**Probability and Statistical Inference**
**Continuous Assessment Part II**


**Student Number:** D20123666

**Name:** Amit Kumar Gangwar

**Programme Code:** TU059


## Section 1 - Research Question

Whether a student's ability of Critical Reading and English language skills affects their competency as a citizen, their ability to understand the social world.


## Section 2 - Dataset

Dataset is about academic performance evolution for engineering students. The dataset contains the results in national assessments for secondary and university education in engineering students and contains academic, social, economic information for 12,411 students. The data was collected as part of the Master's Degree in Engineering project of the Technological University of Bolívar (UTB) titled Academic Efficiency Analysis in Engineering students

A full descriptor is available at:
https://www.sciencedirect.com/science/article/pii/S2352340920304315#utbl0001

Link To download dataset: https://data.mendeley.com/datasets/83tcx8psxv/1

**Reference**: Delahoz-Dominguez, E., Zuluaga, R., & Fontalvo-Herrera, T. (2020). Dataset of academic performance evolution for engineering students. *Data in Brief*, *30*, 105537. https://doi.org/10.1016/j.dib.2020.105537


- Missingness: There are no missing values in the Dataset
- The variables of interest are **Nature of School (SCHOOL_NAT), Critical Reading (CR_PRO), Citizen Competencies SPRO(CC_PRO), English (ENG_PRO)**
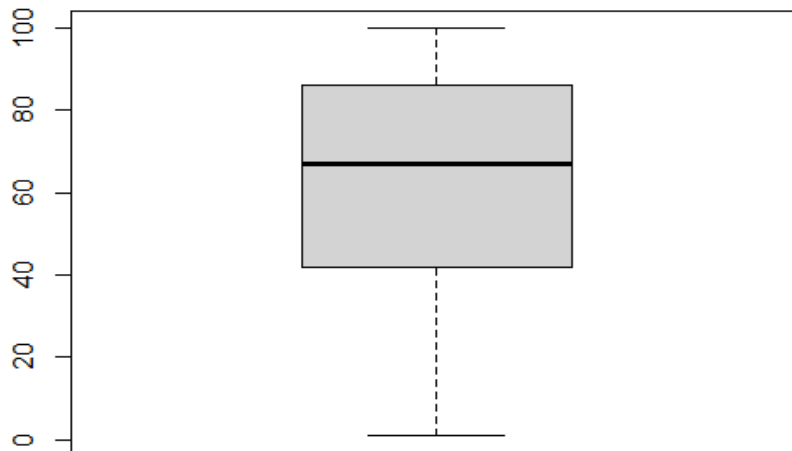- Relevant descriptive statistics and visuals:
  1. **Critical Reading (CR_PRO)**
     The sample as a whole has above average scores in critical reading (M = 62.2, SD = 27.66)

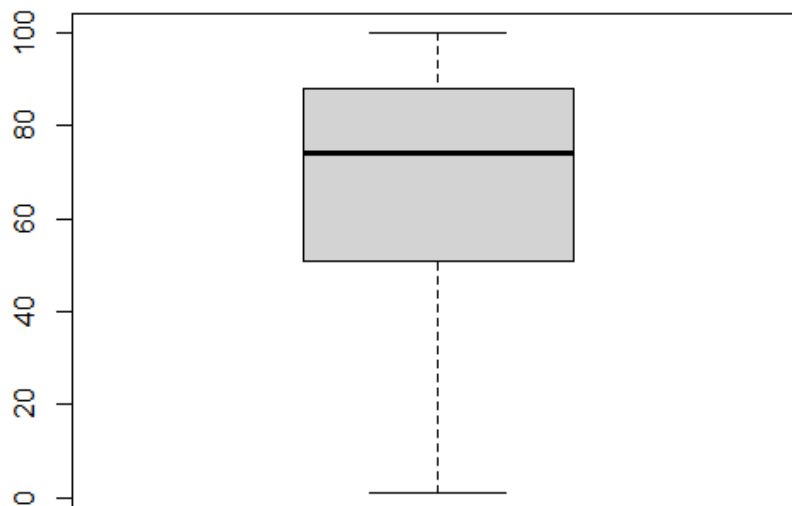| Min. | 1st Qu. | Median | Mean | 3rd Qu. | Max. | Std. Dev. | Variance | Outliers |
|------|---------|--------|------|---------|------|-----------|----------|----------|
| 1.0 | 42.0 | 67.0 | 62.2 | 86.0 | 100.0 | 27.66656 | 765.4384 | NA |

2. **English (ENG_PRO)**

The sample as a whole has above average scores in critical reading (M = 67.5, SD = 25.5)

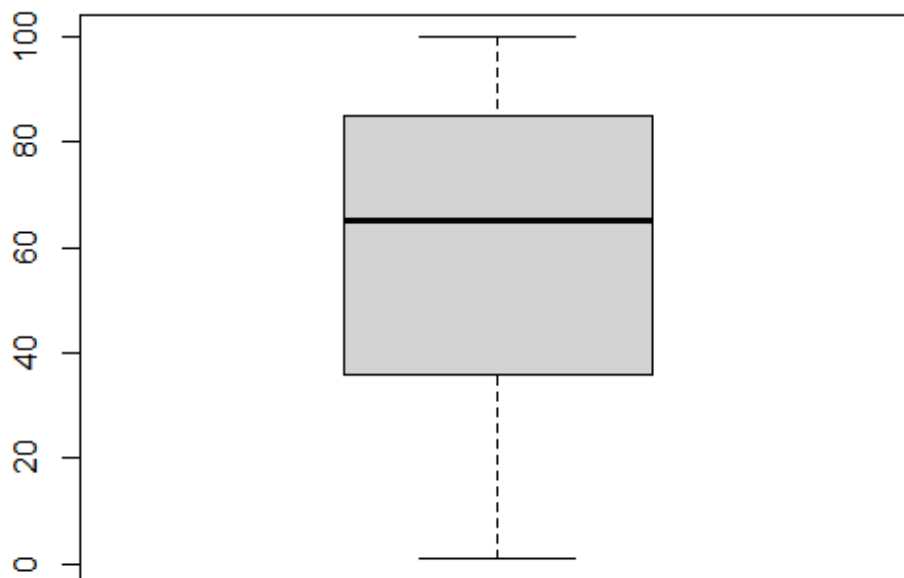| Min. | 1st Qu. | Median | Mean | 3rd Qu. | Max. | Std. Dev. | Variance | outliers |
|------|---------|--------|------|---------|------|-----------|----------|----------|
| 1.0 | 51.0 | 74.0 | 67.5 | 88.0 | 100.0 | 25.4951 | 649.9999 | NA |



3. **Citizen Competencies SPRO(CC_PRO)**

The sample as a whole has average scores in critical reading (M = 59.2, SD = 28.99)

| Min. | 1st Qu. | Median | Mean | 3rd Qu. | Max. | Std. Dev. | Variance | Outliers |
|------|---------|--------|-------|---------|--------|-----------|----------|----------|
| 1.00 | 36.00 | 65.00 | 59.19 | 85.00 | 100.00 | 28.99184 | 840.5268 | NA |



4. **Nature of School (SCHOOL_NAT)**
   Nearly half of the students studied in private school (52.89%) and rest half studied in public school (47.1).

