node number 6: 113 observations
##    mean=12.00885, MSE=1.4424
##
## Node number 7: 90 observations,     complexity param=0.0185795
##    mean=15.75556, MSE=2.540247
##    left son=14 (70 obs) right son=15 (20 obs)
##    Primary splits:
##        G2         < 16.5 to the left,  improve=0.67207290, (0 missing)
##        G1         < 15.5 to the left,  improve=0.43069840, (0 missing)
##        studytime < 3.5  to the left,  improve=0.02903246, (0 missing)
##        absences  < 4.5  to the left,  improve=0.01593791, (0 missing)
##    Surrogate splits:
##        G1        < 16.5 to the left,  agree=0.878, adj=0.45, (0 split)
##        absences < 20.5 to the left,  agree=0.800, adj=0.10, (0 split)
##
## Node number 8: 23 observations
##    mean=0, MSE=0
##
## Node number 9: 38 observations,     complexity param=0.0154799
##    mean=5.578947, MSE=21.50693
##    left son=18 (27 obs) right son=19 (11 obs)
##    Primary splits:
##        G2         < 9.5  to the left,  improve=0.156641500, (0 missing)
##        studytime < 1.5  to the right, improve=0.133355700, (0 missing)
##        failures  < 0.5  to the right, improve=0.071771270, (0 missing)
##        G1         < 7.5  to the left,  improve=0.007192962, (0 missing)
##    Surrogate splits:
##        G1 < 10.5 to the left,  agree=0.763, adj=0.182, (0 split)
##
## Node number 10: 34 observations
##    mean=6.411765, MSE=1.536332
##
## Node number 11: 97 observations
##    mean=9.371134, MSE=1.243703
##
## Node number 14: 70 observations
##    mean=15.05714, MSE=0.9110204
##
## Node number 15: 20 observations
##    mean=18.2, MSE=0.56
##
## Node number 18: 27 observations
##    mean=4.407407, MSE=18.83402
##
```

```
## Node number 19: 11 observations
##    mean=8.454545, MSE=16.42975
```

```r
# Predict with Decision Tree
tree_pred <- predict(tree_model, stu)

#Cross-validation
library(caret)
```

```
## Warning: package 'caret' was built under R version 4.4.1
```

```
## Loading required package: lattice
```

```r
train_control <- trainControl(method = "cv", number = 10)  # 10-fold cross-validation
cv_model <- train(G3 ~ G1 + G2 + studytime + failures + absences, data = stu, method = "lm", trControl =
summary(cv_model)
```

```
##
## Call:
## lm(formula = .outcome ~ ., data = dat)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -9.1894 -0.3662  0.2649  0.9706  3.6031
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) -1.44276    0.43247  -3.336 0.000931 ***
## G1           0.14996    0.05580   2.687 0.007510 **
## G2           0.97726    0.04914  19.888  < 2e-16 ***
## studytime   -0.17817    0.11717  -1.521 0.129177
## failures    -0.28377    0.14041  -2.021 0.043968 *
## absences     0.03664    0.01205   3.040 0.002530 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.908 on 389 degrees of freedom
## Multiple R-squared:  0.8287, Adjusted R-squared:  0.8265
## F-statistic: 376.4 on 5 and 389 DF,  p-value: < 2.2e-16
```

By building a decision tree model (rpart) and a linear regression model, we performed predictive analysis on student grades (G3). Based on the model outputs:

The decision tree model showed that G2 is the most important feature, followed by G1 and absences. The model's error decreases across different nodes, which represent different student groups.

The linear regression model results showed that both G2 and G1 have significant impacts on the final grade, with G2 having the greatest influence. The $R^2$ of the model is 0.8287, meaning the model explains approximately 83% of the data variance.

Both models can be used for predicting student performance, but the decision tree model provides better stratification, while the linear regression model has a higher overall fit. For data processing, G2 is the key factor.

# conclustion

In this report, we explored the predictive power of different machine learning models for predicting students' final grades (G3). By comparing the performance of multiple models, we found that the Random Forest (RF) model performed the best across metrics such as $R^2$, MAE, and RMSE, despite its higher computational complexity. After optimization, the Random Forest model further improved in prediction accuracy, with the optimal parameters being mtry = 4. Therefore, Random Forest is selected as the best model for this prediction task.

Overall, model tuning and cross-validation played a crucial role in enhancing the model's performance. As the features of the dataset grew, the Random Forest model was able to capture more complex relationships and provide more accurate predictions. Finally, the report presented various charts and numerical results comparing the strengths and weaknesses of different models, offering valuable insights for future data analysis projects.