

Statistical Inference: Project Phase I

Sarmad Zandi Goharrizi - 810199181

Question 0

a.

Student's Performance includes various information about a sample of students studying in two different schools.

A sense of responsibility towards one's education and academic future is a notable information which can be mined from each individual's *study time* and their rate of *going out* which has an effect on their *failures* and their *grades*.

This dataset also contains some semi-relevant factors like each student's parent's job as well as their love life.

X	school	sex	age	Fjob	Mjob	goout	internet	romantic	studytime	failures	health	absences	G1	G2	G3
0	GP	F	18	teacher	at_home	4	no	no	2	0	3	6	5.000000	7.529856	9.289229
1	GP	F	17	other	at_home	3	yes	no	2	0	3	4	5.000000	7.192039	9.424835
2	GP	F	15	other	at_home	2	yes	no	2	3	3	10	3.807703	8.000000	7.354029
3	GP	F	15	services	health	2	yes	yes	3	0	5	2	15.000000	16.373208	17.796916
4	GP	F	16	other	other	2	no	no	2	0	5	4	6.000000	12.138542	12.800024
5	GP	M	16	other	services	2	yes	no	2	0	5	10	15.000000	16.804680	18.347259
6	GP	M	16	other	other	4	yes	no	2	0	3	0	12.000000	13.691091	14.187810
7	GP	F	17	teacher	other	4	no	no	2	0	1	6	6.000000	6.794185	9.012740
8	GP	M	15	other	services	2	yes	no	2	0	1	0	16.000000	19.852952	20.000000
9	GP	M	15	other	other	1	yes	no	2	0	5	0	14.000000	17.180466	18.073614
10	GP	F	15	health	teacher	3	yes	no	2	0	2	0	10.000000	9.609179	11.950918

Figure 1: Head of the dataset

b.

We have a dataset of 395 students. Each student have 16 features (some of them where mentioned in part a).

```
> summary(StudentsPerformance)
```

X	school	sex	age	Fjob	Mjob	goout	internet	romantic
Min. : 0.0	GP:349	F:208	Min. :15.0	at_home : 20	at_home : 59	Min. :1.000	no : 66	no :263
1st Qu.: 98.5	MS: 46	M:187	1st Qu.:16.0	health : 18	health : 34	1st Qu.:2.000	yes:329	yes:132
Median :197.0			Median :17.0	other :217	other :141	Median :3.000		
Mean :197.0			Mean :16.7	services:111	services:103	Mean :3.109		
3rd Qu.:295.5			3rd Qu.:18.0	teacher : 29	teacher : 58	3rd Qu.:4.000		
Max. :394.0			Max. :22.0			Max. :5.000		
studytime	failures	health	absences	G1	G2	G3		
Min. :1.000	Min. :0.0000	Min. :1.000	Min. : 0.000	Min. : 1.714	Min. : 0.000	Min. : 0.00		
1st Qu.:1.000	1st Qu.:0.0000	1st Qu.:3.000	1st Qu.: 0.000	1st Qu.: 8.000	1st Qu.: 9.988	1st Qu.:10.00		
Median :2.000	Median :0.0000	Median :4.000	Median : 4.000	Median :11.000	Median :12.244	Median :13.37		
Mean :2.035	Mean :0.3342	Mean :3.554	Mean : 5.709	Mean :10.783	Mean :12.273	Mean :12.64		
3rd Qu.:2.000	3rd Qu.:0.0000	3rd Qu.:5.000	3rd Qu.: 8.000	3rd Qu.:13.000	3rd Qu.:15.076	3rd Qu.:16.47		
Max. :4.000	Max. :3.0000	Max. :5.000	Max. :75.000	Max. :19.000	Max. :20.000	Max. :20.00		

Figure 2: Summary of the dataset

c.

As *Figure 3* and *4* suggests, there were no missing values in our dataset.

If so, there are a multitude of methods to handle missing data like, list-wise deletion, estimating them using other similar variables and ...

X	missing value	0
school	missing value	0
sex	missing value	0
age	missing value	0
Fjob	missing value	0
Mjob	missing value	0
goout	missing value	0
internet	missing value	0
romantic	missing value	0
studytime	missing value	0
failures	missing value	0
health	missing value	0
absences	missing value	0
G1	missing value	0
G2	missing value	0
G3	missing value	0

Figure 3: proportion of missing value in each feature

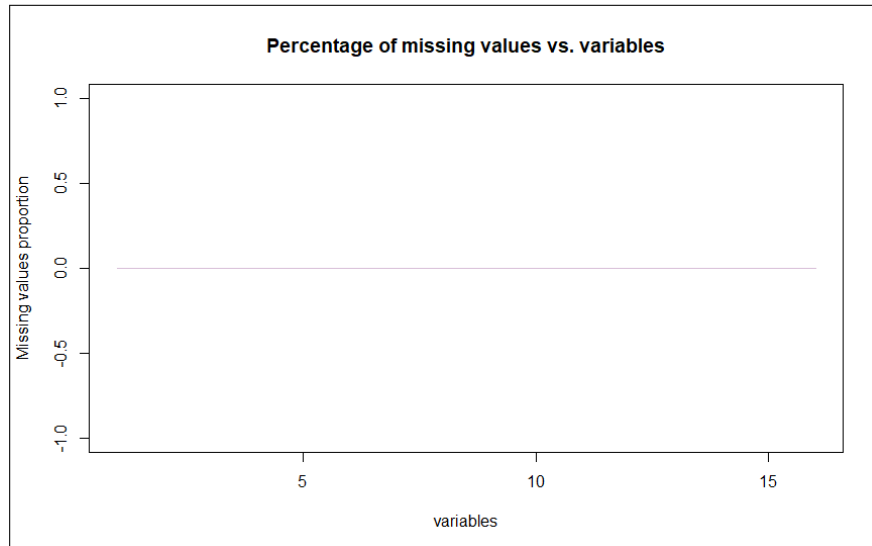


Figure 4: Line plot of missing value proportion

d.

Each student's performance is influenced highly from many different factors and cannot be decided using 3 grades, however we have to work with what we have and as was mentioned in part a, *study time* plays an important role in each individual's grades.

Question 1

Chosen Numerical Variable : *G1*

a.

The appropriate bin width is computed using *Freedman–Diaconis rule* , which leads to a normally distributed histogram of *G1*.

$$\text{Bin Width} : 2 \frac{IQR(x)}{\sqrt[3]{n}}$$

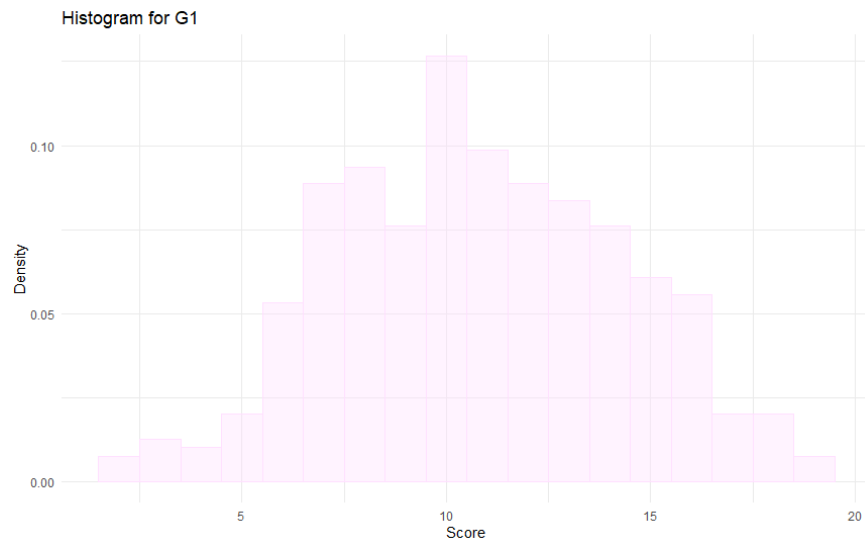


Figure 5: Histogram of *G1*

Figure 6 describes a *unimodal*.

A *unimodal* distribution is a distribution that has one clear peak (as can be seen in *Figure 6*). The values increase at first, rising to a single highest point where they then start to decrease. A *unimodal* distribution can either be symmetrical or non-symmetrical (more about this in part c).

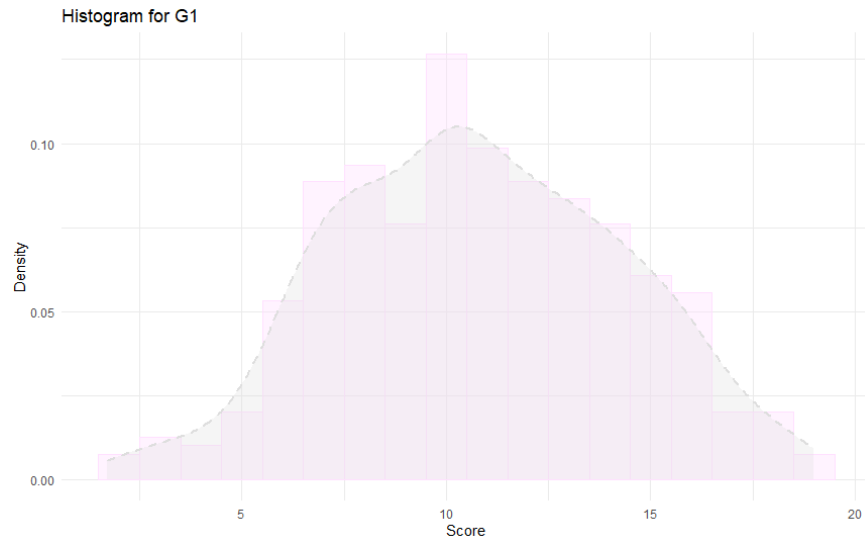


Figure 6: Histogram of G1 overlaid with density plot

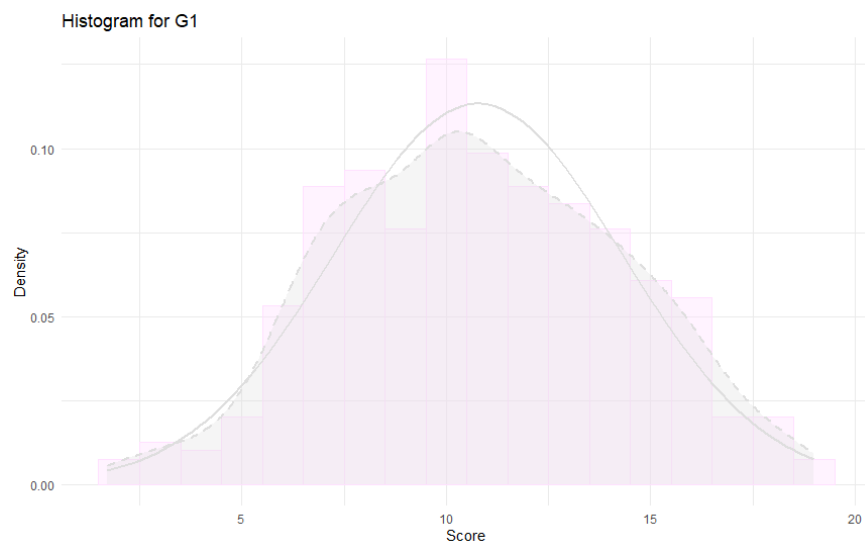


Figure 7: Histogram of G1 overlaid with fitted density plot and MLE density plot

b.

There are 3 basic properties of a distribution that we have to address: *location*, *spread*, and *shape*.

The *location* refers to the typical value of the distribution, such as the *mean* (10.783) or *median* (11.00).

The *spread* of the distribution is the amount by which smaller values differ from larger ones. The *standard deviation* (3.521) or *variance* (12.39) are measures of distribution spread.

The *shape* of a distribution is its pattern—peakedness, symmetry, etc. A given phenomenon may have any one of a number of distribution shapes, e.g., the distribution may be bell-shaped, rectangular-shaped,

etc which in our case is nearly *bell-shaped symmetrical (unimodal)* as was mentioned in part a and will be discussed in part c.

```
> summary(StudentsPerformance)
  X      school sex      age      math score      reading score      writing score
Min.   : 0.0    GP:349  F:208 Min.   :15.0    1st Qu.:16.0    1st Qu.:16.0    1st Qu.:16.0
1st Qu.: 98.5    MS: 46  M:187 1st Qu.:16.0    1st Qu.:16.0    1st Qu.:16.0
Median :197.0    Mean :16.7    Mean :16.7    Mean :16.7
Mean   :197.0    3rd Qu.:18.0    3rd Qu.:18.0    3rd Qu.:18.0
3rd Qu.:295.5    Max.   :22.0    Max.   :22.0    Max.   :22.0
Max.   :394.0

  studytime failures health
Min.   :1.000 Min.   :0.0000 Min.   :1.00
1st Qu.:1.000 1st Qu.:0.0000 1st Qu.:13.00
Median :2.000 Median :0.0000 Median :4.00
Mean   :2.035 Mean   :0.3342 Mean :13.55
3rd Qu.:2.000 3rd Qu.:0.0000 3rd Qu.:15.00
Max.   :4.000 Max.   :3.0000 Max.   :5.00

  G1
Min.   : 1.714
1st Qu.: 8.000
Median :11.000
Mean   :10.783
3rd Qu.:13.000
Max.   :19.000
```

Figure 2

of the dataset

Figure 3 and 4 suggests, there were no missing values. Also, there are a multitude of methods to handle missing data like, list-wise deletion, estimating them using

Figure 8: G1 under magnifier

It can be clearly seen that this distribution is very similar to the normal distribution but to be more precise, we use *normal Q-Q plot*.

