

CAB-220

Fundamentals of Data Science

Portfolio 2

September 30, 2019

Kaspar Kielland [n10588281]

CONTENTS

1 Task 1	3
1.1 Task description	3
1.2 Code	3
1.3 Output	10
1.4 Summary	11
2 Task 2	12
2.1 Task description	12
2.2 Code	12
2.3 Output	13
2.4 Statistical test	13
2.5 Interpreting the results	14
3 Task 3	15
3.1 Task description	15
3.2 Code	15
3.3 Output	15
3.4 Interpreting the result	15
4 Task 4	16
4.1 Task description	16
4.2 Graphical analysis	16
Visualising the relationships • Checking for outliers • Check If Response Variable Is Close To Normal	
4.3 Correlation	18
4.4 The linear regression model	18

5 Task 5	20
5.1 Task description	20
5.2 Model fitting	20
5.3 Interpreting the result	21
5.4 ROC Curve	21
AUC	

PACKAGES AND LIBRARY'S

```

1 org_data <- read.csv("Portfolio\ 2\ data.csv", header = TRUE)
2
3 install.packages("ggpubr", "tidyverse", "hexbin")
4 library("ggpubr", "e1071", "corrplot", "DAAG", "ROCR", "gridExtra", "dplyr", "
  tidyverse", "hexbin")
5
6 attach(data)

```

1 TASK 1

1.1 Task description

Summarise the information in each variable (except case ID) using a table or an appropriate statistical graph.

1.2 Code

```
1 #----- Task 1 -----
2 ##### GPA #####
3 density_GPA <- data %>%
4   ggplot( aes(x=GPA)) +
5     xlab("GPA") +
6     ylab("Density") +
7     geom_density(fill="#13BFC4",
8                   color="#e9ecef",
9                   alpha=0.9) +
10    geom_vline(xintercept = mean(GPA),
11               color="#13BFC4") +
12    geom_text(aes(round(mean(GPA), 2), 0.3,
13                  label = paste("Mean: ", round(mean(GPA), 2)),
14                  hjust = 1.5),
15              size = 4,
16              color = "#13BFC4") +
17    theme_minimal()
18 density_GPA <- density_GPA +
19   ggtitle("GPA")
20
21 ##### OP score #####
22 density_OP_score <- data %>%
23   ggplot( aes(x=OP_Score)) +
24     xlab("OP score") +
25     ylab("Density") +
26     geom_density(fill="#F7766D",
27                   color="#e9ecef",
28                   alpha=0.9) +
29     geom_vline(xintercept = mean(OP_Score),
30                 color="#F7766D") +
31     geom_text(aes(round(mean(OP_Score), 2), 0.06,
32                     label = paste("Mean: ", round(mean(OP_Score), 2)),
33                     hjust = -0.2),
34               size = 4, color = "#F7766D") +
35     theme_minimal()
36 density_OP_score <- density_OP_score +
37   ggtitle("OP Score")
38
39 ##### Achived credit points #####
40 histogram_Achieved_Credit_Points <- data %>%
41   ggplot(aes(x=Achieved_Credit_Points)) +
42     xlab("Achieved credit points") +
43     ylab("Number of occurences") +
44     geom_histogram(fill="#02BC38",
45                     color="#e9ecef",
46                     alpha=0.9) +
47     geom_vline(xintercept = mean(Achieved_Credit_Points),
```

```

48         color="#02BC38") +
49     geom_text(aes(round(mean(Achieved_Credit_Points), 2), 1000,
50         label = paste("Mean: ", round(mean(Achieved_Credit_Points), 2)
51     ),
52         hjust = -0.1),
53         size = 4,
54         color = "#02BC38") +
55     theme_minimal()
56 histogram_Achieved_Credit_Points <- histogram_Achieved_Credit_Points +
57     ggtitle("Achived Credit Points")
58 ##### Failed credit points #####
59 histogram_Failed_Credit_Points <- data %>%
60     ggplot(aes(x=Failed_Credit_Points)) +
61     xlab("Failed credit points") +
62     ylab("Number of occurences") +
63     geom_histogram(fill="#F7766D",
64         color="#e9ecef",
65         alpha=0.9) +
66     geom_vline(xintercept = mean(Failed_Credit_Points),
67         color="#F7766D") +
68     geom_text(aes(round(mean(Failed_Credit_Points), 2), 1750,
69         label = paste("Mean: ", round(mean(Failed_Credit_Points), 2)),
70         hjust = -0.1),
71         size = 4, color = "#F7766D") +
72     theme_minimal()
73 histogram_Failed_Credit_Points <- histogram_Failed_Credit_Points +
74     ggtitle("Failed Credit Points")
75
76 ##### Age #####
77 density_Age <- data %>%
78     ggplot(aes(x=Age)) +
79     coord_cartesian(xlim = c(min(Age),
80         max(Age)),
81         expand = FALSE) +
82     xlab("Age") +
83     ylab("Density") +
84     geom_density(fill="#6599FF",
85         color="#e9ecef",
86         alpha=0.9) +
87     geom_vline(xintercept = mean(Age),
88         color="#6599FF") +
89     geom_text(aes(round(mean(Age), 0), 0.25,
90         label = paste("Mean: ", round(mean(Age), 0)),
91         hjust = -0.1,
92         vjust = 1.5),
93         size = 4,
94         color = "#6599FF") +
95     theme_minimal()
96 density_Age <- density_Age +
97     ggtitle("Age")
98
99 ##### Attrition Distrubiation #####
100 df_Attrition <- data %>%
101     group_by(Attrition) %>%