   Total values: 12411

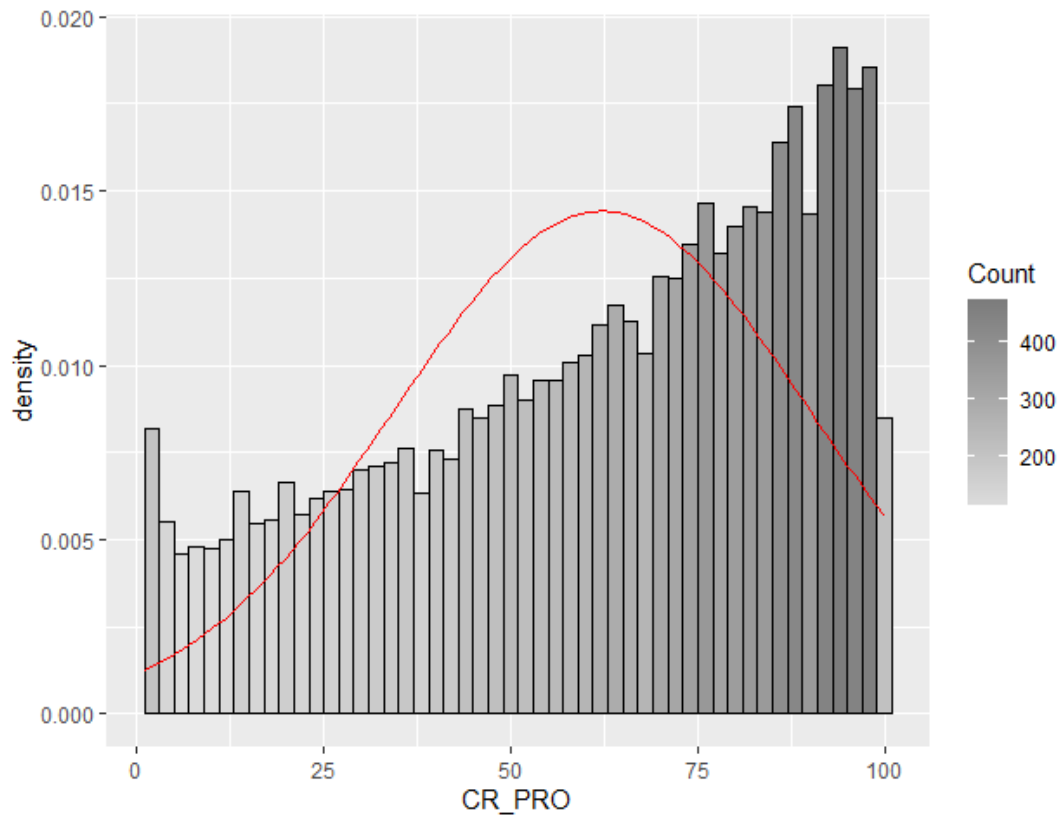   Distinct values: PRIVATE     PUBLIC

                    6565         5846
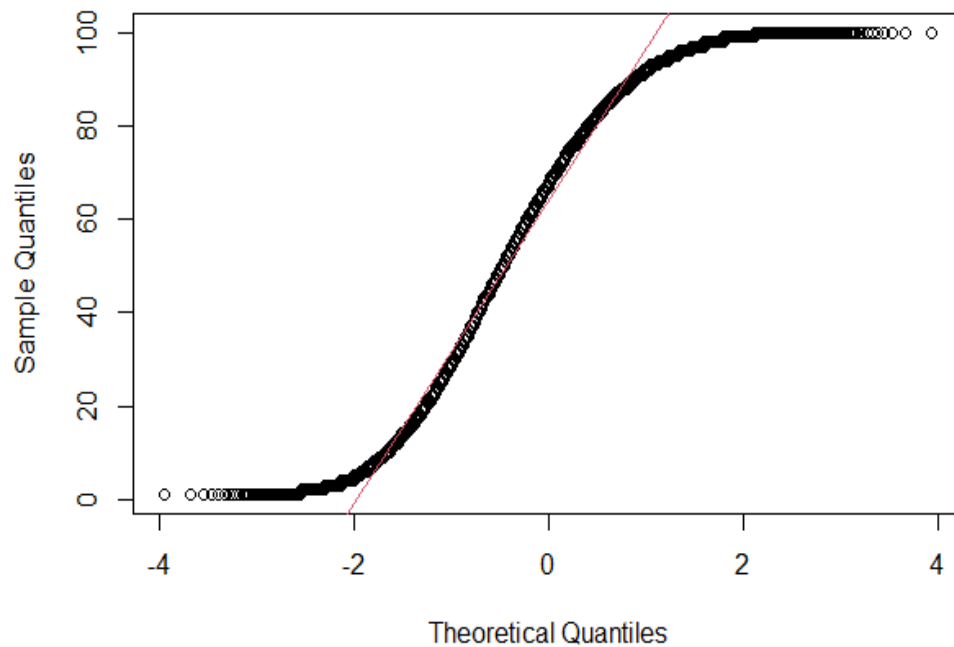
- Assessing normality:
  1. **Critical Reading (CR_PRO)**
     The input variable (CR_PRO) is approximately normally distributed as 99.9% of z-scores lie between $-3.29$ and $3.29$.