The main purpose of a *normal probability plot (normal Q-Q plot)* is to assess normality.

A one-to-one relationship (straight line in *Figure 8*) between the data and the theoretical quantiles can be considered, so the data follow a nearly normal distribution. In other words, the closer the points to the straight line, the more confident we can be that the data follow the normal model.)

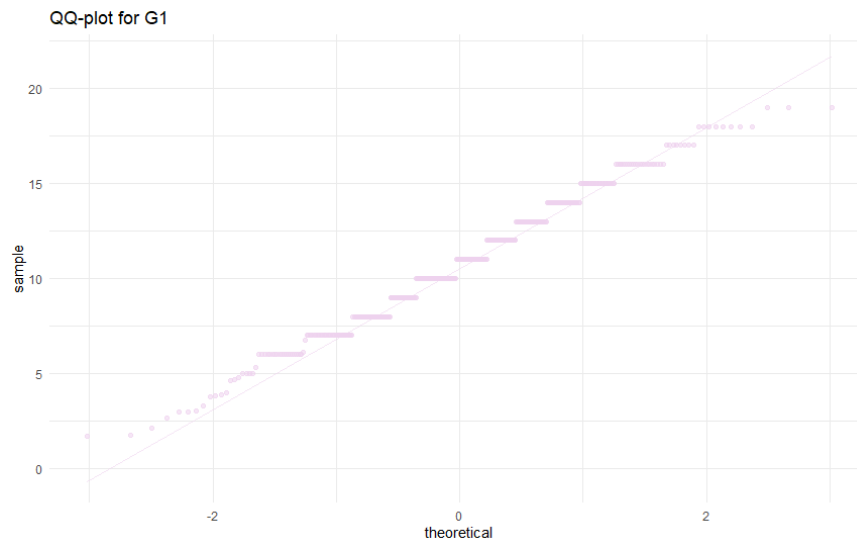


Figure 9: Normal Q-Q plot of G1

c.

Skewness is a statistical numerical method to measure the asymmetry of the distribution or data set. It tells about the position of the majority of data values in the distribution around the mean value.

$$Skewness = \frac{mean - median}{sd}$$

One method to address the skewness is to compare the mean and the median.

If :

1. $mean > median$: right skewed (negatively skewed)

2. $mean = median$: Symmetric

3. $mean < median$: left skewed (positively skewed)

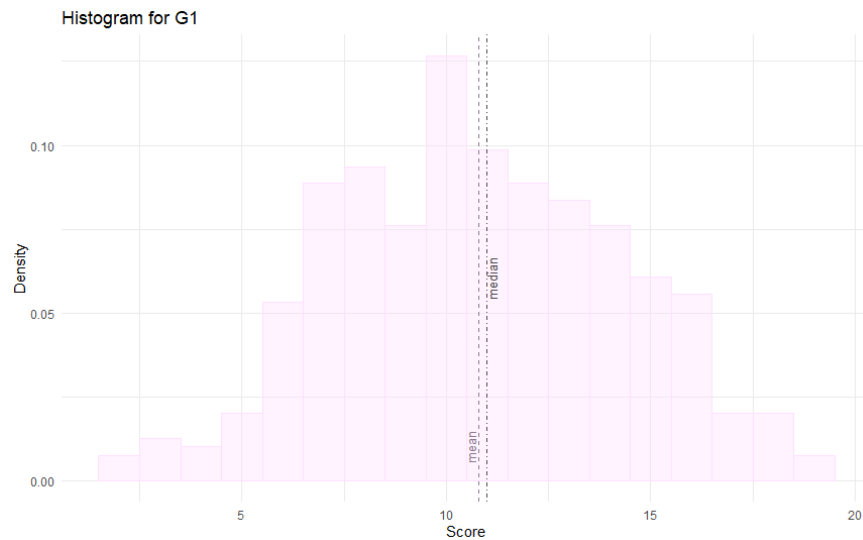


Figure 10: Median and mean marked on histogram of G1

As can be deduced from *Figure 10*, G1 (barely) falls under the third category. This conclusion can also be supported by calculating the skewness of G1 :

```
>
> skewness(StudentsPerformance$G1)
[1] 0.01764784
>
```

Figure 11: Calculated skewness of G1

The coefficient of skewness is greater than 0, meaning the graph is positively skewed with the majority of data values less than mean. In other words, most of the values are concentrated on the left side of the graph.

d.

An outlier is a value or an observation that is distant from other observations, that is to say, a data point that differs significantly from other data points.

Boxplots provide a useful visualization of the distribution of data. Typically, Boxplots show the *median*, *1st quartile*, *3rd quartile*, *maximum datapoint*, and *minimum datapoint* for a dataset (more to it in part h) and also, last but not least, *outliers*. Fortunately, my chosen variable didn't have any outliers and the *Figures 12* and *13* below are the proof.

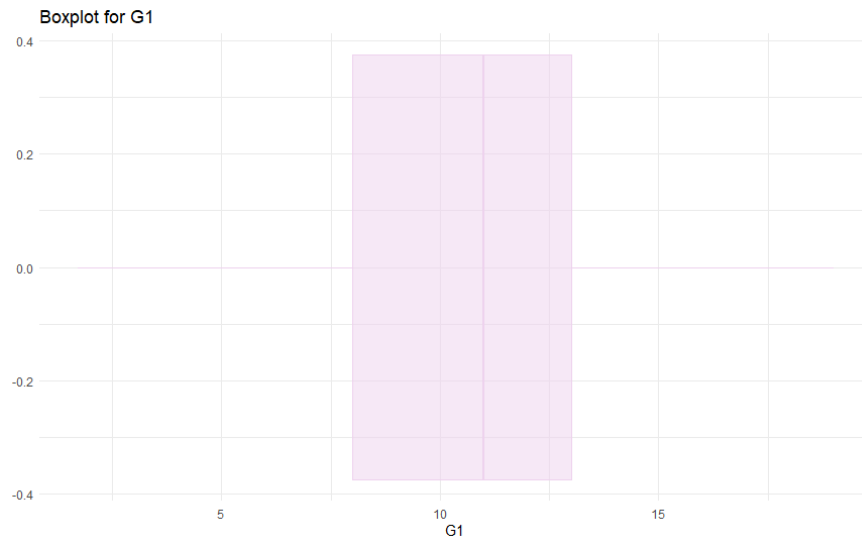


Figure 12: Boxplot of G1 to visualize outliers

```
>
> boxplot.stats(StudentsPerformance$G1)$out
numeric(0)
>
```

Figure 13: Using stats of boxplot to visualize outliers

e.

Mean : The mean identifies the average value of the set of numbers.

$$\mu = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

Median : The median identifies the midpoint or middle value of a set of numbers.

Variance : Variance measures the variability of the data set. It indicate how far individuals in the group are spread out, in the set of data from the mean.

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

Standard deviation : Standard deviation measures the dispersion of the data set. A smaller standard deviation indicates less variability. Standard deviation is expressed in the same unit as the values in the dataset so it measures how much observations of the data set differs from its mean.

$$\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2}$$

```
>
> mean(StudentsPerformance$G1)
[1] 10.78285
> median(StudentsPerformance$G1)
[1] 11
> var(StudentsPerformance$G1)
[1] 12.39784
> sd(StudentsPerformance$G1)
[1] 3.521057
>
```

Figure 14: Statistics: Mean-Median-Variance-Standard Deviation

f.

The perfect description of the relationship between *mean*, *median* and *density* is that the *median* of a density curve is the point that divides the area under the curve in half, the *mean* is the point at which the curve would balance if made out of solid material.

In a perfectly symmetrical distribution, the mean and the median are the same.

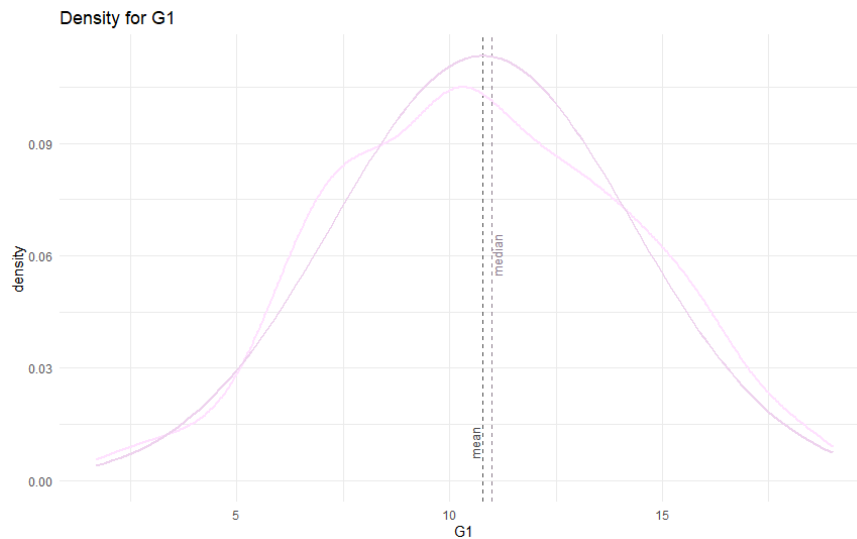


Figure 15: Median and mean marked on density of G1 - darker one is drawn using *dnorm*

g.

Pie charts are best to use when you are trying to compare parts of a whole. For this question, two different courses of action were taken :

First Method : Categorizing data by a range of values

In this approach categories are created according to logical cut-off values in the scores or measured values.

$$\left\{ \begin{array}{ll} G1 < \frac{\mu}{2} & \text{Very Low} \\ \frac{\mu}{2} < G1 < \mu & \text{Low} \\ \mu < G1 < \frac{\mu + \max(G1)}{2} & \text{High} \\ G1 > \frac{\mu + \max(G1)}{2} & \text{Very High} \end{array} \right.$$

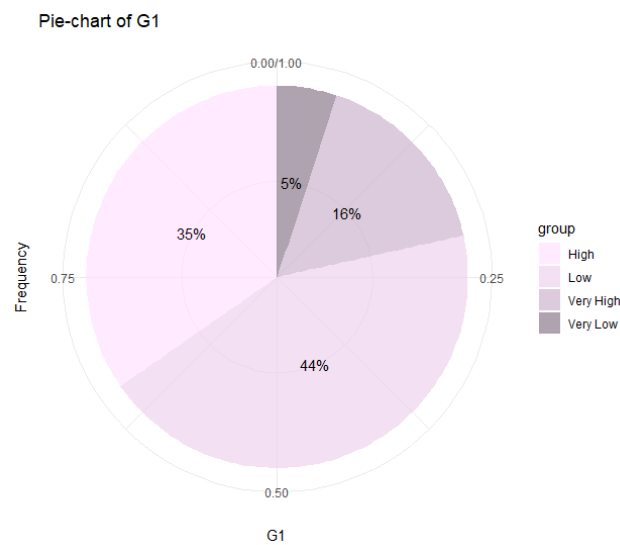
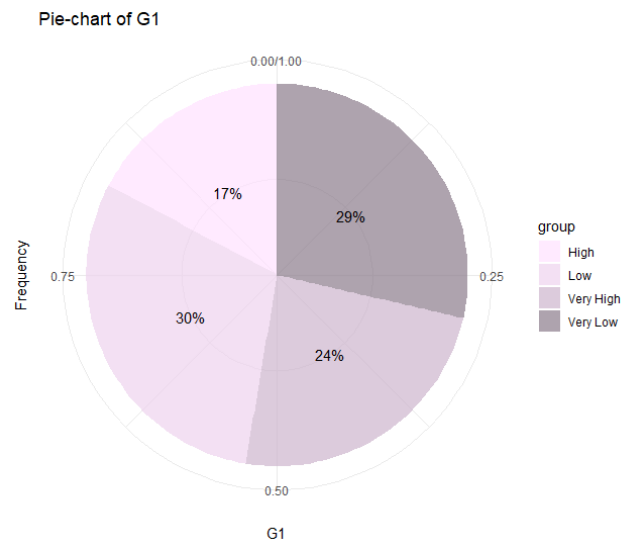


Figure 16: Piechart of G1 - 1st method

Second Method : Categorizing data by percentiles (since mean and median are close)

A second approach is to use percentiles to categorize data. The advantage to this approach is that it does not rely on the scoring system being meaningful in its absolute values

$$\left\{ \begin{array}{ll} G1 < 25^{th} percentile & \text{Very Low} \\ 25^{th} percentile < G1 < 50^{th} percentile & \text{Low} \\ 50^{th} percentile < G1 < 75^{th} percentile & \text{High} \\ G1 > 75^{th} percentile & \text{Very High} \end{array} \right.$$

Figure 17: Piechart of G1 - 2nd method

In this approach, there are approximately an equal number of respondents in each category.

h.

```
>
> G1.quant
      0%      25%      50%      75%     100%
1.713843  8.000000 11.000000 13.000000 19.000000
>
>
```

Figure 18: 0th, 25th, 50th, 75th, and 100th percentiles of G1

```
>
> G1.iqr
[1] 5
>
>
```

Figure 19: IQR of G1

Box plots are a five-number summary that includes the minimum and maximum data values, the median and lower and upper quartiles. They can be useful in understanding how is data distributed in a given set and give information about the spread of the data.