```

```

102 summarise(counts = n())
103 # Updates df with percentage distribution and a vector with the cumulative
    sum as element
104 df_Attrition <- df_Attrition %>%
105   arrange(desc(Attrition)) %>%
106   mutate(prop = round(counts/sum(counts), 3),
107          lab.ypos = cumsum(prop) - 0.5*prop)
108 # Creates pie chart for distribution
109 pie_chart_Attrition <- ggplot(df_Attrition,
110                               aes(x = "",
111                                   y = prop,
112                                   fill = Attrition)) +
113   geom_bar(width = 1,
114            stat = "identity",
115            color = "transparent") +
116   geom_text(aes(y = lab.ypos,
117                label = prop*100),
118            color = "white") +
119   coord_polar("y", start = 0) +
120   theme_void()
121 pie_chart_Attrition <- pie_chart_Attrition +
122   labs(fill = "Attrition") +
123   ggtitle("Attrition")
124
125 ##### Degree_Type Distribution #####
126 df_Degree_Type <- data %>%
127   group_by(Degree_Type) %>%
128   summarise(counts = n())
129 # Updates df with percentage distribution and a vector with the cumulative
    sum as element
130 df_Degree_Type <- df_Degree_Type %>%
131   arrange(desc(Degree_Type)) %>%
132   mutate(prop = round(counts/sum(counts), 3),
133          lab.ypos = cumsum(prop) - 0.5*prop)
134 # Creates pie chart for distribution
135 pie_chart_Degree_Type <- ggplot(df_Degree_Type,
136                                 aes(x = "",
137                                     y = prop,
138                                     fill = Degree_Type)) +
139   geom_bar(width = 1,
140            stat = "identity",
141            color = "transparent") +
142   geom_text(aes(y = lab.ypos,
143                label = prop*100),
144            color = "white") +
145   coord_polar("y", start = 0) +
146   theme_void()
147 pie_chart_Degree_Type <- pie_chart_Degree_Type +
148   labs(fill = "Degree type") +
149   ggtitle("Degree type")
150
151 ##### Attendance_Type Distribution #####
152 df_Attendance_Type <- data %>%
153   group_by(Attendance_Type) %>%
154   summarise(counts = n())

```

```

155 # Updates df with percantage distrubiation and a vector with the cumulative
    sum as element
156 df_Attendance_Type <- df_Attendance_Type %>%
157   arrange(desc(Attendance_Type)) %>%
158   mutate(prop = round(counts/sum(counts), 3),
159          lab.ypos = cumsum(prop) - 0.5*prop)
160
161 # Creates pie chart for gender distrubiation
162 pie_chart_Attendance_Type <- ggplot(df_Attendance_Type,
163                                   aes(x = "",
164                                       y = prop,
165                                       fill = Attendance_Type)) +
166   geom_bar(width = 1,
167            stat = "identity",
168            color = "transparent") +
169   geom_text(aes(y = lab.ypos,
170                label = prop*100),
171            color = "white") +
172   coord_polar("y", start = 0) +
173   theme_void()
174 pie_chart_Attendance_Type <- pie_chart_Attendance_Type +
175   labs(fill = "Attendance type") +
176   ggtitle("Attendance type")
177
178 ##### International_student Distrubiation #####
179 df_International_student <- data %>%
180   group_by(International_student) %>%
181   summarise(counts = n())
182 # Updates df with percantage distrubiation and a vector with the cumulative
    sum as element
183 df_International_student <- df_International_student %>%
184   arrange(desc(International_student)) %>%
185   mutate(prop = round(counts/sum(counts), 3),
186          lab.ypos = cumsum(prop) - 0.5*prop)
187 # Creates pie chart for gender distrubiation
188 pie_chart_International_student <- ggplot(df_International_student,
189                                           aes(x = "",
190                                               y = prop,
191                                               fill = International_student)) +
192   geom_bar(width = 1,
193            stat = "identity",
194            color = "transparent") +
195   geom_text(aes(y = lab.ypos,
196                label = prop*100),
197            color = "white") +
198   coord_polar("y", start = 0) +
199   theme_void()
200 pie_chart_International_student <- pie_chart_International_student +
201   labs(fill = "International student") +
202   ggtitle("International student")
203
204 ##### Gender Distrubiation #####
205 df_Gender <- data %>%
206   group_by(Gender) %>%
207   summarise(counts = n())

```

```

208 # Updates df with percentage distrubiation and a vector with the cumulative
    sum as element
209 df_Gender <- df_Gender %>%
210   arrange(desc(Gender)) %>%
211   mutate(prop = round(counts/sum(counts), 3),
212          lab.ypos = cumsum(prop) - 0.5*prop)
213 # Creates pie chart for distrubiation
214 pie_chart_Gender <- ggplot(df_Gender,
215                            aes(x = "",
216                               y = prop,
217                               fill = Gender)) +
218   geom_bar(width = 1,
219            stat = "identity",
220            color = "transparent") +
221   geom_text(aes(y = lab.ypos,
222                label = prop*100),
223            color = "white") +
224   coord_polar("y",
225              start = 0) +
226   theme_void()
227 pie_chart_Gender <- pie_chart_Gender +
228   scale_fill_discrete(name = "Gender",
229                      labels = c("Female", "Male")) +
230   ggtitle("Gender")
231
232 ##### Socio_Economic_Status Distrubiation #####
233 df_Socio_Economic_Status <- data %>%
234   group_by(Socio_Economic_Status) %>%
235   summarise(counts = n())
236 # Updates df with percentage distrubiation and a vector with the cumulative
    sum as element
237 df_Socio_Economic_Status <- df_Socio_Economic_Status %>%
238   arrange(desc(Socio_Economic_Status)) %>%
239   mutate(prop = round(counts/sum(counts), 3),
240          lab.ypos = cumsum(prop) - 0.5*prop)
241 # Creates pie chart for gender distrubiation
242 pie_chart_Socio_Economic_Status <- ggplot(df_Socio_Economic_Status,
243                                           aes(x = "",
244                                              y = prop,
245                                              fill = Socio_Economic_Status)) +
246   geom_bar(width = 1,
247            stat = "identity",
248            color = "transparent") +
249   geom_text(aes(y = lab.ypos,
250                label = prop*100),
251            color = "white") +
252   coord_polar("y", start = 0) +
253   theme_void()
254
255 pie_chart_Socio_Economic_Status <- pie_chart_Socio_Economic_Status +
256   labs(fill = "Socio economic status") +
257   ggtitle("Socio economic status")
258
259
260 ##### Teaching._Period_Admitted Distrubiation #####