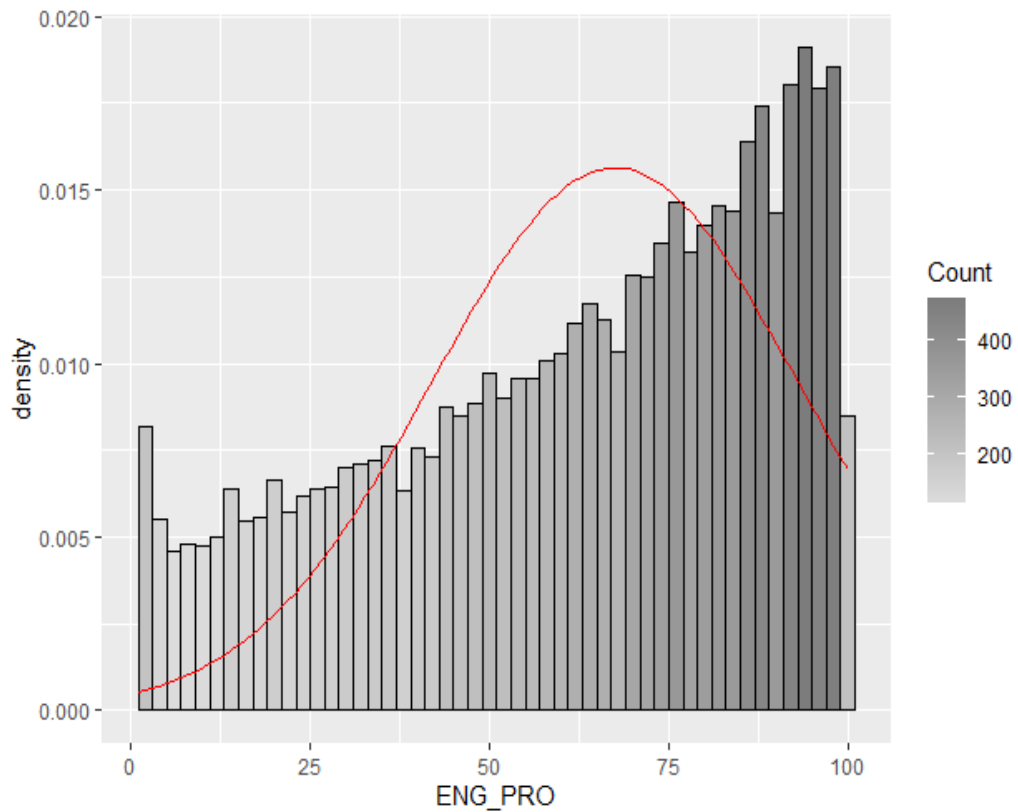2.  **English (ENG_PRO)**

The input variable (ENG_PRO) is approximately normally distributed as 99.9% of z-scores lie between −3.29 and 3.29.
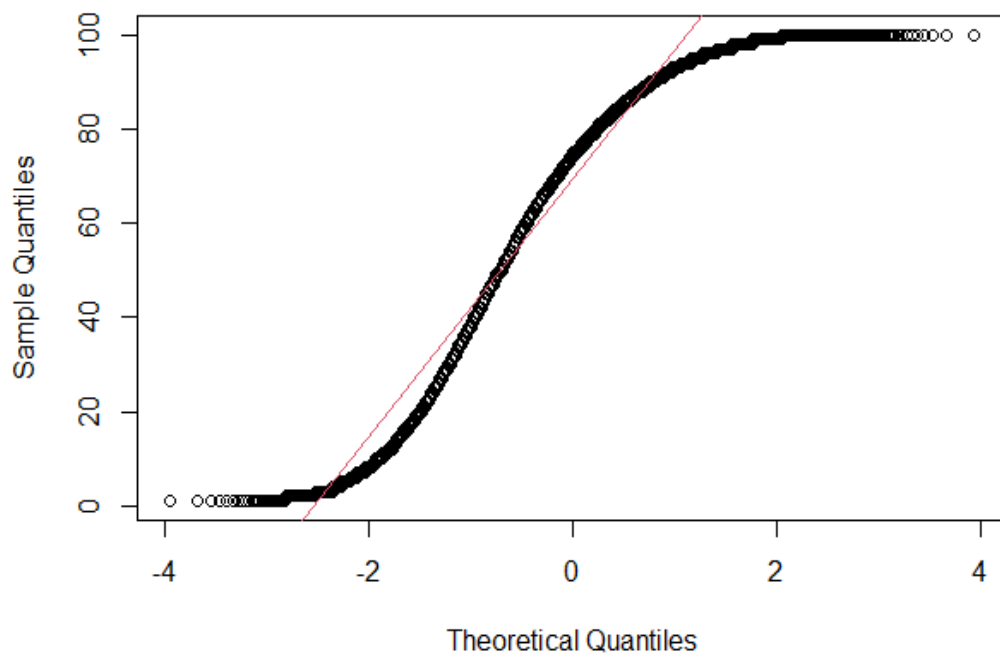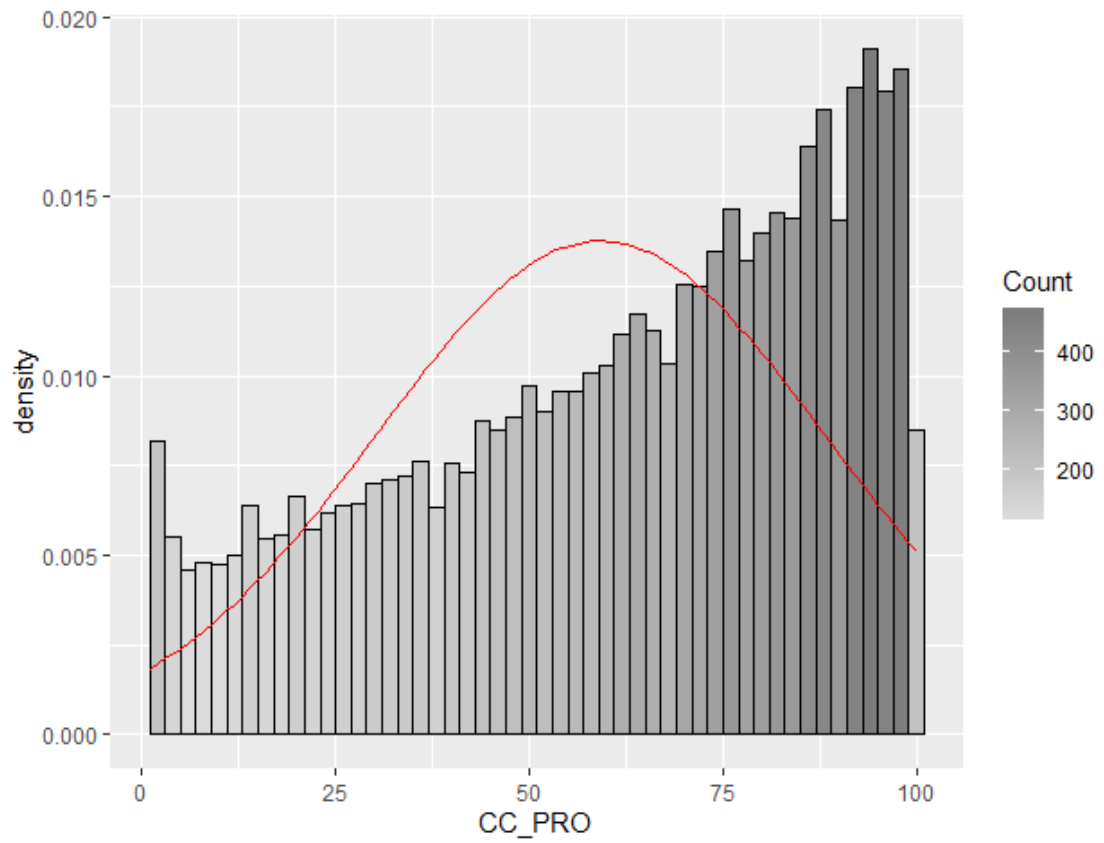
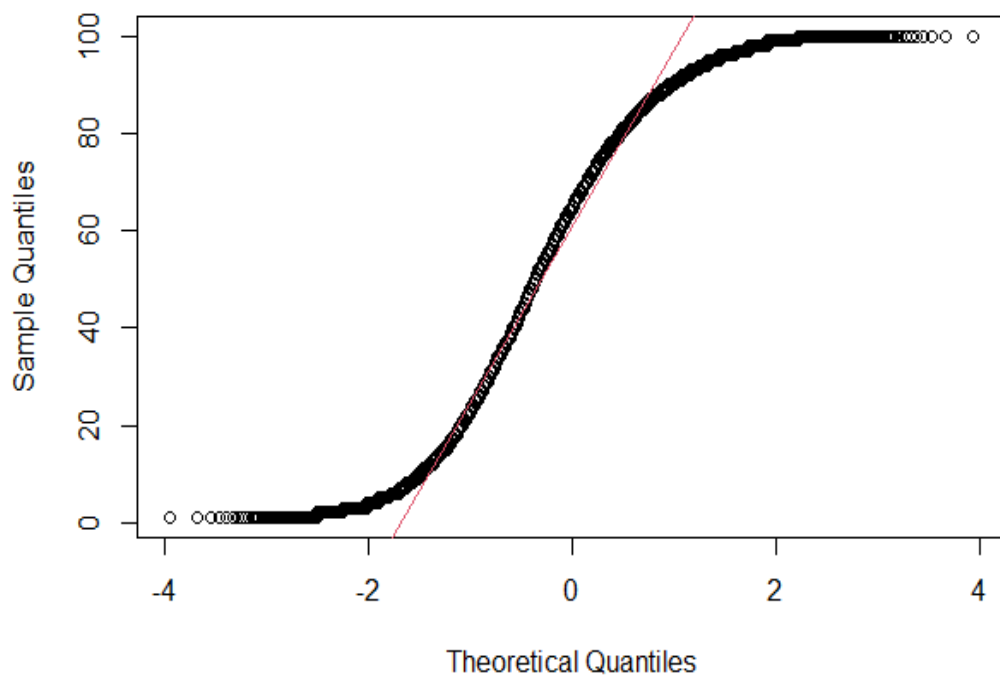3. **Citizen Competencies SPRO(CC_PRO)**

The response variable (CC_PRO) is approximately normally distributed as 99.9% of z-scores lie between $-3.29$ and $3.29$.