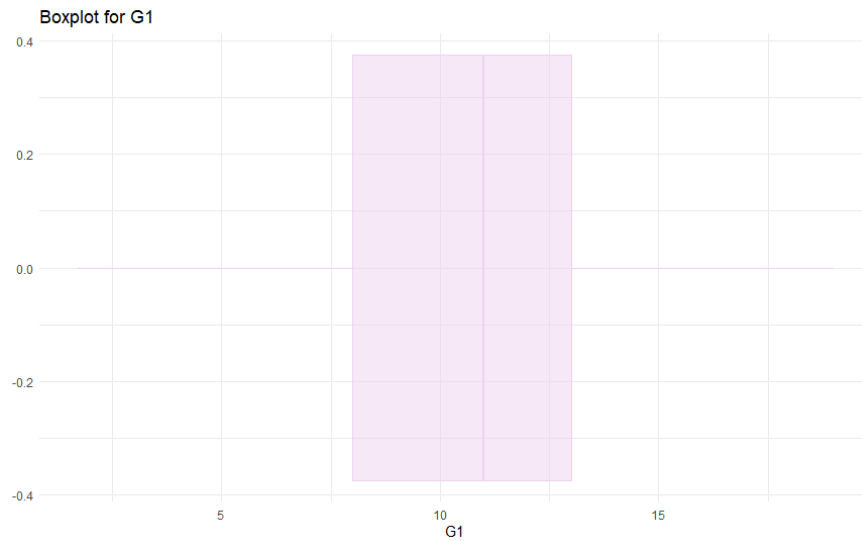


Figure 20: Boxplot of G1

```
> boxplot.stats(StudentsPerformance$G1)
$stats
[1]  1.713843  8.000000 11.000000 13.000000 19.000000

$n
[1] 395

$conf
[1] 10.60251 11.39749

$out
numeric(0)
```

Figure 21: Stats of Boxplot of G1

From *Figure 20*, G1 being (barely) LS is also clear.

Question 2

Chosen Categorical Variable : *sex*

a.

Most of them are female students.

```
> female.freq  
[1] 0.5265823  
> male.freq  
[1] 0.4734177  
>
```

Figure 22: Frequency of each category and its percentage

b.

A stacked barplots is a variant of the bar chart.

A standard barplots compares individual data points with each other. In a stacked barplots, parts of the data are adjacent (in the case of horizontal bars) or stacked (in the case of vertical bars); each bar displays a total amount, broken down into sub amounts.

Stacked barplots are useful for visualizing conditional frequency distributions.(But in general, it is better to avoid them.)

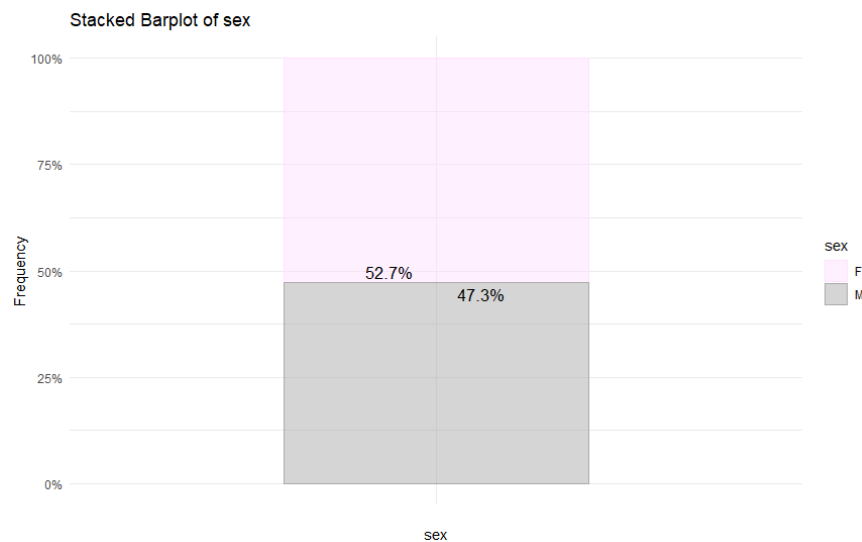


Figure 23: Stacked barplot of sex

c.

Barplots for categorical variables are like histograms for numerical variables.

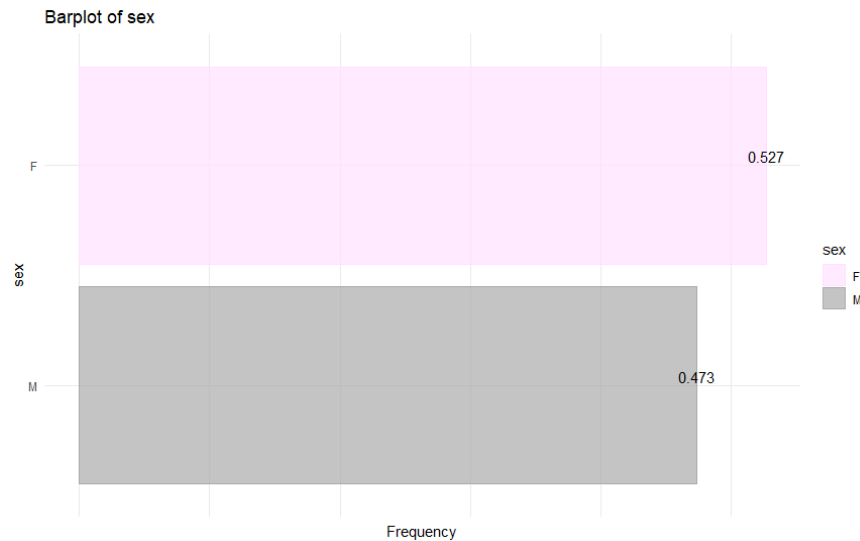


Figure 24: Horizontal barplot of sex

d.

A violinplot is a method of plotting numeric data. It is similar to a boxplot, with the addition of a rotated kernel density plot on each side.

A violinplot is more informative than a plain boxplot. While a boxplot only shows summary statistics such as mean/median and inter-quartile ranges, the violin plot shows the full distribution of the data. Wider sections of the violin plot represent a higher probability that members of the population will take on the given value; the skinnier sections represent a lower probability.

Violin plots are used to represent comparison of a variable distribution (or sample distribution) across different "categories".

In our case, *Female* students are around 16 to 18 years old and the distribution of *Male* is wider than *Female* and continues until the age of 22 years.

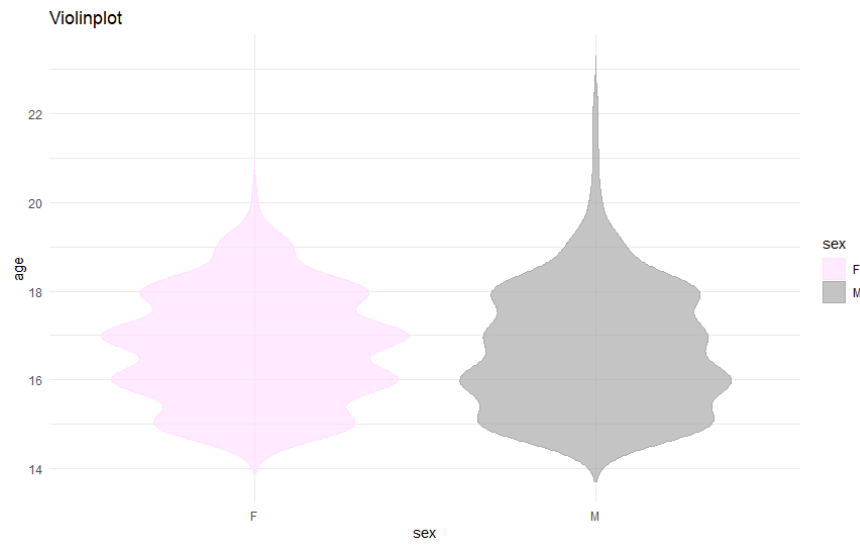


Figure 25: Violin plot of sex

Question 3

Chosen Numerical Variables : *goout* and *absences*

a.

The data points might follow an overall positive trend, the more you go out, the less you can show up to class.

My guess is a positive non-linear relationship between these two.

b.

A clear relationship cannot be described. It seems like a bell-shaped relationship, also an outlier in *goout* = 1 is detected.

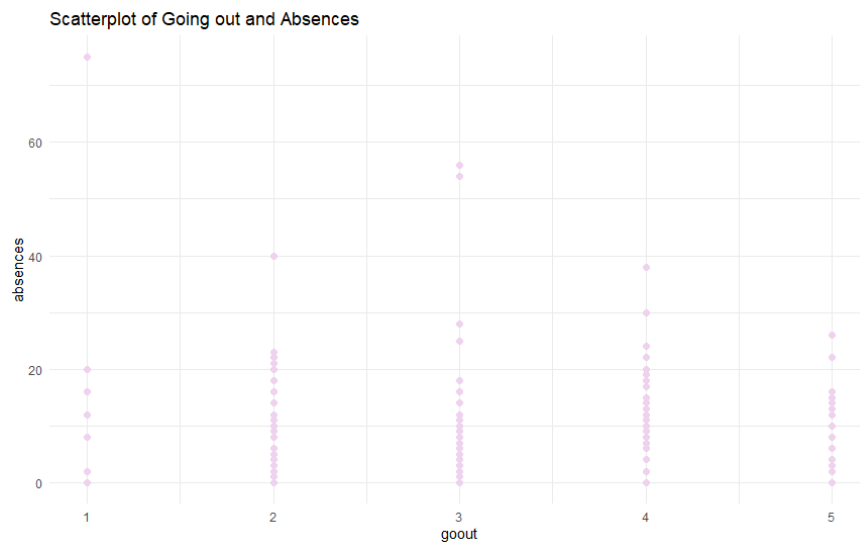


Figure 26: Scatterplot of goout and absences

c.

A correlation coefficient is a numerical measure of some type of correlation, meaning a statistical relationship between two variables.

Correlation is computed using *Pearson correlation coefficient*.

Pearson's correlation coefficient, when applied to a sample, is commonly represented by r_{xy} and may be referred to as the sample correlation coefficient or the sample Pearson correlation coefficient.

$$r_{xy} = \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sqrt{\sum x_i^2 - n \bar{x}^2} \sqrt{\sum y_i^2 - n \bar{y}^2}}$$


```
>
> goout_absences.correlation
[1] 0.04430222
>
>
```

Figure 27: Correlation coefficient of goout and absences

d.

The correlation coefficient ranges from -1 to 1 . A value of 1 implies that a *linear equation* describes the relationship between X and Y perfectly (a.k.a perfect positive correlation), with all data points lying on a line for which Y increases as X increases. A value of -1 implies that all data points lie on a line for which Y decreases as X increases (a.k.a perfect negative correlation). A value of 0 implies that there is no linear correlation between the variables.