```

```

261 df_Teaching._Period_Admitted <- data %>%
262   group_by(Teaching._Period_Admitted) %>%
263   summarise(counts = n())
264 # Updates df with percentage distribution and a vector with the cumulative
    sum as element
265 df_Teaching._Period_Admitted <- df_Teaching._Period_Admitted %>%
266   arrange(desc(Teaching._Period_Admitted)) %>%
267   mutate(prop = round(counts/sum(counts), 3),
268          lab.ypos = cumsum(prop) - 0.5*prop)
269 # Creates pie chart for gender distribution
270 pie_chart_Teaching._Period_Admitted <- ggplot(df_Teaching._Period_Admitted,
271                                               aes(x = "",
272                                                  y = prop,
273                                                  fill = Teaching._Period_
274
    Admitted)) +
275   geom_bar(width = 1,
276            stat = "identity",
277            color = "transparent") +
278   geom_text(aes(y = lab.ypos,
279                 label = prop*100),
280            color = "white") +
281   coord_polar("y", start = 0) +
282   theme_void()
283 pie_chart_Teaching._Period_Admitted <- pie_chart_Teaching._Period_Admitted +
284   labs(fill = "Teaching period admitted") +
285   ggtitle("Teaching period admitted")
286 ##### Faculty Distribution #####
287 df_Faculty <- data %>%
288   group_by(Faculty) %>%
289   summarise(counts = n())
290 # Updates df with percentage distribution and a vector with the cumulative
    sum as element
291 df_Faculty <- df_Faculty %>%
292   arrange(desc(Faculty)) %>%
293   mutate(prop = round(counts/sum(counts), 3),
294          lab.ypos = cumsum(prop) - 0.5*prop)
295 # Creates pie chart for distribution
296 pie_chart_Faculty <- ggplot(df_Faculty,
297                             aes(x = "",
298                                y = prop,
299                                fill = Faculty)) +
300   geom_bar(width = 1,
301            stat = "identity",
302            color = "transparent") +
303   geom_text(aes(y = lab.ypos,
304                 label = prop*100),
305            color = "white") +
306   coord_polar("y", start = 0) +
307   theme_void()
308 pie_chart_Faculty <- pie_chart_Faculty +
309   labs(fill = "Faculty") +
310   ggtitle("Faculty")
311
312

```



```

313 ##### First_in_family #####
314 df_First_in_family <- data %>%
315   group_by(First_in_family) %>%
316   summarise(counts = n())
317 # Updates df with percentage distribution and a vector with the cumulative
   sum as element
318 df_First_in_family <- df_First_in_family %>%
319   arrange(desc(First_in_family)) %>%
320   mutate(prop = round(counts/sum(counts), 3),
321          lab.ypos = cumsum(prop) - 0.5*prop)
322 # Creates pie chart for distribution
323 pie_chart_First_in_family <- ggplot(df_First_in_family,
324                                   aes(x = "",
325                                       y = prop,
326                                       fill = First_in_family)) +
327   geom_bar(width = 1,
328            stat = "identity",
329            color = "transparent") +
330   geom_text(aes(y = lab.ypos,
331                label = prop*100),
332            color = "white") +
333   coord_polar("y", start = 0) +
334   theme_void()
335 pie_chart_First_in_family <- pie_chart_First_in_family +
336   labs(fill = "First in family") +
337   ggtitle("First in family")
338
339
340 ##### All charts #####
341 combined_charts <- grid.arrange(
342   arrangeGrob(pie_chart_Gender, pie_chart_Attrition, pie_chart_Degree_Type,
343              ncol = 3),
344   arrangeGrob(pie_chart_Teaching_Period_Admitted, pie_chart_Attendance_Type,
345              pie_chart_International_student, ncol = 3),
346   arrangeGrob(pie_chart_Socio_Economic_Status, pie_chart_First_in_family, pie_
347              chart_Faculty, ncol = 3),
348   arrangeGrob(histogram_Achieved_Credit_Points, histogram_Failed_Credit_Points
349              , ncol = 2),
350   arrangeGrob(density_Age, ncol = 1),
351   arrangeGrob(density_OP_score, density_GPA, ncol = 2),
352   nrow = 6)

```

Listing 1. Code behind Figure 1

1.3 Output

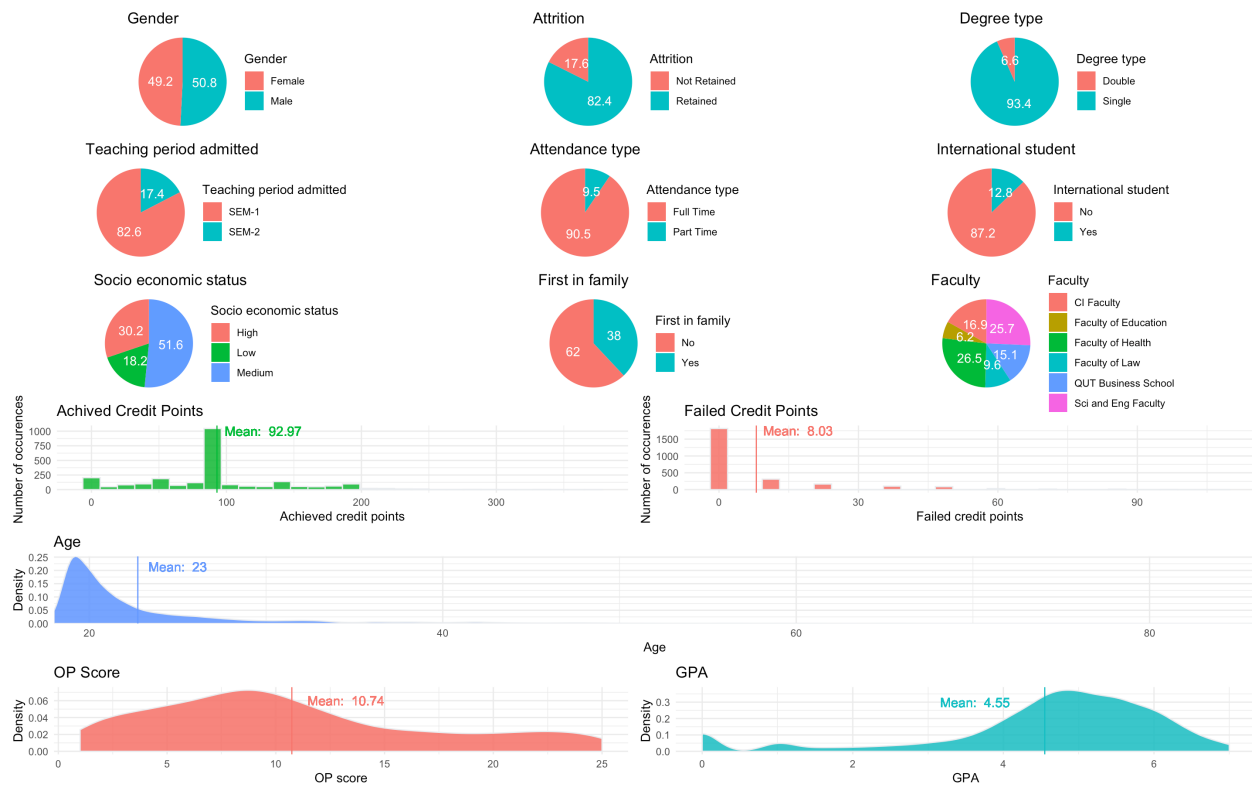


Figure 1. Data summarise

1.4 Summary

By using the following code we can get a summarised overview of our data.

```
1 summary(data)
```

The output in R gives us information about categories' frequencies of categorical variables and educational statistics, and are listed in Table 1.