## Probability and Statistical Inference
## Continuous Assessment Part II



## Normal Q-Q Plot

## Section 3 - Results

### Section 3.1 - Statistical Evidence

1.  **Correlation Test between CR_SPRO and CC_PRO**:
    **Null Hypothesis**: There is a significant difference between CR_SPRO scores and CC_PRO scores of students.
    **Alternate Hypothesis**: There is no significant difference between CR_SPRO scores and CC_PRO scores of students.

    **Results:**
    A Pearson correlation coefficient was computed to assess the relationship between Critical Reading (CR_PRO) and Citizen Competencies SPRO(CC_PRO) exam marks. There was a significant positive correlation between the two variables, $r = 0.608$, $n = 12409$, $p < 2.2e\text{-}16$. A scatterplot below summarizes the results.



2.  **Correlation Test between ENG_PRO and CC_PRO:**
    **Null Hypothesis**: There is a significant difference between ENG_PRO scores and CC_PRO scores of students.
    **Alternate Hypothesis**: There is no significant difference between ENG_PRO scores and CC_PRO scores of students.

    **Results:**
    A Pearson correlation coefficient was computed to assess the relationship between English (ENG_PRO) and Citizen Competencies SPRO(CC_PRO) exam marks. There was a significant positive correlation between the two variables, $r = 0.48$, $n = 12409$, $p < 2.2e\text{-}16$. A scatterplot below summarizes the results.
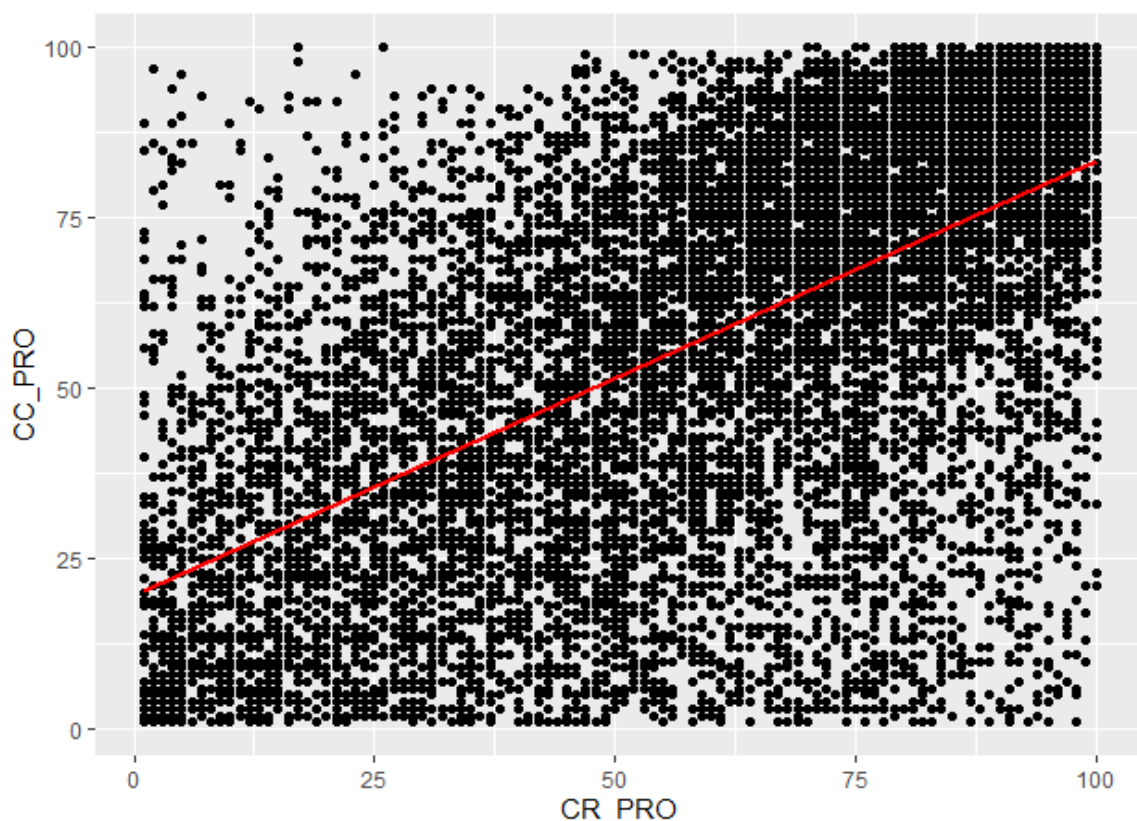
3. **Correlation Test between ENG_PRO and CR_SPRO:**
   **Null Hypothesis**: There is a significant difference between ENG_PRO scores and CR_SPRO scores of students.
   **Alternate Hypothesis**: There is no significant difference between ENG_PRO scores and CR_SPRO scores of students.

   **Results:**
   A Pearson correlation coefficient was computed to assess the relationship between English (ENG_PRO) and Critical Reading (CR_PRO) exam marks. There was a significant positive correlation between the two variables, r = 0.53, n = 12409, p < 2.2e-16. A scatterplot below summarizes the results. The collinearity between these independent variables doesn't exist since correlation is not very strong (r <8) so these variables can be used as input variables in the prediction model.

**Probability and Statistical Inference**
**Continuous Assessment Part II**



## Section 3.2 - Model1 Results

A multiple linear regression model was calculated to predict student's competency as a citizen (denoted by CC_PRO) based on their ability of Critical Reading (denoted by CR_PRO) and their English language skills (denoted by ENG_PRO). A significant regression equation was found **(F (2, 12408) = 4,201.48, p<0.01),** with an **R2 of 0.404**. Participants predicted competency score is equal to **10.333704+ 0.247549 (ENG_PRO) + 0.516789 (CR_PRO)**, where ENG_PRO and CR_PRO are scores which can vary from 0 to 100. Participant's CC_PRO score increases by 0.247549 for each unit in ENG_PRO and by 0.516789 unit for each unit increase in CR_PRO. Both CR_PRO and ENG_PRO were significant predictors of CC_PRO.