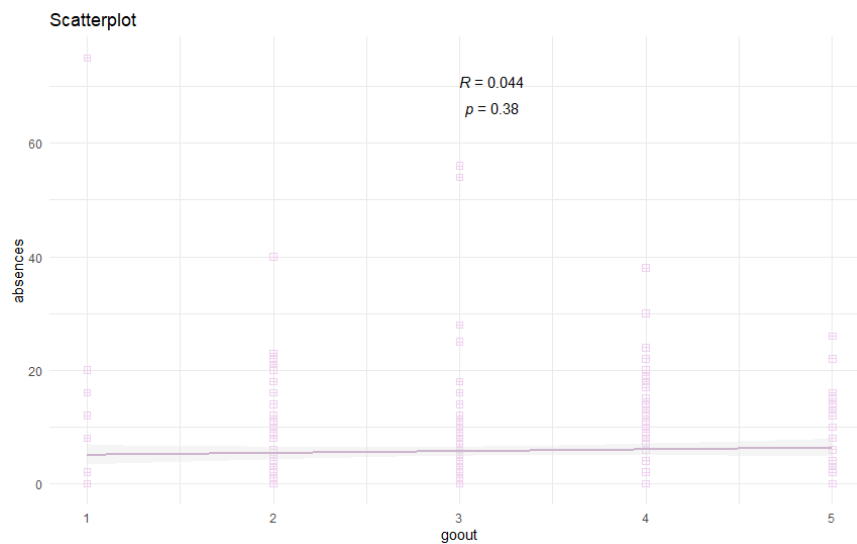


Figure 28: Scatterplot of goout and absences

In our case, $R = 0.044$ means no or negligible (positive) relationship. (So the assumption made in part a was somewhat true.)

e.

Statistical inference based on *Pearson's correlation coefficient* often focuses on one of the following two aims:

- One aim is to test the null hypothesis that the true correlation coefficient ρ is equal to 0, based on the value of the sample correlation coefficient r .
- The other aim is to derive a confidence interval that, on repeated sampling, has a given probability of containing ρ .

In this part, the first aim is our target. A p-value is the probability that the null hypothesis is true. When using *Pearson's correlation coefficient*, it represents the probability that the *correlation* between x and y in the sample data occurred by chance.

In our case, ρ a.k.a p -value is 0.38.

A p -value of 0.38 means that there is 38% chance (!) that results from the sample occurred due to chance. Comparing to significant level of 5%, we fail to reject the null hypothesis.

We conclude that the correlation is not statically significant. Or in other words *we conclude that there is not a significant linear correlation between x and y in the population whatsoever.*

f.

Chosen Categorical Variable : *romantic*

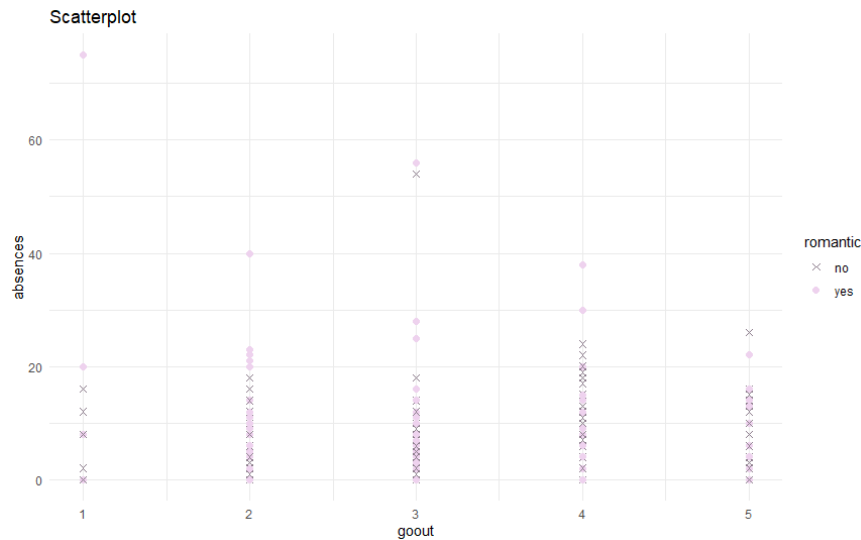


Figure 29: Scatterplot of goout and absences categorized by romantic

g.

Hexbin map uses hexagons to split the area into several parts and attribute a color to it. The graphic area is divided into a multitude of hexagons and the number of data points in each is counted and represented using a *color gradient*.

Hexbin plot is helpful in situations where :

- Creating an unbiased density distribution is needed
- Representing discrete categorical information is needed (Better than heatmaps in visualizing categorical information)
- Showing complete information by eliminating the edge effects is needed (Circle is the lowest ratio, but cannot form a continuous grid, and hexagons are the closest shape to a circle that can still form a grid.)

Hexbin plot should be avoided in situations where simplicity of definition and data storage is needed.

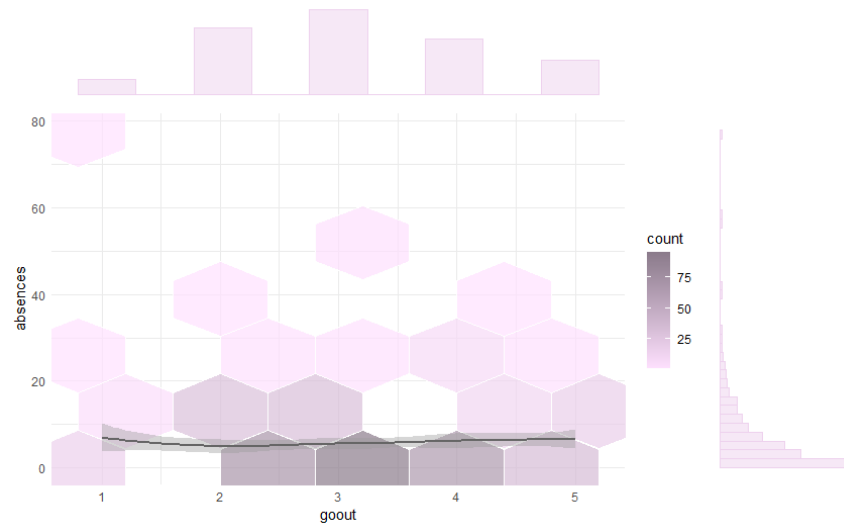


Figure 30: Hexbin plot, binsize = 5

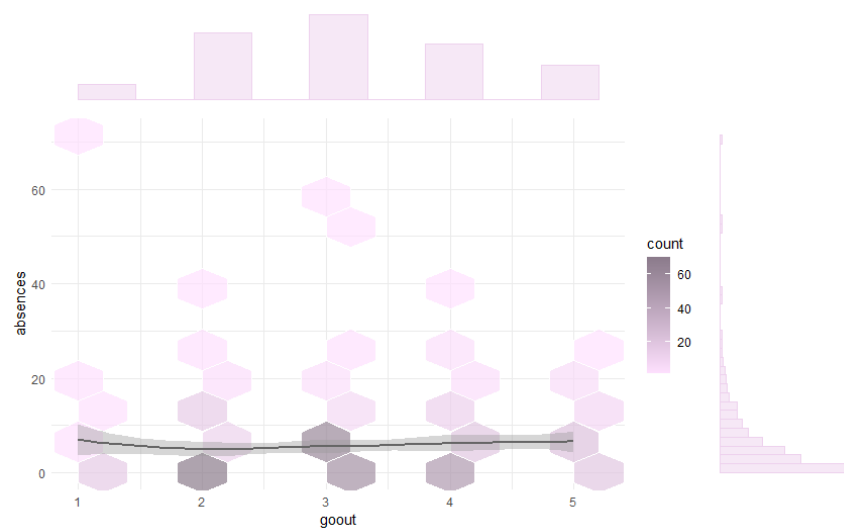


Figure 31: Hexbin plot, binsize = 10

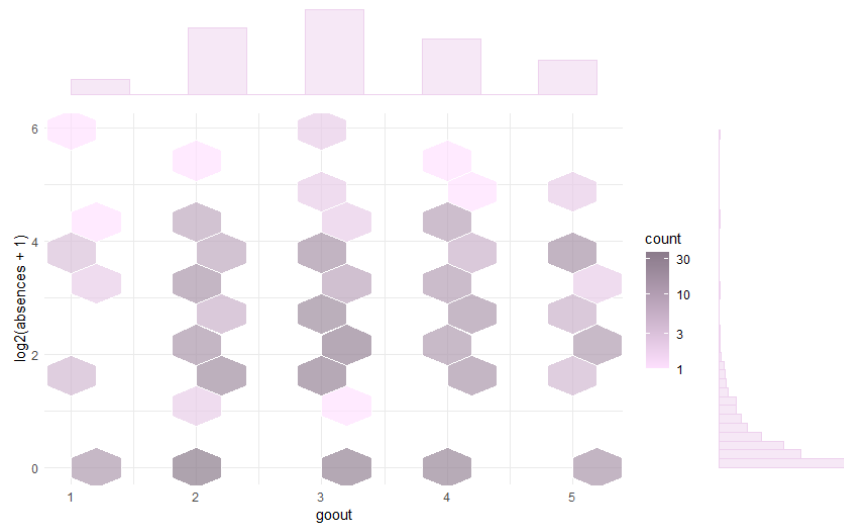


Figure 32: Hexbin plot, binsize = 10 (logarithmic)

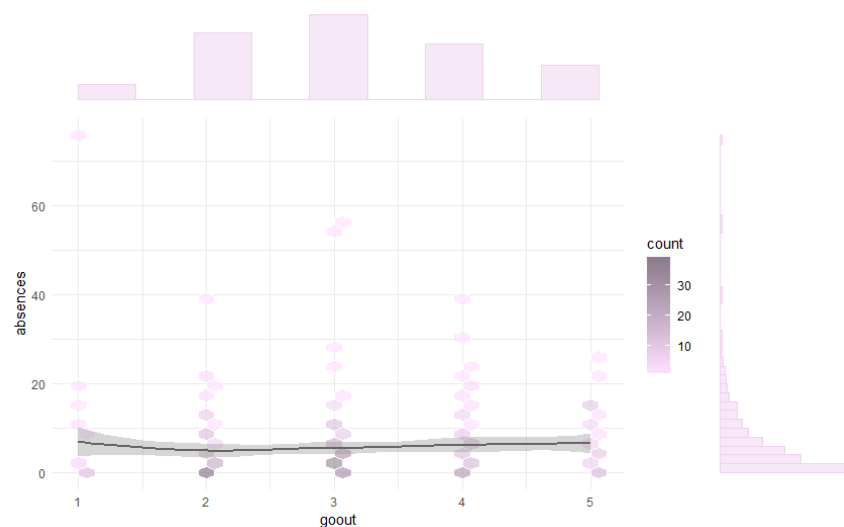


Figure 33: Hexbin plot, binsize = 30

It can be seen that by decreasing the binsize, each hexagon contains more amount of samples. Binsize about 10 is fairly good and can be informative. Bigger Binsizes will be misleading and not robust to noisy datas. Logarithmic plot was also plotted to have a better visualization.

h.

A 2D density plot displays the relationship between 2 numeric variables, where one variable is represented on the X-axis, the other on the Y axis. The number of observations within a particular area of the 2D space is counted and represented by a *color gradient* to indicate differences in the distribution of data in one region with respect to the other.

2D density plot is helpful in situations where :

- Sample size is huge and a clearer picture of the distribution is needed
- A nuanced visualization of density is needed (Better than heatmaps in visualizing categorical information)
- Visualize several distributions at once is needed

2D density plot should be avoided in situations where not enough data points are present, therefore risk of overplotting is low (using scatterplot is a more effective visualization).

The biggest disadvantage of 2D density plots and Hexbin maps are their sensitivity to bin size/bandwidth, inaccurate bin size/bandwidth and can lead to different and/or wrong conclusions.

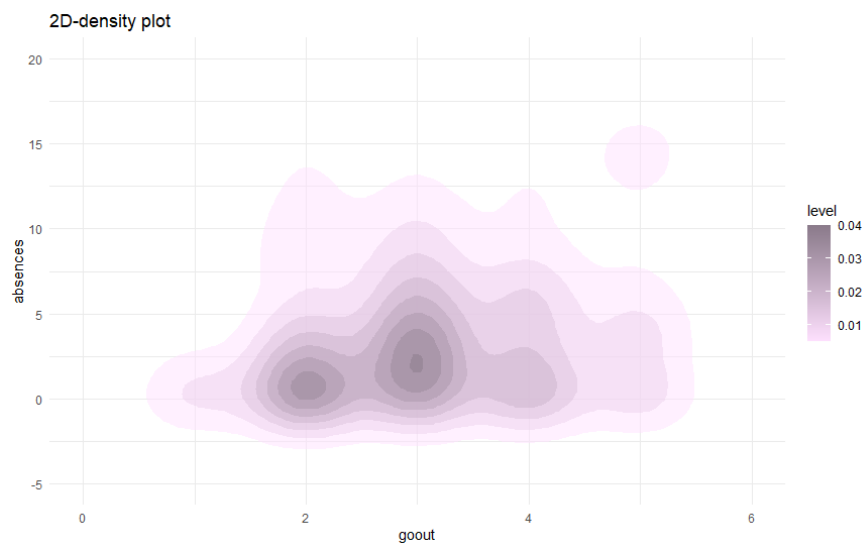


Figure 34: 2D density plot of goout and absences

As can be concluded from *Figure 30*, the densest part of the plot is when students goout 3 times and are absent for 5 times.