Attrition		Degree type		Attendance type		International student	
Not retained	448	Double	169	Full Time	2308	No	2223
Ratained	2102	Single	2381	Part Time	242	Yes	327

First in family		Gender		Teaching period admitted		Socio economic status	
No	1580	Female	1254	SEM-1	2107	High	771
Yes	978	Male	1296	SEM-2	443	Medium	1316
						Low	463

Achieved Credit Points		Failed Credit Points		Age		GPA	
Min.	0	Min.	0	Min.	18	Min.	0
1st Qu.	60	1st Qu.	0	1st Qu.	19	1st Qu.	4.13
Median	96	Median	0	Median	20	Median	4.88
Mean	92.97	Mean	8.033	Mean	22.74	Mean	4.55
3rd Qu.	108	3rd Qu.	12	3rd Qu.	23	3rd Qu.	5.63
Max.	378	Max.	108	Max.	86	Max.	7

OP Score		Faculty	
Min.	1	CI Faculty	430
1st Qu.	6	Faculty of Education	158
Median	9	Faculty of Health	677
Mean	10.74	Faculty of Law	244
3rd Qu.	15	QUT Business School	385
Max.	25	Sci and Eng Faculty	656

Table 1. Summary of data

2 TASK 2

2.1 Task description

Compare average GPA between male and female students using a graph, conduct a statistical test, and interpret its results.

2.2 Code

```
1 #----- Task 2 -----
2 means <- aggregate(GPA ~ Gender, data, mean)
3 medians <- aggregate(GPA ~ Gender, data, median)
4
5 box_plot_GPA_Gender <- ggplot(data,
6                               aes(x=Gender,
7                                   y=GPA,
8                                   fill=Gender)) +
9   geom_boxplot(alpha=0.8,
10               fatten = 0,
11               notch = TRUE,
12               varwidth = FALSE) +
13   stat_summary(fun.y=mean,
14               geom="errorbar",
15               aes(ymax = ..y..,
16                   ymin = ..y..),
17               width = 0.75,
18               size = 1,
19               linetype = "solid") +
20   stat_summary(geom = "text",
21               label = paste("Mean: ", round(means$GPA, 2)),
22               fun.y = mean,
23               vjust=1.2) +
24   stat_summary(geom = "text",
25               label = paste("Median: ", round(medians$GPA, 2)),
26               fun.y = median,
27               vjust=-1.2) +
28   theme(legend.position="right",
29         axis.title.x = element_blank(),
30         axis.text.x = element_blank(),
31         axis.ticks.x = element_blank()) +
32   scale_fill_discrete(name = "Gender",
33                       labels = c("Female", "Male"))
34
35 box_plot_GPA_Gender +
36   ggtitle("Comparison of average GPA between\nfemale and male students",
37           subtitle = "The respective mean is represented by the black
38                       horizontal line\nThe respective median is represented by the notch") +
39   scale_y_continuous(breaks=seq(0,7,1))
```