**CC_PRO = 10.333704+ 0.247549 (ENG_PRO) + 0.516789 (CR_PRO)**

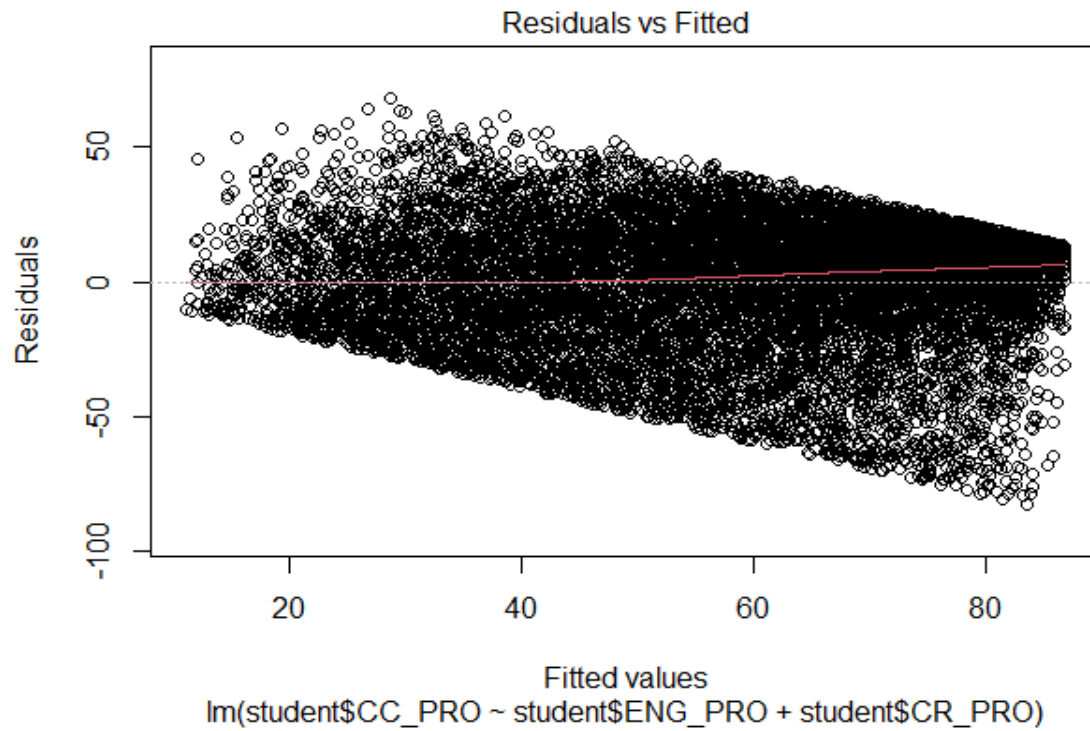Let's take mean values for ENG_PRO (67.5) and CR_PRO (62.2) and calculate CC_PRO

CC_PRO = 10.333704+ 0.247549 (67.5) + 0.516789 (62.2) = 59.09

59.09, Which is very close to the mean value of CC_PRO (59.2), this proves that model is significant with very small error term.
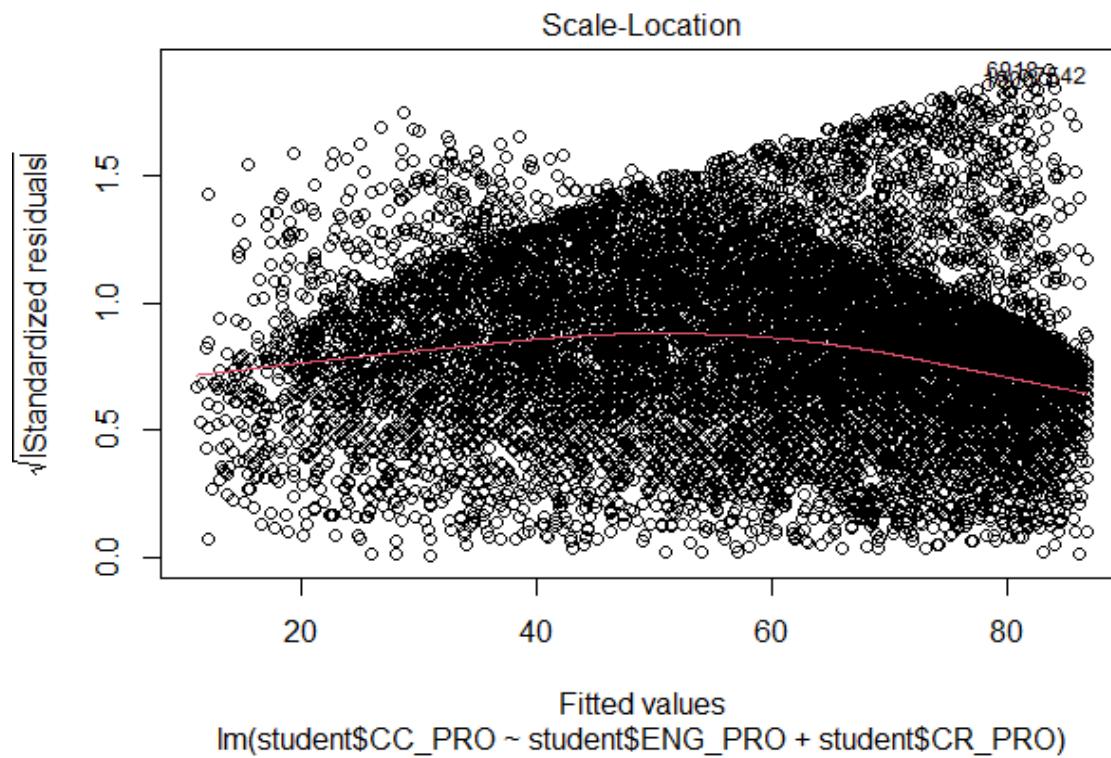
**Assessing the model against assumptions:**

**Homoscedasticity**

PLOT 1



PLOT 2

PLOT1 is the chart of residuals vs fitted values; in PLOT2 the standardised residuals are on the Y axis. If there is absolutely no heteroscedastity, we should see a completely random, equal distribution of points throughout the range of X axis and a flat red line. We really want to see that there is no pattern in the residuals and that they are equally spread around the y = 0 line - the dashed line.

As you can notice the red line is later slightly lifted upward on plot 1 but is not a big problem. Looking at the second plot we can see that while it is a problem it is not huge. Only a concern if there are definite patterns.

**Collinearity**

Tests to see if the data met the assumption of collinearity indicated that multicollinearity was not a concern (ENG_PRO, Tolerance = 0.7221794, VIF = 1.384697; CR_PRO, Tolerance = 0.7221794 VIF = 1.384697).

**Fit and Goodness**

The independent variables explain explains 40.4% (Adjusted $R^2$ = 0.404) variations in the students' CC_PRO scores (competency as a citizen).

# Section 3.3 - Model 2 Results

A multiple linear regression model was calculated to predict student's competency as a citizen (denoted by CC_PRO) based on their ability of Critical Reading (denoted by CR_PRO) and their English language skills (denoted by ENG_PRO) and the type of school they study in (denoted by SCHOOL_NAT). A significant regression equation was found **(F (3, 12407) = 2,804.05, p<0.01),** with an **$R^2$ of 0.404**. Participants predicted competency score is equal to **9.38 + 0.255 (ENG_PRO) + 0.515 (CR_PRO) + 1.042568 (SCHOOL_NAT), where SCHOOL_NAT is coded as 'PUBLIC=1', which is the dummy variable of interest and 'PRIVATE=0'**, ENG_PRO and CR_PRO are scores which can vary from 0 to 100. Participant's CC_PRO score increases by 0.255 for each unit in ENG_PRO and by 0.515 unit for each unit increase in CR_PRO and 1.042568 for School being Public where it is 0 for private school, so **there is a positive differential effect for Public schools**. All variables were significant predictors of CC_PRO while SCHOOL_NAT was the least significant with p=0.0153 (p<0.05).

**CC_PRO = 9.38 + 0.255 (ENG_PRO) + 0.515 (CR_PRO) + 1.042568 (SCHOOL_NAT)** for public schools

**CC_PRO = 9.38 + 0.255 (ENG_PRO) + 0.515 (CR_PRO)** for private schools

Let's take mean values for ENG_PRO (67.5) and CR_PRO (62.2) and calculate CC_PRO for Public school,

CC_PRO = 9.38 + 0.255 (67.5) + 0.515 (62.2) + 1.042568 (1) = 59.663

59.663, Which is very close to the mean value of CC_PRO (59.2), this proves that model is significant with very small error term.