Question 4

a.

Scatterplots of each pair of numeric variable are drawn on the left part of the figure. Pearson correlation is displayed on the right. Variable distribution is available on the diagonal.

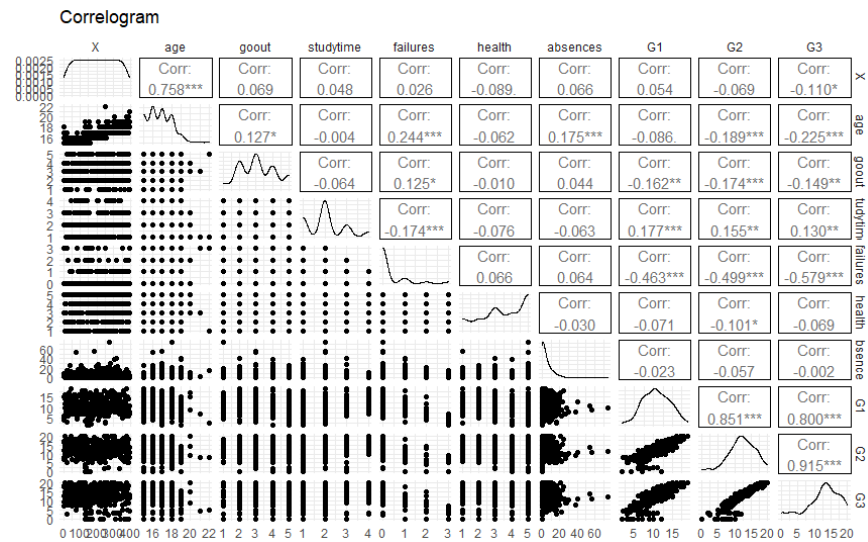


Figure 35: Bivariate Correlogram with Pearson correlation

Density's bandwidth of *Failures* variable was inf, so we had to omit it in order to get a plot:

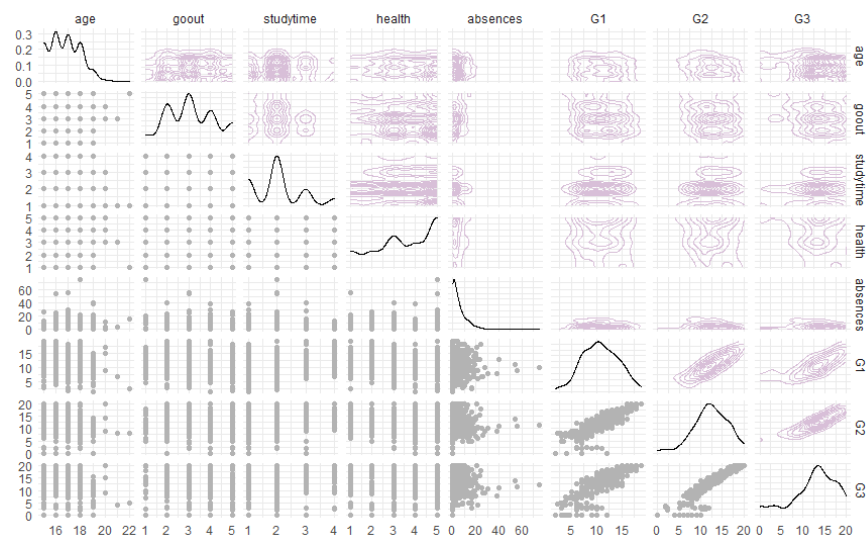


Figure 36: Bivariate Correlogram with density - scatterplot

Judging by *Figure 36*, where the scatterplot of 2 variables is dense, density plot is completely meaningful and where the scatterplot of 2 variables is not dense, density plot is not that informative and it's better to stick to scatterplot as was mentioned in part h of question3, *2D density plot should be avoided in situations where not enough data points are present, therefore risk of over-plotting is low (using scatterplot is a more*

effective visualization)

To have the full view of all of our numerical variables, boxplot, barplot and scatterplot with linear association was also plotted :

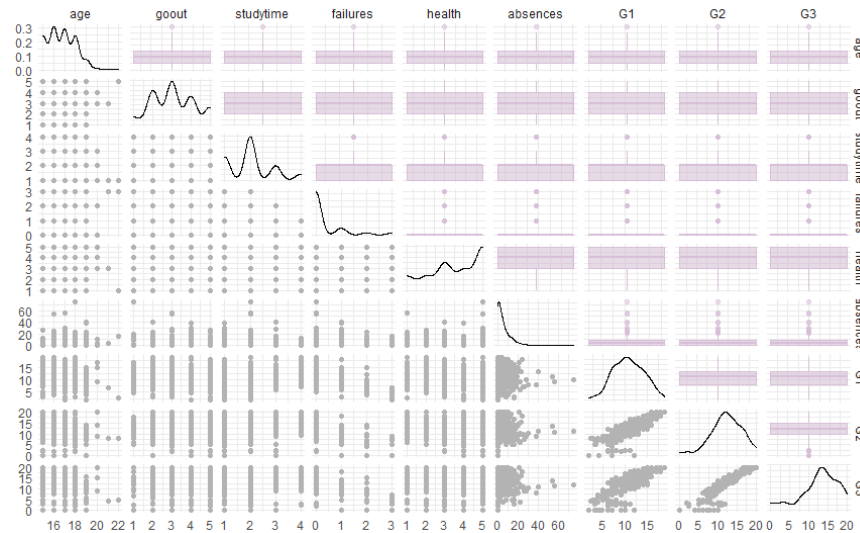


Figure 37: Bivariate Correlogram with barplot - scatterplot

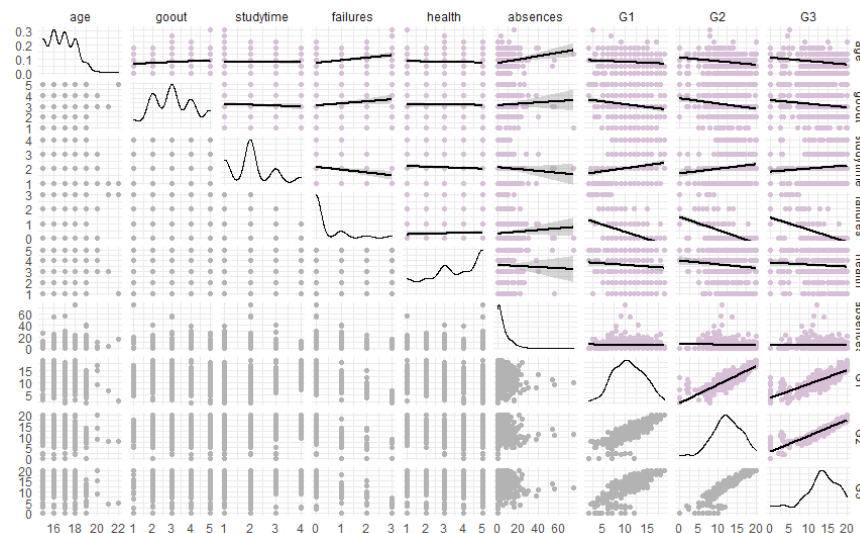


Figure 38: Bivariate Correlogram with linear association - scatterplot

Judging by *Figure 38*, $G1$ and $G2$ and $G3$ have positive linear associations with each other and with *studytime* as expected. *Failure* and *goout* both have a negative linear associations with $G1$, $G2$ and $G3$.

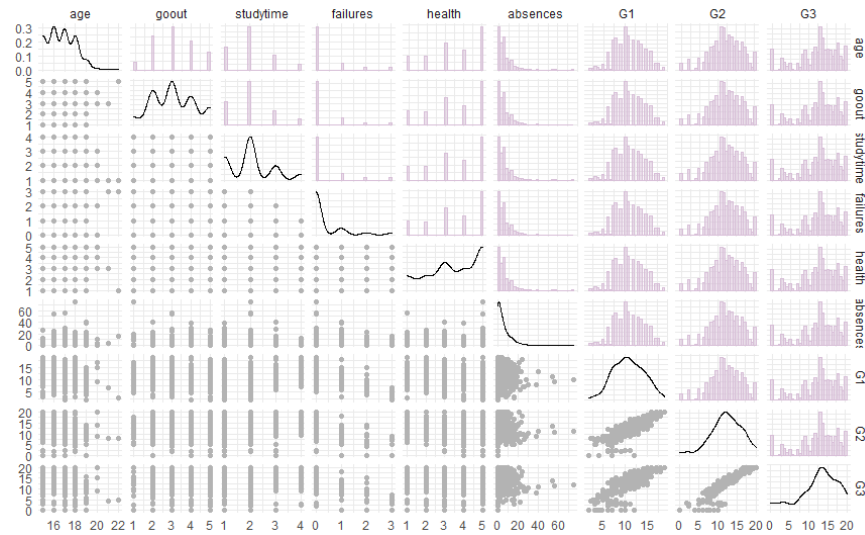


Figure 39: Bivariate Correlogram with barplot - scatterplot

b.

I used black for negative correlation and thistle for positive correlation (hope thats okay :))
 significance level = 0.05 .

The cells that are crossed are rejected by p-value.

(Note : diag. correlations are omitted)

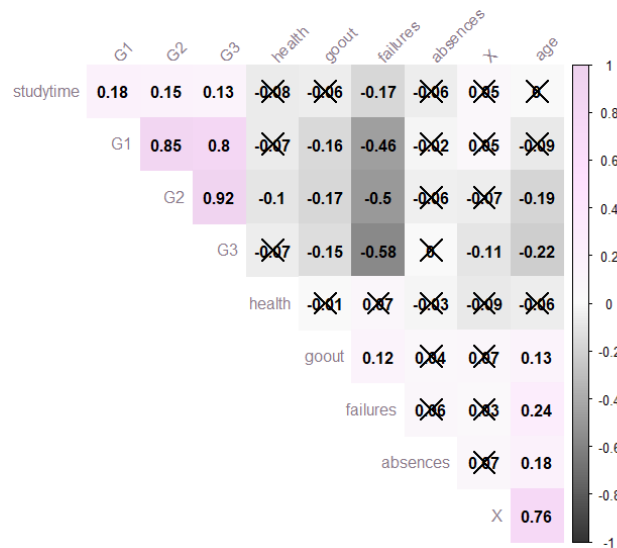


Figure 40: Heatmap correlogram of numerical values

c.

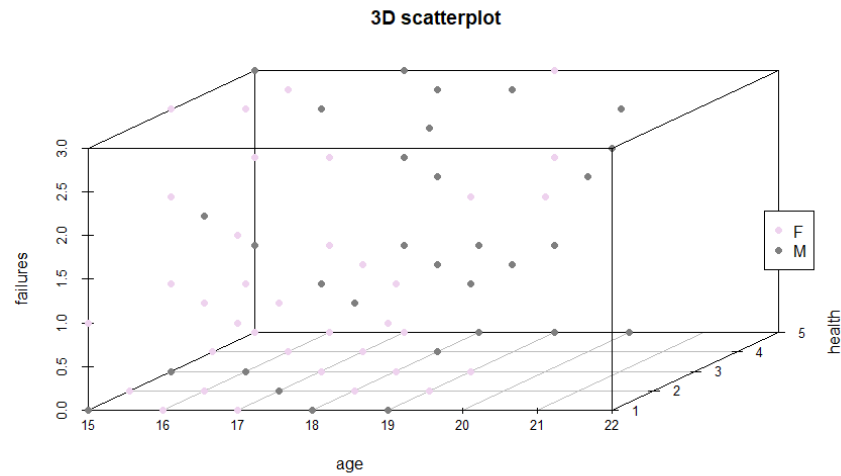


Figure 41: 3D scatterplot of age, failures and health colorized by sex

Unfortunately, it seems like there is not a specific relationship between these 3 variables; but we can see that *Females* have *Females* failures and *Males* and also, *Females* are in the younger *age* group.

Question 5

Chosen Categorical Variables : *sex* and *romantic*

a.

```
>  
> print.table(table)  
  
      F   M Sum  
no  129 134 263  
yes   79  53 132  
Sum  208 187 395  
>
```

Figure 42: Frequency/ Contingency table of sex and romantic

b.