Listing 2. Code behind the box plot chart shown by **Figure 2**

2.3 Output

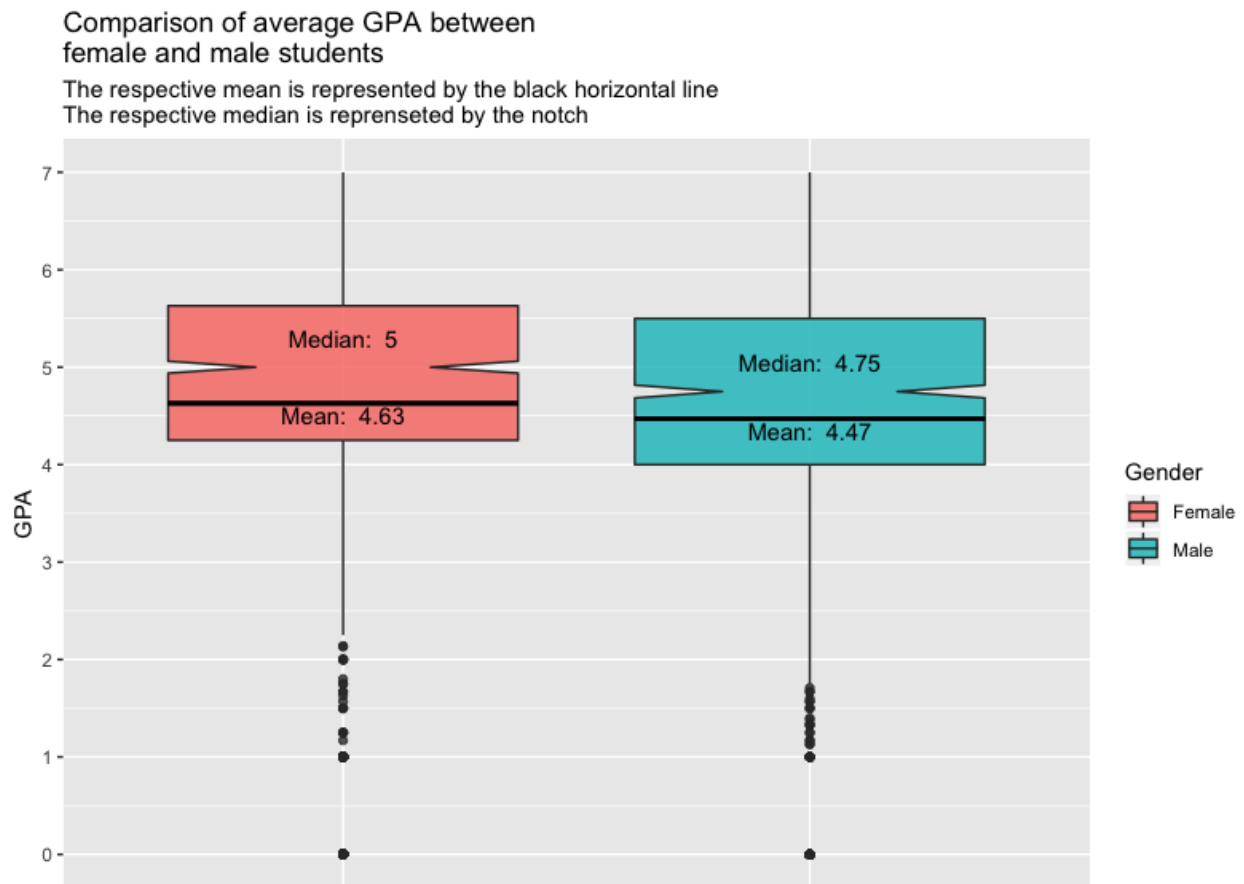


Figure 2. Comparison of GPA between genders

2.4 Statistical test

Before conducting the statistical test we need to define our null and alternative hypotheses. We first define our null hypothesis, H_0 , to be “Female students gets on average the same GPA result as male students”, and our alternative hypothesis, H_A , to be “Female students gets on average a higher GPA result then male students”. Now that we have our null hypothesis, H_0 , alternative hypothesis, H_A , we need to define our significance level, α . For this statistical test we set our significance level, α , to be 5%.

Our data to be used in this statistical test are summarized as followed:

H_0 : Female students gets on average the same GPA result as male students

H_A : Female students gets on average a higher GPA result then male students

$\alpha = 5\%$

```
1 t.test(GPA ~ Gender)
```

By conducting a Welch Two Sample t-test using the above code we get the following result in the console window:

```
1 Welch Two Sample t-test
2
3 data: GPA by Gender
4 t = 2.4454, df = 2539.7, p-value = 0.01453
5 alternative hypothesis: true difference in means is not equal to 0
6 95 percent confidence interval:
7  0.03111718 0.28297210
8 sample estimates:
9 mean in group F mean in group M
10      4.629282      4.472238
```

2.5 Interpreting the results

The result from Figure 2 shows us that female students, on average, gets a higher GPA than male students. Both the mean, median and inter-quartile range (IQR) are higher for the female students than for the male students. However, both genders have a lowest low of 0 and a highest high of 7.

The result from the statistical test gives us a p-value of 0.0145 or 1.45% which is lesser then our significance level; $\alpha > p - value \Leftrightarrow 5\% > 1.45\%$. Therefore we can conclude with saying that our alternative hypothesis, H_A , have passed and accept that female students do get, on average, a higher GPA result than male student.

3 TASK 3

3.1 Task description

Explore the relationship between OP Score and GPA using a graph, describe the relationship.

3.2 Code

```
1 #----- Task 3 -----
2 library(hexbin)
3 library(RColorBrewer)
4
5 bin<-hexbin(GPA, OP_Score, xbins=20)
6 my_colors=colorRampPalette(rev(brewer.pal(11,'Spectral'))))
7 plot(bin, style = "colorscale", legend = 1.2, lcex = 1,
8      mincnt = 1, maxcnt = 40,
9      colorcut = seq(0, 1, length = min(9, max(bin@count))),
10     border = 1,
11     colramp = my_colors,
12     xlab = "GPA", ylab = "OP Score",
13     main = "Relationship between\nOP Score and GPA")
```

Listing 3. Code behind the hexbin chart shown by **Figure 3**

3.3 Output

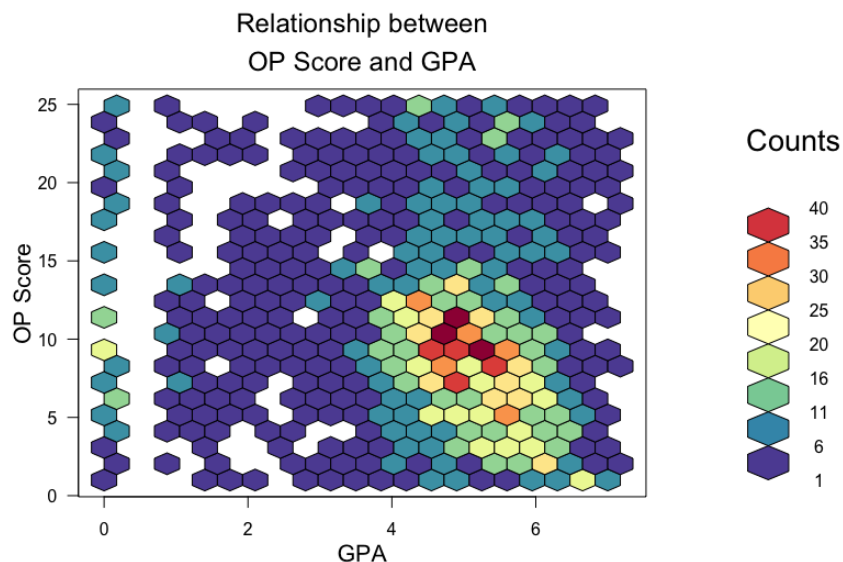


Figure 3. Hexbin chart displaying the relationship between OP Score and GPA

3.4 Interpreting the result

By studying the relationship between OP Score and GPA using the hexbin chart above, we can clearly see that the majority of students have a GPA between 4 and 6. These students have an OP Score lying between 1 and 15.

4 TASK 4

4.1 Task description

Develop a linear regression model of GPA using the given data. You need to describe your choice of predictors, examine your model's assumptions, assess model fit, and interpret the final model's regression coefficients.

4.2 Graphical analysis

This section have been conducted by following Selva Prabhakarans' "Complete Introduction to Linear Regression in R" [1].

4.2.1 Visualising the relationships

Before we can develop a linear regression model we should visualise the relationship between our response variable, GPA, and predictor variables. Since we have multiple potential predictor variables we need to draw a scatter plot along with the line of best fit for each of them.