**Assessing the model against assumptions:**

**Homoscedasticity**

Residuals vs Fitted



Fitted values
lm(student$CC_PRO ~ student$ENG_PRO + student$CR_PRO + student$SCHOOL_N/

**PLOT1**

Scale-Location



Fitted values
lm(student$CC_PRO ~ student$ENG_PRO + student$CR_PRO + student$SCHOOL_N/

**PLOT 2**

Probability and Statistical Inference
Continuous Assessment Part II

PLOT1 is the chart of residuals vs fitted values; in PLOT2 the standardised residuals are on the Y axis. If there is absolutely no heteroscedastity, we should see a completely random, equal distribution of points throughout the range of X axis and a flat red line. We really want to see that there is no pattern in the residuals and that they are equally spread around the y = 0 line - the dashed line.

As you can notice the red line is later slightly lifted upward on plot 1 but is not a big problem. Looking at the second plot we can see that while it is a problem it is not huge. Only a concern if there are definite patterns.

**Collinearity**

Tests to see if the data met the assumption of collinearity indicated that multicollinearity was not a concern (ENG_PRO, Tolerance = 0.6484227, VIF = 1.542204; CR_PRO, Tolerance = 0.7210604 VIF = 1.386846; SCHOOL_NAT, Tolerance = 0.8769922, VIF = 1.140261).

**Fit and Goodness**

The independent variables explain explains 40.4% (Adjusted R2 = 0.404) variations in the students' CC_PRO scores (competency as a citizen).

## Section 3. Model Comparison

Though both the models (Model 1 and Model 2) have the same fit as both have R2=0.404.

The model comparison (ANOVA) results show a Df of 1 (indicating that the more complex model has one additional parameter), and a p-value (p=0.0153, i.e., $p < 0.05$). This means that adding the SCHOOL_NAT to the model2 did lead to a slightly improved fit over the model 1, but not very significant.