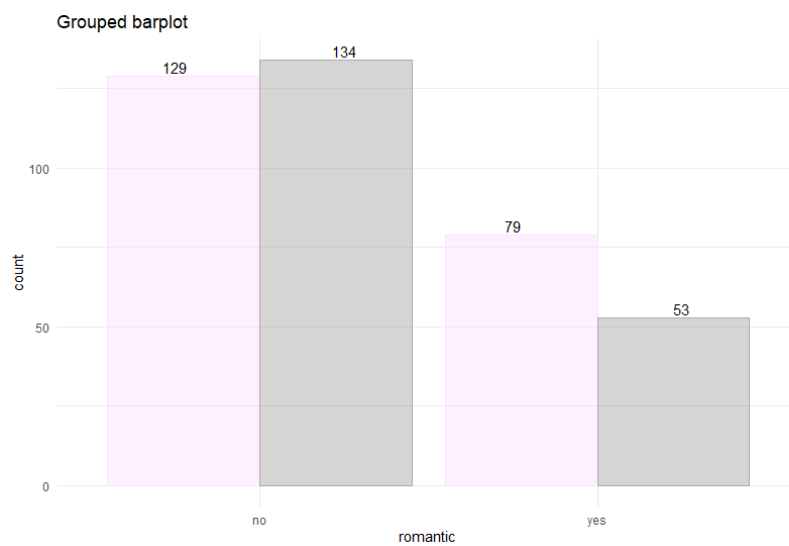


Figure 43: Grouped barplot of sex and romantic

c.

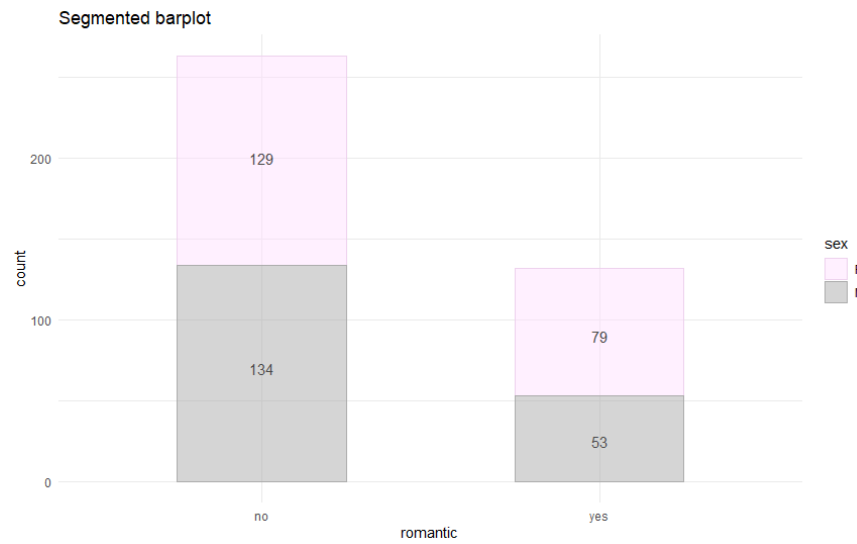


Figure 44: Segmented barplot of sex and romantic

d.

The segmented barplot does well in informing about the percent of each category within each group. The information that is missing is the size of each group.

A mosaic plot allows us to see these group sizes by scaling on the x-axis!

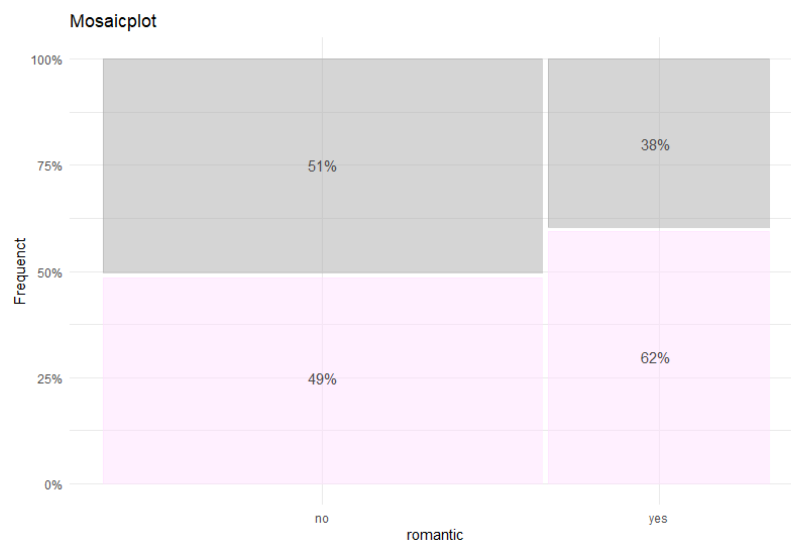


Figure 45: Mosaicplot of sex and romantic

Question 6

Chosen Numerical Variable : *goout*

Check Condition :

- Independent Observations :
 - Random sample/assignment
 - sampling without replacement, $395 < 10\%$ all of the students
- Sample size / skew :
 - $n < 30 \rightarrow$ t-test , $n > 30 \rightarrow$ z-test
 - skewness : *Figure 46* shows no skewness and also by checking mean and median of *age* in *Figure 2*, we can see that mean and median are pretty much the same so we are good to go.

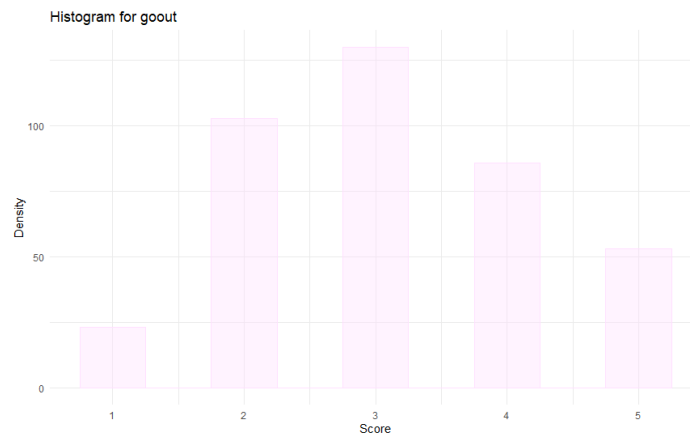


Figure 46: Histogram of goout

a.

Confidence intervals include the point estimate for the sample with a margin of error around the point estimate. The point estimate is the most likely value of the parameter and equals the sample value. The margin of error accounts for the amount of doubt involved in estimating the population parameter. The more variability there is in the sample data, the less precise the estimate, which causes the margin of error to extend further out from the point estimate.

Sample size = 25, t-test :

```
"Confidence Interval(using t-test) : ( 2.693 , 3.387 )"
```

Figure 47: Confidence Interval of goout using $\alpha = 5\%$

Sample size = 200, z-test :

```
"Confidence Interval(using z-test) : ( 2.92 , 3.23 )"
```

Figure 48: Confidence Interval of goout using $\alpha = 5\%$

b.

We are 95% confident that the the times these students goout are on average between 2.92 and 3.23 (according to *z-test*).

In other words, 95% of random samples of 395 students will yeild CIs that capture the true population mean of the times they goout.

c.

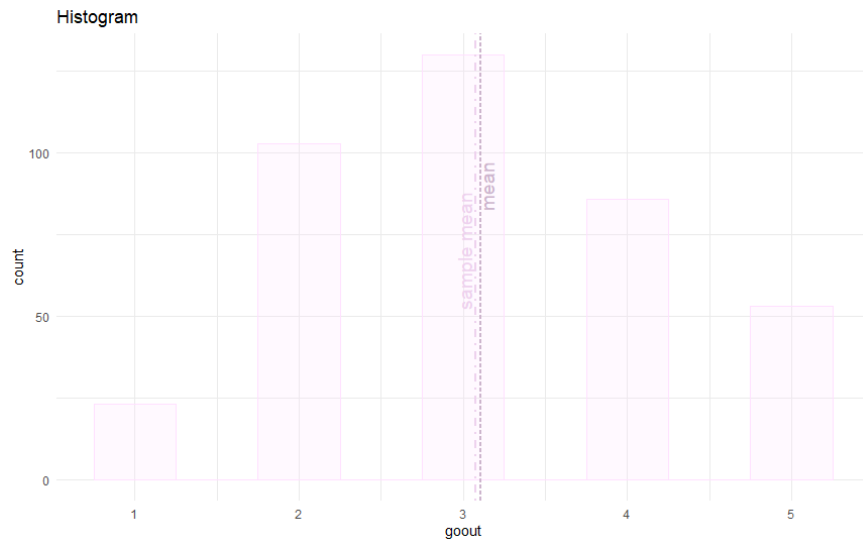


Figure 49: Histogram of goout marked with actual mean and sample mean

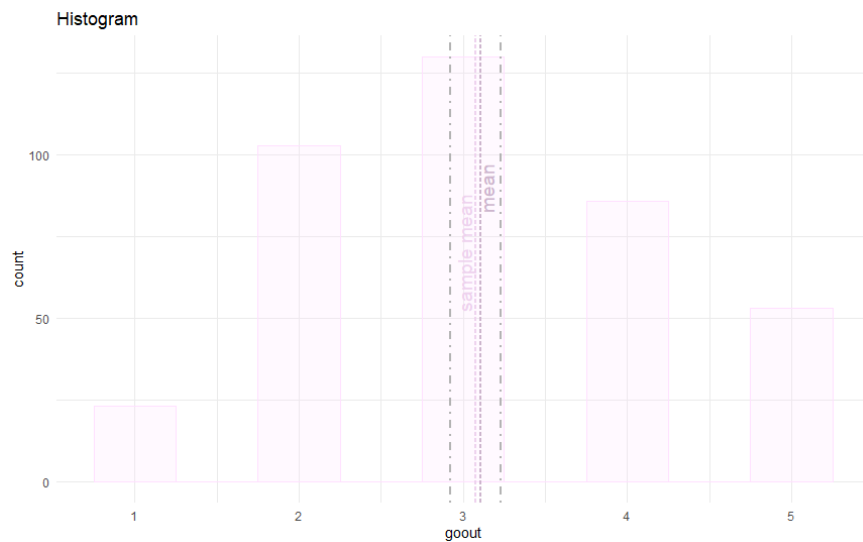


Figure 50: Histogram of goout marked with CI, actual mean and sample mean

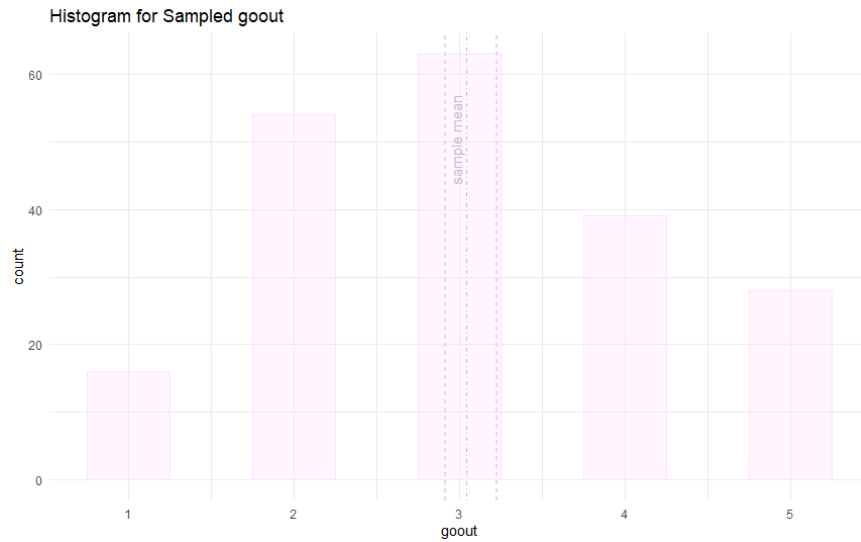


Figure 51: Histogram of sampled goout marked with CI and sample mean

d.

Hypothesis test :

$$H_0 : \mu = 2.8$$

$$H_A : \mu \neq 2.8$$

Sample size = 25, t-test :

```
>
> Hypothesis.test(goout.sampled.t, null.value = 2.8)
[1] "Null Hypothesis: mean = 2.8"
[1] "Alternative Hypothesis: mean /= 2.8"
[1] "Using t-distribution"
[1] "p-value = 0.0219829970441023"
[1] "Reject Null Hypothesis."
>
>
```

Figure 52: Hypothesis test of goout using $\alpha = 5\%$

Since *p-value* is 5% and is higher than 0.021, we should reject the null hypothesis in favor of the alternative hypothesis.

Sample size = 200, z-test :

```
>
> Hypothesis.test(goout.sampled, null.value = 2.8)
[1] "Null Hypothesis: mean = 2.8"
[1] "Alternative Hypothesis: mean /= 2.8"
[1] "Using z-distribution"
[1] "p-value = 0.000135695892379579"
[1] "Reject Null Hypothesis."
>
>
```

Figure 53: Hypothesis test of goout using $\alpha = 5\%$

Since *p-value* is 5% and is higher than 0.00013, we should reject the null hypothesis in favor of the alternative hypothesis.