Code

```
1 #Using Scatter Plot to visualise the relationship
2 par(mfrow = c(2, 2)) # Set up a 2 x 2 plotting space
3
4 smoothScatter(Achieved_Credit_Points, y=GPA, xlab = "Achived Credit Points",
5               main = "GPA ~ Achived Credit Points")+
6               with(data, lines(loess.smooth(Achieved_Credit_Points, GPA), col = "red"))
7
8 smoothScatter(Age, y=GPA, main = "GPA ~ Age")+
9               with(data, lines(loess.smooth(Age, GPA), col = "red"))
10
11 smoothScatter(Failed_Credit_Points, y=GPA, xlab = "Failed Credit Points", main =
12               "GPA ~ Failed Credit Points")+
13               with(data, lines(loess.smooth(Failed_Credit_Points, GPA), col = "red"))
14
15 smoothScatter(OP_Score, y=GPA, xlab = "OP Score", main = "GPA ~ OP Score")+
16               with(data, lines(loess.smooth(OP_Score, GPA), col = "red"))
```

Listing 4. Code behind the box plot chart shown by **Figure 4**

Output

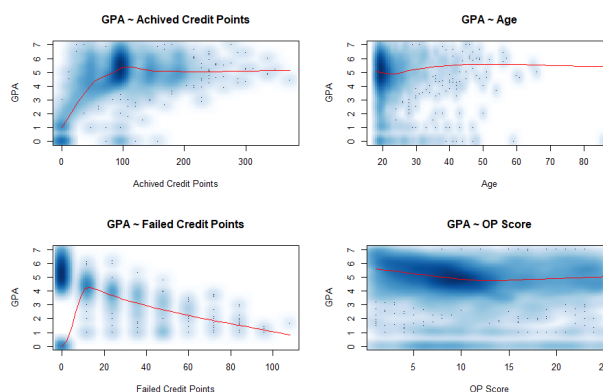


Figure 4. Smooth scatter plots for potentially useful predictors variables along with a line of best fit in red

From figure 4 we can see that our predictor variable to best predict the GPA will be the “Failed Credit Points” variable. These two variables present, overall, a negative relationship.

4.2.2 Checking for outliers

Code

```
1 # Using Boxplot to check for outliers
2 par(mfrow = c(1, 2)) # Set up a 1 x 2 plotting space
3 boxplot(Failed_Credit_Points, main = "Failed Credit Points", sub=paste("
  Outlier rows: ", boxplot.stats(Failed_Credit_Points)$out))
4 boxplot(GPA, main = "GPA", sub=paste("Outlier rows: ", boxplot.stats(GPA)$out)
  )
```

Listing 5. Code behind the box plot chart shown by **Figure 5**

Output

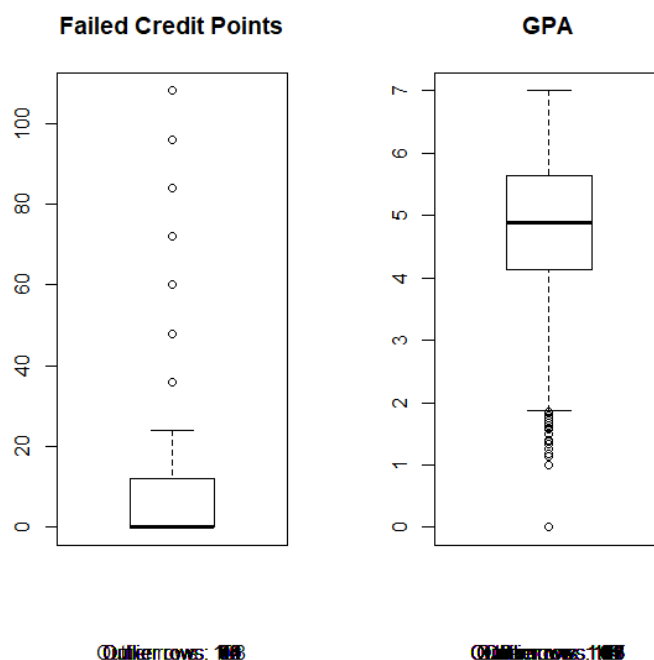


Figure 5. Boxplot of Failed Credit Points and GPA

4.2.3 Check If Response Variable Is Close To Normal

Code

```
1 # Using Density Plot to check of response variable is close to normal
2 library(e1071)
3 par(mfrow=c(1,2))
4 plot(density(GPA), main="Density Plot: GPA", ylab="Frequency", sub=paste("
  Skewness:", round(e1071::skewness(GPA), 2))
5 polygon(density(GPA), col="red")
6
7 plot(density(Failed_Credit_Points), main="Density Plot: Failed Credit Points",
  ylab="Frequency", sub=paste("Skewness:", round(e1071::skewness(Failed_
  Credit_Points), 2))
```

```
8 polygon(density(Failed_Credit_Points), col="red")
```