## Code with Output:

```
1.  > #First unzip the Dataset file and place the xlsx file in the Working directory
2.  > #setting wroking directory
3.  > setwd('C:/Users/toami/OneDrive/Desktop/Data Science/SEM1/Probability & Stats/CA2')
4.  > par(mar=c(1,1,1,1))
5.  > #clearing environment
6.  > rm(list = ls())
7.  >
8.  > #installing required packages
9.  > needed_packages <- c("pastecs", "ggplot2", "semTools", "FSA","dplyr",
10. +                      "tidyr","outliers","ggplot2", "readxl")
11. > not_installed <- needed_packages[!(needed_packages %in% installed.packages()[ , "Package"])]
12. > if(length(not_installed)) install.packages(not_installed)
13. > library(pastecs)
14. > library(semTools)
15. >
16. > #loading library
17. > library(dplyr)
18. > library(tidyr)
19. > library(outliers)
20. > library(ggplot2)
21. > library(lm.beta)
22. > library(stargazer)
```

```
23. > library(readxl)
24. >
25. > #reading xlsx file
26. > student=read_excel("data_academic_performance.xlsx", sheet = "SABER11_SABERPRO")
27. New names:
28. * `` -> ...10
29. >
30. > #DATA CLEANING
31. > #removing blank column
32. > student = subset(student, select = -c(...10) )
33. >
34. > #checking for NA values
35. > student[is.na(student)==TRUE]
36. <unspecified> [0]
37. >
38. > #Checking for duplicate rows
39. > which(duplicated(student))
40. integer(0)
41. >
42. > #removing duplicates
43. > student=student[!duplicated(student),]
44. >
45. > #variables of interest
46. > #Nature of School(SCHOOL_NAT)
47. > #Critical Reading(CR_PRO)
48. > #Citizen Competencies SPRO(CC_PRO)
49. > #English(ENG_PRO)
50. >
51. > #summary (CR_PRO)
52. > summary(student$CR_PRO)
53.    Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
54.     1.0    42.0    67.0    62.2    86.0   100.0
55. > sd(student$CR_PRO)
56. [1] 27.66656
57. > var(student$CR_PRO)
58. [1] 765.4384
59. > Outlier = boxplot(student$CR_PRO)$out
60. > print(unique(Outlier))
61. numeric(0)
62. >
63. > #summary (ENG_PRO)
64. > summary(student$ENG_PRO)
65.    Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
66.     1.0    51.0    74.0    67.5    88.0   100.0
67. > sd(student$ENG_PRO)
68. [1] 25.4951
69. > var(student$ENG_PRO)
70. [1] 649.9999
71. > Outlier = boxplot(student$ENG_PRO)$out
72. > print(unique(Outlier))
73. numeric(0)
74. >
75. > #summary (CC_PRO)
76. > summary(student$CC_PRO)
77.    Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
78.     1.00   36.00   65.00   59.19   85.00  100.00
79. > sd(student$CC_PRO)
80. [1] 28.99184
81. > var(student$CC_PRO)
82. [1] 840.5268
83. > Outlier = boxplot(student$CC_PRO)$out
84. > print(unique(Outlier))
85. numeric(0)
86. >
87. > #summary SCHOOL_NAT
88. > summary(student$SCHOOL_NAT)
89.    Length     Class      Mode
90.     12411 character character
91. > table(student$SCHOOL_NAT)
92.
```

```
 93. PRIVATE  PUBLIC
 94.    6565    5846
 95. >
 96. > #===========================================================================
 97. > #NORMALITY CHECK OF VARIABLES OF INTEREST
 98. > #NORMALITY TEST of CR_PRO
 99. > gg <- ggplot(student, aes(x=student$CR_PRO))
100. > gg <- gg + labs(x="CR_PRO")
101. > gg <- gg + geom_histogram(binwidth=2, colour="black", aes(y=..density.., fill=..count..))
102. > gg <- gg + scale_fill_gradient("Count", low="#DCDCDC", high="#7C7C7C")
103. > gg <- gg + stat_function(fun=dnorm, color="red",
104. +                          args=list(mean=mean(student$CR_PRO, na.rm=TRUE),
105. +                                    sd=sd(student$CR_PRO, na.rm=TRUE)))
106. > #Plotting Histogram
107. > gg
108. >
109. > #Creating qqplot
110. > qqnorm(student$CR_PRO)
111. > qqline(student$CR_PRO, col=2)
112. >
113. > #Summary Statistics
114. > pastecs::stat.desc(student$CR_PRO, basic=F)
115.       median           mean      SE.mean CI.mean.0.95          var      std.dev
116.   67.0000000     62.1993393    0.2483429    0.4867906   765.4384234   27.6665579
117.      coef.var
118.     0.4448047
119. >
120. > #skew
121. > tpskew<-semTools::skew(student$CR_PRO)
122. > tpskew[1]/tpskew[2]
123. skew (g1)
124. -23.52302
125. >
126. > #kurtosis
127. > tpkurt<-semTools::kurtosis(student$CR_PRO)
128. > tpkurt[1]/tpkurt[2]
129. Excess Kur (g2)
130.        -19.24846
131. >
132. > ztpCR_PRO<- abs(scale(student$CR_PRO))
133. >
134. > #how much data fall outside 95% of region
135. > FSA::perc(as.numeric(ztpCR_PRO), 1.96, "gt")
136. [1] 3.658045
137. >
138. > #how much data fall outside 99.9% of region
139. > FSA::perc(as.numeric(ztpCR_PRO), 3.29, "gt")
140. [1] 0
141. >
142. > #---------------------------------------------------------------------------
143. > #NORMALITY TEST of ENG_PRO
144. > gg <- ggplot(student, aes(x=student$CR_PRO))
145. > gg <- gg + labs(x="ENG_PRO")
146. > gg <- gg + geom_histogram(binwidth=2, colour="black", aes(y=..density.., fill=..count..))
147. > gg <- gg + scale_fill_gradient("Count", low="#DCDCDC", high="#7C7C7C")
148. > gg <- gg + stat_function(fun=dnorm, color="red",
149. +                          args=list(mean=mean(student$ENG_PRO, na.rm=TRUE),
150. +                                    sd=sd(student$ENG_PRO, na.rm=TRUE)))
151. > #Plotting Histogram
152. > gg
153. >
154. > #Creating qqplot
155. > qqnorm(student$ENG_PRO)
156. > qqline(student$ENG_PRO, col=2)
157. >
158. > #Summary Statistics
159. > pastecs::stat.desc(student$ENG_PRO, basic=F)
160.       median           mean      SE.mean CI.mean.0.95          var      std.dev
161.   74.0000000     67.4983482    0.2288512    0.4485839   649.9998965   25.4950955
162.      coef.var
```

```
163.     0.3777144
164.  >
165.  > #skew
166.  > tpskew<-semTools::skew(student$ENG_PRO)
167.  > tpskew[1]/tpskew[2]
168.  skew (g1)
169.  -36.31304
170.  >
171.  > #kurtosis
172.  > tpkurt<-semTools::kurtosis(student$ENG_PRO)
173.  > tpkurt[1]/tpkurt[2]
174.  Excess Kur (g2)
175.      -7.222944
176.  >
177.  > ztpENG_PRO<- abs(scale(student$ENG_PRO))
178.  >
179.  > #how much data fall outside 95% of region
180.  > FSA::perc(as.numeric(ztpENG_PRO), 1.96, "gt")
181.  [1] 5.752961
182.  >
183.  > #how much data fall outside 99.9% of region
184.  > FSA::perc(as.numeric(ztpENG_PRO), 3.29, "gt")
185.  [1] 0
186.  >
187.  > #------------------------------------------------------------------------
188.  > #NORMALITY TEST of CC_PRO
189.  > gg <- ggplot(student, aes(x=student$CR_PRO))
190.  > gg <- gg + labs(x="CC_PRO")
191.  > gg <- gg + geom_histogram(binwidth=2, colour="black", aes(y=..density.., fill=..count..))
192.  > gg <- gg + scale_fill_gradient("Count", low="#DCDCDC", high="#7C7C7C")
193.  > gg <- gg + stat_function(fun=dnorm, color="red",
194.  +                     args=list(mean=mean(student$CC_PRO, na.rm=TRUE),
195.  +                                 sd=sd(student$CC_PRO, na.rm=TRUE)))
196.  > #Plotting Histogram
197.  > gg
198.  >
199.  > #Creating qqplot
200.  > qqnorm(student$CC_PRO)
201.  > qqline(student$CC_PRO, col=2)
202.  >
203.  > #Summary Statistics
204.  > pastecs::stat.desc(student$CC_PRO, basic=F)
205.      median        mean    SE.mean CI.mean.0.95         var      std.dev
206.   65.0000000   59.1867698   0.2602390   0.5101088   840.5267581   28.9918395
207.     coef.var
208.    0.4898365
209.  >
210.  > #skew
211.  > tpskew<-semTools::skew(student$CC_PRO)
212.  > tpskew[1]/tpskew[2]
213.  skew (g1)
214.  -19.08356
215.  >
216.  > #kurtosis
217.  > tpkurt<-semTools::kurtosis(student$CC_PRO)
218.  > tpkurt[1]/tpkurt[2]
219.  Excess Kur (g2)
220.      -24.02213
221.  >
222.  > ztpCC_PRO<- abs(scale(student$CC_PRO))
223.  >
224.  > #how much data fall outside 95% of region
225.  > FSA::perc(as.numeric(ztpCC_PRO), 1.96, "gt")
226.  [1] 1.232777
227.  >
228.  > #how much data fall outside 99.9% of region
229.  > FSA::perc(as.numeric(ztpCC_PRO), 3.29, "gt")
230.  [1] 0
231.  > #========================================================================
232.  > #CORRELATION CHECK AMONG VARIABLES OF INTEREST
```

```
233.  > #correlation between CR_PRO and CC_PRO
234.  > #Scatterplot
235.  > scatter <- ggplot(student, aes(CR_PRO, CC_PRO))
236.  > scatter + geom_point() + geom_smooth(method = "lm", colour = "Red", se = F) + labs(x =
      "CR_PRO", y = "CC_PRO")
237.  `geom_smooth()` using formula 'y ~ x'
238.  >
239.  > #Pearson Correlation
240.  > stats::cor.test(student$CR_PRO, student$CC_PRO, method='pearson')
241.
242.          Pearson's product-moment correlation
243.
244.  data:  student$CR_PRO and student$CC_PRO
245.  t = 85.287, df = 12409, p-value < 2.2e-16
246.  alternative hypothesis: true correlation is not equal to 0
247.  95 percent confidence interval:
248.   0.5966959 0.6188821
249.  sample estimates:
250.        cor
251.  0.6079076
252.
253.  >
254.  > #correlation between ENG_PRO and CC_PRO
255.  > #Scatterplot
256.  > scatter <- ggplot(student, aes(ENG_PRO, CC_PRO))
257.  > scatter + geom_point() + geom_smooth(method = "lm", colour = "Red", se = F) + labs(x =
      "ENG_PRO", y = "CC_PRO")
258.  `geom_smooth()` using formula 'y ~ x'