According to *Figure 52*, if the null hypothesis were true, there is only 2.1% (very tiny) chance that we would take a sample of size 25 and obtain a sample mean of 3.07.

According to *Figure 53*, if the null hypothesis were true, there is a tiny chance that we would take a sample of size 25 and obtain a sample mean of 3.07.

e.

P-value and *Confidence Interval* are two equivalent methods of interpreting results of a statistical analysis and their results *always agree*.

Both of these concepts specify a distance from the mean to a limit and these distances are precisely the same length.

f. and g.

The error that occurs when one accepts a null hypothesis that is actually false is the type II error. A type II error produces a false negative, also known as an error of omission.

$$\beta = P(H_0 \text{ is true} \mid H_0 \text{ is actually false})$$

```
>
> TypeIIerr(goout.sampled, null.value = 2.8)
[1] "TypeII error = % 2.4"
[1] "Power = % 97.6"
>
>
```

Figure 54: Power and typeII error of goout

Using **R**'s built-in function :

```
one-sample t test power calculation

      n = 200
  delta = 0.3088608
    sd = 1.111837
sig.level = 0.05
  power = 0.974388
alternative = two.sided

> |
```

Figure 55: Power and typeII error of goout

An effect size is closely related to a power of a statistical test because when *difference* of two groups is big, it is easy to reject the null hypothesis.

In other words, as the effect size gets larger, it is more likely to reject the null hypothesis; less likely to fail to reject the null hypothesis, thus the power of the test increases.

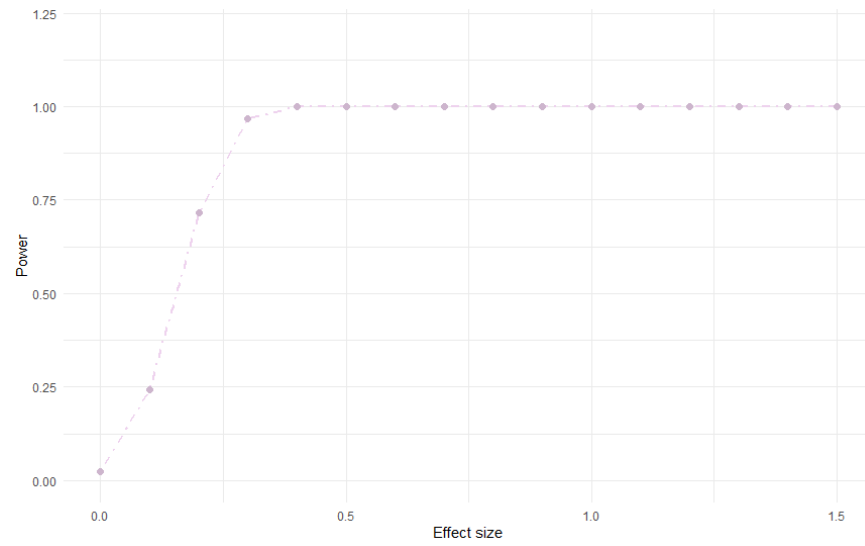


Figure 56: Relationship between effect size and power

Question 7

a.

Chosen Numerical Variable : *health* and *goout*

When two sets of observations have a special correspondence (they were chosen from one X in the dataset), they are said to be paired. To analyze paired data, it is useful to look at the difference in outcomes of each paired observation.

Check Condition :

- Independent Observations :
 - Random sample/assignment
 - sampling without replacement, $395 < 10\%$ all of the students
- Sample size / skew :
 - $n = 25 < 30 \rightarrow$ t-test. The Central Limit Theorem states that when the sample size is small, the normal approximation may not be very good. However, as the sample size becomes large, the normal approximation improves. Usually, t-tests are more appropriate when dealing with problems with a limited sample size .
 - skewness : As was mentioned in question6 (part a) *goout* is not skewed, *health* is a bit leftskewed.

```
>
> Hypothesis.test(StudentsPerformance.sampled$health, StudentsPerformance.sampled$goout, paired = TRUE)
[1] "Null Hypothesis: diff mean = 0"
[1] "Alternative Hypothesis: diff mean /= 0"
[1] "Using t-distribution"
[1] "p-value = 0.0384069445168378"
[1] "Reject Null Hypothesis."
>
>
```

Figure 57: Paired t-test between health and goout

Using **R**'s built-in function :

```
Paired t-test

data: StudentsPerformance.sampled$health and StudentsPerformance.sampled$goout
t = 2.1909, df = 24, p-value = 0.03841
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 0.0463708 1.5536292
sample estimates:
mean of the differences
              0.8
```

Figure 58: Paired t-test between health and goout

Since *p-value* is 5% and is higher than 0.038, we should reject the null hypothesis in favor of the alternative hypothesis. There is strong evidence that the null hypothesis is invalid.

b.

Check Condition :

- Independent Observations :
 - Random sample/assignment
 - sampling without replacement, $395 < 10\%$ all of the students
- Sample size / skew :
 - $n = 100 > 30 \rightarrow$ z-test
 - skewness : As was mentioned in question6 (part a) *goout* is not skewed, *health* is a bit leftskewed but our sample size is big enough so we can ignore it.

```
>
> Hypothesis.test(health.sampled, goout.sampled)
[1] "Null Hypothesis: diff mean = 0"
[1] "Alternative Hypothesis: diff mean /= 0"
[1] "Using Z-distribution"
[1] "p-value = 0.0355069327255375"
[1] "Reject Null Hypothesis."
>
>
```

Figure 59: z-test between health and goout

Using **R**'s built-in function :

```
      welch Two Sample t-test

data:  health.sampled and goout.sampled
t = 2.1025, df = 186.83, p-value = 0.03685
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 0.02469156 0.77530844
sample estimates:
mean of x mean of y
   3.48      3.08

> |
```

Figure 60: z-test between health and goout

Confidence Interval : $0 \notin [0.024, 0.77] \rightarrow$ Reject the null hypothesis.

P-value and *Confidence Interval* are two equivalent methods of interpreting results of a statistical analysis and their results *always agree*.

Question 8

Chosen Numerical Variable : *absences*

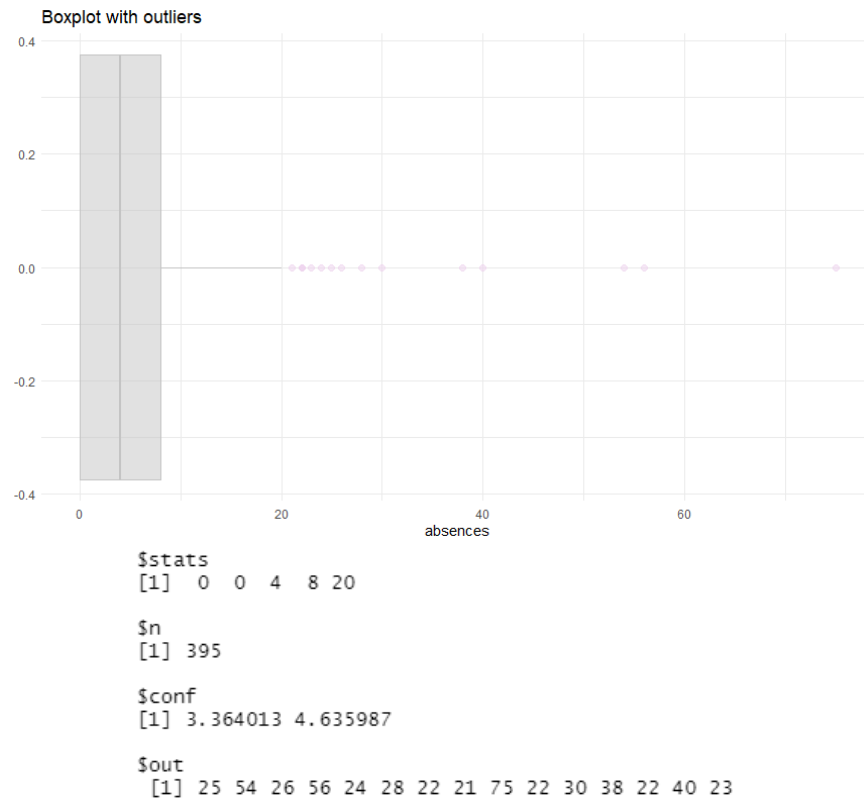


Figure 61: boxplot of absences with stat

Using the normal approximation might not be good in all applications where the sample size is at least 30. Generally, the more skewed a population distribution or the more common the frequency of outliers, the larger the sample required to guarantee the distribution of the sample mean is nearly normal.

a.

Using quantile doesn't seem like a good idea as can be deducted from *Figure 62*, so 100 samples were chosen and replicated 1000 times, and the interval for their mean can be seen in *Figure*.

```

>
> quantile(StudentsPerformance$absences, c(0.025, 0.975))
2.5% 97.5%
0.00 23.15
>
>
  
```

Figure 62: Simple percentile method

```
"confidence Interval: ( 5.18 , 7.74 )"
```

Figure 63: Percentile method

b.

A random sample with replacement was taken from the original sample. Bootstrap statistic (*mean* in our case) was computed on bootstrap samples and these steps was repeated to create a bootstrap distribution. The middle 95% of the bootstrap distribution was calculated for CI :

```
"confidence Interval: ( 5.56 , 6.22 )"
```

Figure 64: Percentile method (bootstrapped)

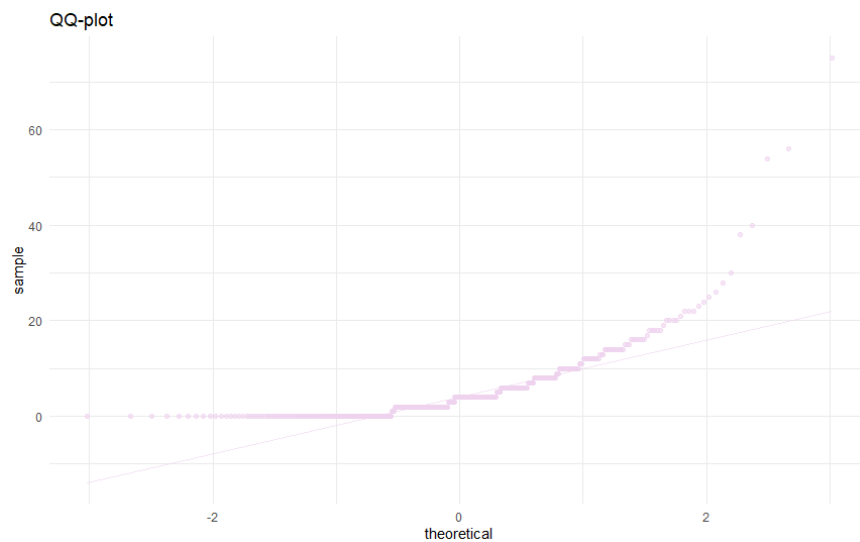
c.

Figure 65: QQ-plot of absences

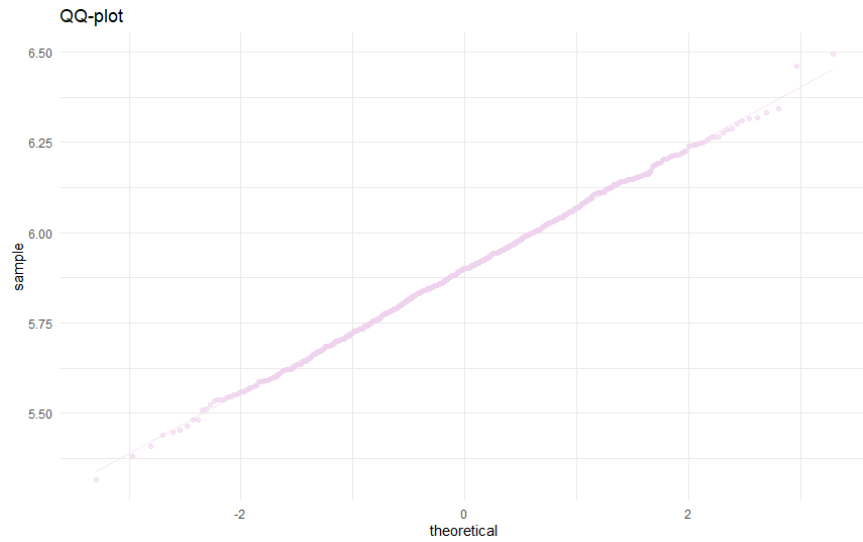


Figure 66: QQ-plot of mean of bootstrapped samples

Percentile method is a method which is sensitive to outliers; so, the calculated interval might not be as informative as we desired. Therefore, bootstrapping method was used in order to remove outliers and result a (approximately) normal distribution (*Figure 66*). (A better approach is using SD method which is more robust when facing outliers)

Knowing these facts and figures, we can conclude that *bootstrapping* is a stronger procedure and a more informative CI is the proof of it.