Listing 6. Code behind the box plot chart shown by **Figure 6**

Output

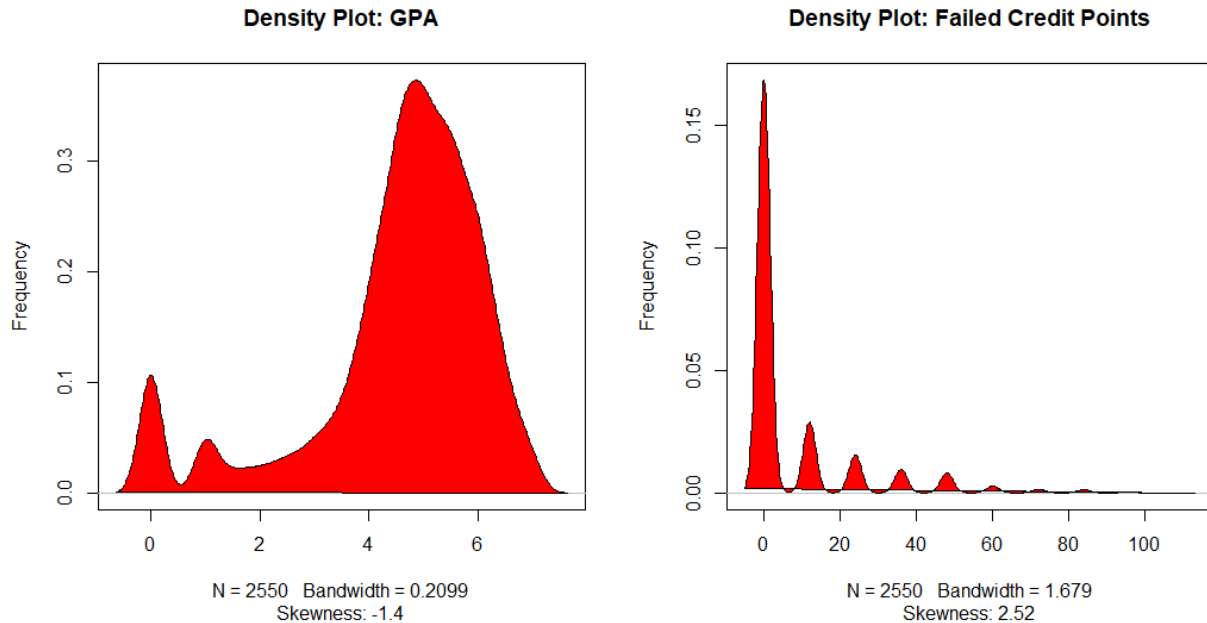


Figure 6. Density plot for GPA and Failed Credit Points

4.3 Correlation

By using the following code we find that Pearson's linear correlation to be -0.473, which is a weak negative correlation.

```
1 cor(GPA, Failed_Credit_Points)
```

4.4 The linear regression model

```
1 linearMod <- lm(GPA ~ Failed_Credit_Points)
2 summary(linearMod)
```

By executing the above code we will get the following output:

```
1 Call:
2 lm(formula = GPA ~ Failed_Credit_Points)
3
4 Residuals:
5     Min       1Q   Median       3Q      Max
6 -4.9293 -0.2593  0.2081  0.7725  2.4056
7
8 Coefficients:
9             Estimate Std. Error t value Pr(>|t|)
10 (Intercept)    4.929342   0.031576  156.11  <2e-16 ***
11 Failed_Credit_Points -0.047290   0.001743  -27.13  <2e-16 ***
12 ---
```

```

13 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
14
15 Residual standard error: 1.429 on 2548 degrees of freedom
16 Multiple R-squared:  0.2241,    Adjusted R-squared:  0.2238
17 F-statistic:    736 on 1 and 2548 DF,  p-value: < 2.2e-16

```

As we can see from the output above, both our p-values are less than the pre-determined statistical significance level of 0.05. We can therefore say that our model will be statistically significant. However, our R-squared are fairly low with only 0.224 which means that we should not rely solely on this test to estimate possible GPA outcomes.

By using the two components “Intercept” and “Failed_Credit_Points”, also known as the *beta coefficients*, we can develop a function to calculate the GPA by using Failed Credit Points as a parameter.

$$GPA = Intercept + (\beta * Failed_Credit_Points)$$

⇓

$$GPA = 4.929 + (-0.047 * Failed_Credit_Points)$$

Now that we have our linear regression line we can draw it up with our scatter plot to visualise using the following code.

```

1 plot(Failed_Credit_Points, GPA)
2 abline(4.929, -0.047)

```

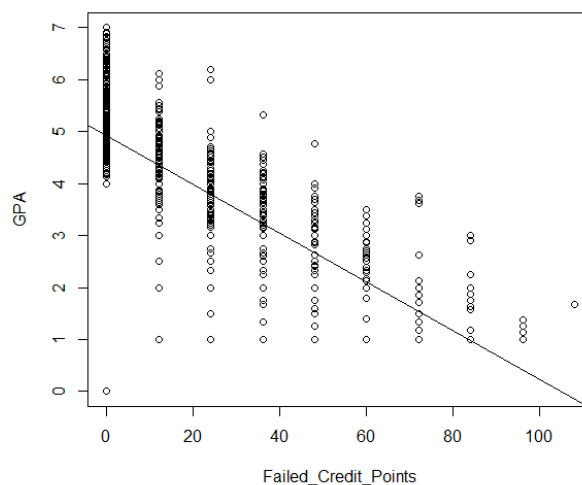


Figure 7. Scatter plot of relation between GPA and Failed Credit Points along with our linear regression line

With this visualisation we can clearly see that the GPA will be lower for students who have many failed credit points, and higher for students with few or none failed credit points.

5 TASK 5

5.1 Task description

Develop a logistic regression model to predict Attrition. You need to describe your choice of predictors, assess model fit, and interpret the final model's regression coefficients.

5.2 Model fitting

We first split our data in to two parts; Training data and Test data. Where Training data contains 80% and the remaining 20% are to be used with our Test data. This is done by executing the below code.

```
1 # Create Training data and Test data
2 set.seed(100) # setting seed to reproduce results of random sampling
3 trainingRowIndex <- sample(1:nrow(data), 0.8*nrow(data)) # row indices for
  training data
4 trainingData <- data[trainingRowIndex, ] # model training data
5 testData <- data[-trainingRowIndex, ] # test data
```

Now we need to fit our model, this can be done by executing the below code.

```
1 logistic <- glm(Attrition ~ ., data=trainingData, family=binomial(link='logit'
  '))
2 summary(logistic)
```

We have now obtained the results from our fitted model which is shown below.

```
1 Call:
2 glm(formula = Attrition ~ ., family = binomial(link = "logit"),
3     data = trainingData)
4
5 Deviance Residuals:
6     Min       1Q   Median       3Q      Max
7 -3.4995  0.2041  0.4308  0.5803  1.8828
8
9 Coefficients:
10              Estimate Std. Error z value Pr(>|z|)
11 (Intercept)    1.030e+00  5.151e-01   1.999  0.04557 *
12 ID            -1.502e-04  8.888e-05  -1.690  0.09103 .
13 Degree_TypeSingle -8.894e-01  3.469e-01  -2.564  0.01035 *
14 Achieved_Credit_Points  1.637e-02  2.027e-03   8.079 6.52e-16 ***
15 Attendance_TypePart Time  1.356e+00  2.881e-01   4.706 2.53e-06 ***
16 Age            -2.126e-02  1.101e-02  -1.930  0.05355 .
17 Failed_Credit_Points -1.614e-02  3.848e-03  -4.193 2.75e-05 ***
18 International_studentYes  7.979e-01  2.826e-01   2.824  0.00474 **
19 First_in_familyYes    7.267e-02  1.350e-01   0.538  0.59031
20 GenderM           1.670e-01  1.402e-01   1.191  0.23357
21 GPA              9.047e-02  4.752e-02   1.904  0.05692 .
22 OP_Score        -4.433e-03  1.058e-02  -0.419  0.67523
23 Socio_Economic_StatusLow -5.796e-03  1.921e-01  -0.030  0.97593
24 Socio_Economic_StatusMedium -9.001e-03  1.503e-01  -0.060  0.95223
25 Teaching._Period_AdmittedSEM-2  3.504e-01  1.881e-01   1.862  0.06254 .
26 FacultyFaculty of Education  7.109e-01  3.221e-01   2.207  0.02730 *
27 FacultyFaculty of Health    2.103e-01  1.945e-01   1.081  0.27954
28 FacultyFaculty of Law     -3.151e-01  2.528e-01  -1.246  0.21260
29 FacultyQUT Business School  5.658e-01  2.365e-01   2.392  0.01675 *
30 FacultySci and Eng Faculty  4.465e-01  2.046e-01   2.182  0.02910 *
31 ---
```

```

32 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
33
34 (Dispersion parameter for binomial family taken to be 1)
35
36 Null deviance: 1913.5  on 2039  degrees of freedom
37 Residual deviance: 1550.8  on 2020  degrees of freedom
38 AIC: 1590.8
39
40 Number of Fisher Scoring iterations: 5