259.  >
260.  > #Pearson Correlation
261.  > stats::cor.test(student$ENG_PRO, student$CC_PRO, method='pearson')
262.
263.          Pearson's product-moment correlation
264.
265.  data:  student$ENG_PRO and student$CC_PRO
266.  t = 60.561, df = 12409, p-value < 2.2e-16
267.  alternative hypothesis: true correlation is not equal to 0
268.  95 percent confidence interval:
269.   0.4639377 0.4910993
270.  sample estimates:
271.        cor
272.  0.4776327
273.
274.  >
275.  > #correlation between ENG_PRO and CR_PRO
276.  > #Scatterplot
277.  > scatter <- ggplot(student, aes(ENG_PRO, CR_PRO))
278.  > scatter + geom_point() + geom_smooth(method = "lm", colour = "Red", se = F) + labs(x =
      "ENG_PRO", y = "CR_PRO")
279.  `geom_smooth()` using formula 'y ~ x'
280.  >
281.  > #Pearson Correlation
282.  > stats::cor.test(student$ENG_PRO, student$CR_PRO, method='pearson')
283.
284.          Pearson's product-moment correlation
285.
286.  data:  student$ENG_PRO and student$CR_PRO
287.  t = 69.092, df = 12409, p-value < 2.2e-16
288.  alternative hypothesis: true correlation is not equal to 0
289.  95 percent confidence interval:
290.   0.5142623 0.5396758
291.  sample estimates:
292.        cor
293.  0.5270869
294.
295.  >
296.  > #=================================================================================
297.  > #MULTIPLE REGRESSION MODEL 1
298.  > #building multiple regression model
299.  > model1<-lm(student$CC_PRO~student$ENG_PRO+student$CR_PRO)
```

```
300.  >
301.  > #Analysis of Variance Table
302.  > anova(model1)
303.  Analysis of Variance Table
304.
305.  Response: student$CC_PRO
306.                    Df  Sum Sq Mean Sq F value    Pr(>F)
307.  student$ENG_PRO    1 2379641 2379641  4747.7 < 2.2e-16 ***
308.  student$CR_PRO     1 1832117 1832117  3655.3 < 2.2e-16 ***
309.  Residuals      12408 6219179     501
310.  ---
311.  Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
312.  >
313.  > #plots
314.  > plot(model1,1)
315.  > plot(model1,2)
316.  > plot(model1,3)
317.  >
318.  > #Calculate Collinearity
319.  > vifmodel1<-car::vif(model1)
320.  > vifmodel1
321.  student$ENG_PRO  student$CR_PRO
322.         1.384697        1.384697
323.  >
324.  > #Calculate tolerance
325.  > 1/vifmodel1
326.  student$ENG_PRO  student$CR_PRO
327.        0.7221794       0.7221794
328.  >
329.  > #summary
330.  > summary(model1)
331.
332.  Call:
333.  lm(formula = student$CC_PRO ~ student$ENG_PRO + student$CR_PRO)
334.
335.  Residuals:
336.      Min      1Q  Median      3Q     Max
337.  -82.506 -13.001   4.772  15.257  68.304
338.
339.  Coefficients:
340.                   Estimate Std. Error t value Pr(>|t|)
341.  (Intercept)     10.333704   0.603445   17.12   <2e-16 ***
342.  student$ENG_PRO  0.247549   0.009276   26.69   <2e-16 ***
343.  student$CR_PRO   0.516789   0.008548   60.46   <2e-16 ***
344.  ---
345.  Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
346.
347.  Residual standard error: 22.39 on 12408 degrees of freedom
348.  Multiple R-squared:  0.4038, Adjusted R-squared:  0.4037
349.  F-statistic:  4201 on 2 and 12408 DF,  p-value: < 2.2e-16
350.
351.  >
352.  > #coefficients
353.  >   #Will allow us to isolate the beta co-efficients
354.  > lm.beta(model1)
355.
356.  Call:
357.  lm(formula = student$CC_PRO ~ student$ENG_PRO + student$CR_PRO)
358.
359.  Standardized Coefficients::
360.     (Intercept) student$ENG_PRO  student$CR_PRO
361.       0.0000000       0.2176918       0.4931652
362.
363.  >
364.  > #Tidy output of all the required stats
365.  > #For formatting outputs/tables
366.  > stargazer(model1, type="text")
367.
368.  =================================================
369.                       Dependent variable:
```

```
370.                    ---------------------------
371.                              CC_PRO
372.    -----------------------------------------------
373.    ENG_PRO                       0.248***
374.                                 (0.009)
375.
376.    CR_PRO                        0.517***
377.                                 (0.009)
378.
379.    Constant                     10.334***
380.                                 (0.603)
381.
382.    -----------------------------------------------
383.    Observations                 12,411
384.    R2                            0.404
385.    Adjusted R2                   0.404
386.    Residual Std. Error    22.388 (df = 12408)
387.    F Statistic         4,201.478*** (df = 2; 12408)
388.    ===============================================
389.    Note:                  *p<0.1; **p<0.05; ***p<0.01
390.    >
391.    > #============================================================================
392.    > #MULTIPLE REGRESSION MODEL 2
393.    > #building multiple regression model
394.    > model2<-lm(student$CC_PRO~student$ENG_PRO+student$CR_PRO+student$SCHOOL_NAT)
395.    >
396.    > #Analysis of Variance Table
397.    > anova(model2)
398.    Analysis of Variance Table
399.
400.    Response: student$CC_PRO
401.                         Df  Sum Sq Mean Sq   F value Pr(>F)
402.    student$ENG_PRO        1 2379641 2379641 4749.5335 <2e-16 ***
403.    student$CR_PRO         1 1832117 1832117 3656.7292 <2e-16 ***
404.    student$SCHOOL_NAT     1    2948    2948    5.8834 0.0153 *
405.    Residuals         12407 6216232     501
406.    ---
407.    Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
408.    >
409.    > #plots
410.    > plot(model2,1)
411.    > plot(model2,2)
412.    > plot(model2,3)
413.    >
414.    > #Calculate Collinearity
415.    > vifmodel2<-car::vif(model2)
416.    > vifmodel2
417.       student$ENG_PRO     student$CR_PRO student$SCHOOL_NAT
418.           1.542204           1.386846           1.140261
419.    >
420.    > #Calculate tolerance
421.    > 1/vifmodel2
422.       student$ENG_PRO     student$CR_PRO student$SCHOOL_NAT
423.           0.6484227          0.7210604          0.8769922
424.    >
425.    > #summary
426.    > summary(model2)
427.
428.    Call:
429.    lm(formula = student$CC_PRO ~ student$ENG_PRO + student$CR_PRO +
430.        student$SCHOOL_NAT)
431.
432.    Residuals:
433.       Min      1Q  Median      3Q     Max
434.    -82.164 -12.929   4.782  15.292  68.727
435.
436.    Coefficients:
437.                         Estimate Std. Error t value Pr(>|t|)
438.    (Intercept)          9.381324   0.719839  13.033   <2e-16 ***
439.    student$ENG_PRO      0.255136   0.009787  26.068   <2e-16 ***
```

```
440.   student$CR_PRO            0.515972   0.008553  60.329   <2e-16 ***
441.   student$SCHOOL_NATPUBLIC 1.042568   0.429822   2.426   0.0153 *
442.   ---
443.   Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
444.
445.   Residual standard error: 22.38 on 12407 degrees of freedom
446.   Multiple R-squared:  0.4041, Adjusted R-squared:  0.4039
447.   F-statistic:  2804 on 3 and 12407 DF,  p-value: < 2.2e-16
448.
449.   >
450.   > #coefficients
451.   > lm.beta(model2)
452.
453.   Call:
454.   lm(formula = student$CC_PRO ~ student$ENG_PRO + student$CR_PRO +
455.       student$SCHOOL_NAT)
456.
457.   Standardized Coefficients::
458.              (Intercept)        student$ENG_PRO        student$CR_PRO
459.               0.00000000             0.22436344            0.49238586
460.   student$SCHOOL_NATPUBLIC
461.               0.01795089
462.
463.   >
464.   > #Tidy output of all the required stats
465.   > stargazer(model2, type="text")
466.
467.   ===============================================
468.                      Dependent variable:
469.                  ---------------------------
470.                            CC_PRO
471.   -----------------------------------------------
472.   ENG_PRO                   0.255***
473.                            (0.010)
474.
475.   CR_PRO                    0.516***
476.                            (0.009)
477.
478.   SCHOOL_NATPUBLIC          1.043**
479.                            (0.430)
480.
481.   Constant                 9.381***
482.                            (0.720)
483.
484.   -----------------------------------------------
485.   Observations              12,411
486.   R2                         0.404
487.   Adjusted R2                0.404
488.   Residual Std. Error    22.384 (df = 12407)
489.   F Statistic        2,804.049*** (df = 3; 12407)
490.   ===============================================
491.   Note:                *p<0.1; **p<0.05; ***p<0.01
492.   > #==========================================================================
493.   > #MODEL COMPARISON
494.   > anova(model1,model2)
495.   Analysis of Variance Table
496.
497.   Model 1: student$CC_PRO ~ student$ENG_PRO + student$CR_PRO
498.   Model 2: student$CC_PRO ~ student$ENG_PRO + student$CR_PRO + student$SCHOOL_NAT
499.     Res.Df     RSS Df Sum of Sq      F Pr(>F)
500.   1  12408 6219179
501.   2  12407 6216232  1    2947.8 5.8834 0.0153 *
502.   ---
503.   Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
504.   > #END
505.   >
506.
```