```

5.3 Interpreting the result

From the results from the model fitting we can see that only the variables *Achieved Credit Points*, *Attendance Type*, *Failed Credit Points* and *International Student* are of statistically significance. Of these statistically significant variables we see that *Attendance Type* has the lowest p-value suggesting a strong association of the students attendance with the probability of being retained.

Now we can analyze the table of deviance to see how our model is doing against the null model.

```
1 anova(logistic, test="Chisq")
```

The function call above prompts us with the following output.

```

1 Analysis of Deviance Table
2
3 Model: binomial, link: logit
4
5 Response: Attrition
6
7 Terms added sequentially (first to last)
8
9
10
11          Df Deviance Resid. Df Resid. Dev  Pr(>Chi)
12 ID          1    2.268    2038    1911.3 0.1320472
13 Degree_Type 1   12.632    2037    1898.7 0.0003792 ***
14 Achieved_Credit_Points 1 262.255    2036    1636.4 < 2.2e-16 ***
15 Attendance_Type 1   23.335    2035    1613.1 1.361e-06 ***
16 Age          1    3.003    2034    1610.0 0.0830865 .
17 Failed_Credit_Points 1   18.832    2033    1591.2 1.428e-05 ***
18 International_student 1    8.816    2032    1582.4 0.0029863 **
19 First_in_family 1    0.292    2031    1582.1 0.5892548
20 Gender        1    3.839    2030    1578.3 0.0500845 .
21 GPA           1    5.485    2029    1572.8 0.0191823 *
22 OP_Score       1    0.138    2028    1572.7 0.7105264
23 Socio_Economic_Status 2    0.011    2026    1572.6 0.9945336
24 Teaching._Period_Admitted 1    3.661    2025    1569.0 0.0556895 .
25 Faculty        5   18.152    2020    1550.8 0.0027620 **
26 ---
27 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

5.4 ROC Curve

```

1 p <- predict(logistic, newdata=subset(testData), type="response")
2 pr <- prediction(p, testData$Attrition)
3 prf <- performance(pr, measure = "tpr", x.measure = "fpr")
4 plot(prf)

```

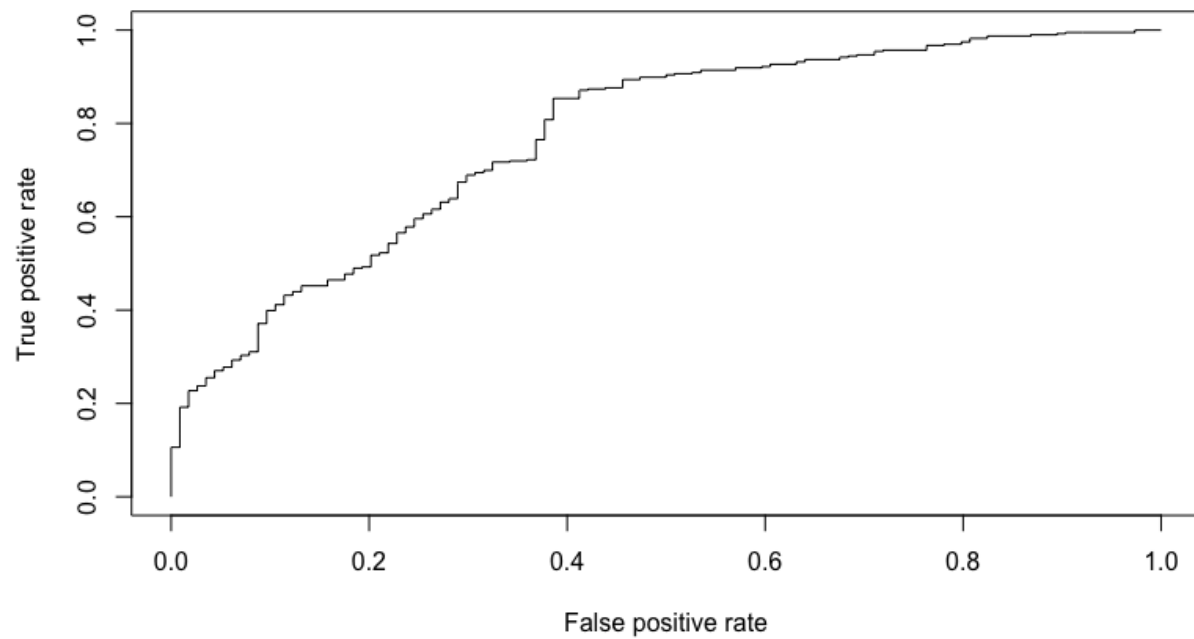


Figure 8. ROC curve for our test data

5.4.1 AUC

```
1 auc <- performance(pr, measure = "auc")
2 auc <- auc@y.values[[1]]
3 auc
```

```
1 [1] 0.7723064
```

Since our AUC are closer to 1 than to 0.5, we can say that our model have an acceptable predictive ability.

REFERENCES

- [1] S. Prabhakaran. (2017). Complete introduction to linear regression in r. Accessed: 21.08.19, [Online]. Available: <https://www.machinelearningplus.com/machine-learning/complete-introduction-linear-regression-r/>.