Question 9

In *ANOVA*, the *null hypothesis* is that there is no difference among group means. If any group differs significantly from the overall group mean, then the *ANOVA* will report a statistically significant result.

In our case :

$$H_0 : \mu_{failure=0, G_1+G_2+G_3} = \mu_{failure=1, G_1+G_2+G_3} = \mu_{failure=2, G_1+G_2+G_3} = \mu_{failure=3, G_1+G_2+G_3}$$

$$H_A : \text{one group differs significantly from the overall group mean}$$

Significant differences among group means are calculated using the *F statistic*, which is the ratio of the mean sum of squares (explained variable) to the mean square error (unexplained variable) .

If the *F statistic* is higher than the alpha value (0.05), then the difference among groups is deemed statistically significant.

Degrees of freedom associated with *ANOVA* :

$$df_T = n - 1 \quad , \quad df_G = k - 1 \quad , \quad df_E = df_T - df_G = n - k$$

$$SST = \sum_{i=1}^n (x_i - \bar{x})^2$$

$$SSG = \sum_{j=1}^k n_j (\bar{x}_j - \bar{x})^2 \Rightarrow MSG = \frac{1}{k-1} \sum_{j=1}^k n_j (\bar{x}_j - \bar{x})^2$$

$$SSE = \sum_{j=1}^k \sum_{i=1}^{n_j} (\bar{x}_j - \bar{x})^2 = \sum_{j=1}^k (n_j - 1) s_j^2 \Rightarrow MSE = \frac{1}{n-k} \sum_{j=1}^k (n_j - 1) s_j^2$$

$$F = \frac{\text{Variability bet. groups}}{\text{Variability w/in groups}} = \frac{MSG}{MSE}$$

Check Condition :

- Independence :
 - within groups: sampled observations are independent
 - between groups: the groups are independent of each other (non-paired)
- Approximate normality : distributions should be nearly normal within each group → we assume they are

- Equal variance : groups should have roughly equal variability

```
>
> sd.df
  groups      sds
1 Group0 10.276562
2 Group1 10.082132
3 Group2 10.556222
4 Group3  6.172193
>
>
```

Figure 67: SD of each group

The standard deviation of group0, group1 and group2 are close to each other, but the one for group3 is different from others. Although this could happen because of the low group size, we can consider these three numbers as almost the same.

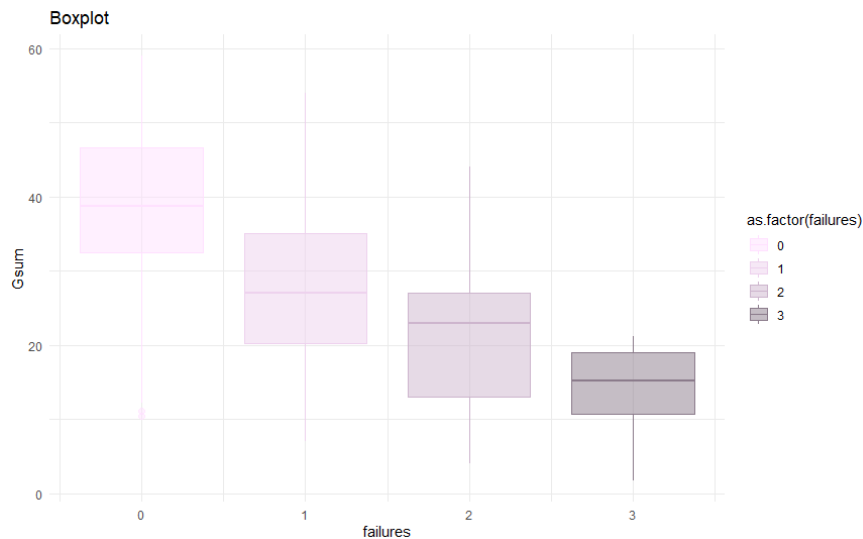


Figure 68: Boxplot grouped by number of failures

```
>
> summary(aov.Gsum_failures)
              Df Sum Sq Mean Sq F value Pr(>F)
as.factor(failures) 3  17949    5983   58.22 <2e-16 ***
Residuals          391  40179     103
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
>
>
```

Figure 69: ANOVA table

Since p -value is smaller than 0.05, we reject the null Hypothesis.

The data provides convincing evidence that at least one pair of population means are different from each other.

ANOVA tells us if there are differences among group means, but *not what the differences* are. To find out which groups are statistically different from one another, you can perform a *Tukey's Honestly Significant*

Difference (Tukey's HSD) post-hoc test for pairwise comparisons.

The significant groupwise differences are any where the 95% confidence interval doesn't include zero. In other words, p-value for these *pairwise differences* is < 0.05 .

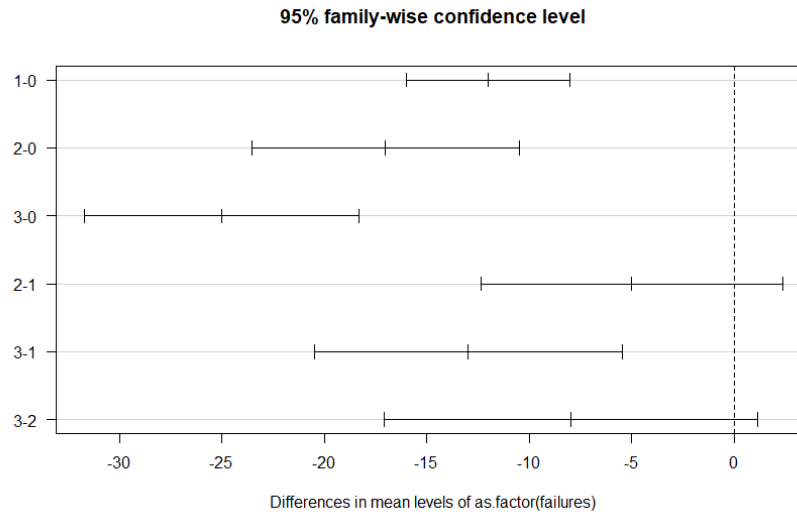


Figure 70: Pairwise confidence level(0.95%)

R Codes

```

1 library(magrittr)
2 library(ggfortify)
3 library(ggplot2)
4 library(plyr)
5 library(gridExtra)
6 require(qqplotr)
7 library(moments)
8 library(hexbin)
9 library(ggmosaic)
10 library("plot3D")
11 library(plotly)
12 library(scatterplot3d)
13 library(RNHANES)
14 library(GGally)
15 library(dplyr)
16 library(Hmisc)
17 require(ggpubr)
18 require(Hmisc)
19 require(corrplot)
20 library(patchwork)
21 library(ggExtra)
22
23
24 theme_set(theme_minimal())
25
26 summary(StudentsPerformance)
27
28 #Question 0
29
30 missingvalues <- colSums(is.na.data.frame(StudentsPerformance))
31 missingvalues.proporion <- missingvalues/nrow(StudentsPerformance)
32 plot(missingvalues.proporion , main = "Percentage of missing values vs. variables",
33      xlab = "variables", ylab = "Missing values proportion" , type = 'l' , col = 'thistle')
34
35 missingvalues.proporion <- data.frame("missing value", missingvalues/nrow(
36      StudentsPerformance))
37
38 #QUESTION 1
39
40 #Numerical value chosen : Grade 1
41
42 StudentsPerformance$G1
43
44 #a.
45 breaks <- pretty(StudentsPerformance$G1, n = nclass.FD(StudentsPerformance$G1), min.n = 0)
46 bwidth <- breaks[2] - breaks[1]
47
48
49 G1_hist <- ggplot(StudentsPerformance, aes(x = G1)) +
50   geom_histogram(aes(y=..density..) , binwidth = bwidth, alpha = 0.4, color="thistle1", fill
51     ="thistle1") +
52   geom_density(color = "gray87", linetype="dashed", fill = "gray87" , alpha = 0.3, size=1) +
53   #stat_function(fun = dnorm, n = 101, args = list(mean = mu, sd = std) , color = "gray87" ,
54     size=1) +
55   labs(title = "Histogram for G1", x = "Score", y="Density")
56
57 G1_hist
58 #

```

```

58
59 #b.
60 G1.qq <- ggplot(StudentsPerformance, aes(sample = G1, color = "", alpha = 0.7)) + geom_qq()
61   +
62   geom_qq_line() + labs(title="QQ-plot for G1")
63
64 #——
65
66 #c.
67 print(skewness(StudentsPerformance$G1))
68
69 G1.hist + geom_vline(xintercept = mean(StudentsPerformance$G1), linetype="dashed", color = "
70   thistle4", size = 0.5) +
71   geom_vline(xintercept = median(StudentsPerformance$G1), linetype="dotdash", color = "
72   gray29", size = 0.5)+
73   annotate("text", x = mu - .2 , label = "mean", y = 0.01, size = 3.4, angle = 90 , color =
74   'thistle4') +
75   annotate("text", x = median + 0.1 , label = "median", y = 0.06, size = 3.4, angle = 90,
76   color = 'gray29')
77
78 #——
79
80 #d.
81 G1.box <- ggplot(StudentsPerformance, aes(x = G1)) + geom_boxplot(color ="thistle2", fill ="
82   thistle2", alpha = 0.5) +
83   labs(title="Boxplot for G1")
84
85 #——
86
87 #e.
88 mu <- mean(StudentsPerformance$G1)
89 median <- median(StudentsPerformance$G1)
90 var <- var(StudentsPerformance$G1)
91 std <- sd(StudentsPerformance$G1)
92
93 #——
94
95 #f.
96 G1.density <- ggplot(StudentsPerformance, aes(x = G1)) +
97   geom_vline(xintercept = mu, linetype="dashed", color = "gray29") +
98   geom_vline(xintercept = median, linetype="dashed", color = "thistle4") +
99   geom_density(color = "thistle1", size = 1) +
100   stat_function(fun = dnorm, n = 101, args = list(mean = mu, sd = std) , color = "thistle2",
101     size = 1) +
102   annotate("text", x = mu - .2 , label = "mean", y = 0.01, size = 3.4, angle = 90 , color =
103     'gray29') +
104   annotate("text", x = median + 0.1 , label = "median", y = 0.06, size = 3.4, angle = 90,
105     color = 'thistle4')+
106   labs(title="Density for G1")
107
108 #——
109
110 #g.
111
112 #method1
113 StudentsPerformance$categorizedG1 <- ifelse(StudentsPerformance$G1 > (mu + max(
114   StudentsPerformance$G1))/2, 'very high', ifelse(StudentsPerformance$G1 > mu, 'high',
115   ifelse(StudentsPerformance$G1 > mu/2, 'low', 'very low')))
116
117 freq_vlow <-length(which(StudentsPerformance[17] == 'very low')) / length(

```

```

    StudentsPerformance$G1)
109 freq_low <- length(which(StudentsPerformance[17] == 'low'))/ length(StudentsPerformance$G1)
110 freq_high <- length(which(StudentsPerformance[17] == 'high'))/ length(StudentsPerformance$G1)
111 freq_vhigh <- length(which(StudentsPerformance[17] == 'very high'))/ length(
    StudentsPerformance$G1)
112
113
114 G1.categorized <- data.frame(group = c("Very Low", "Low", "High", "Very High"),
115     value = c(freq_vlow, freq_low, freq_high, freq_vhigh))
116
117
118
119 G1.pie <- ggplot(G1.categorized, aes(x="", y = value, fill = group)) +
120     geom_bar(stat = "identity", alpha = 0.7) + coord_polar("y")
121
122
123 G1.pie + scale_fill_manual(values = c("thistle1", "thistle2", "thistle3", "thistle4")) +
124     geom_text(aes(label = paste0(round(value*100), "%"), position = position_stack(vjust =
125     0.5)) +
126     labs(title="Pie-chart of G1", x = 'Frequency', y = 'G1')
127 #method2
128
129 G1.quant <- quantile(StudentsPerformance$G1)
130
131 StudentsPerformance$categorizedG1 <- ifelse(StudentsPerformance$G1 > G1.quant[[4]], 'very
    high', ifelse(StudentsPerformance$G1 > G1.quant[[3]], 'high', ifelse(
    StudentsPerformance$G1 > G1.quant[[2]], 'low', 'very low')))
132
133 freq_vlow <-length(which(StudentsPerformance[17] == 'very low')) / length(
    StudentsPerformance$G1)
134 freq_low <- length(which(StudentsPerformance[17] == 'low'))/ length(StudentsPerformance$G1)
135 freq_high <- length(which(StudentsPerformance[17] == 'high'))/ length(StudentsPerformance$G1)
136 freq_vhigh <- length(which(StudentsPerformance[17] == 'very high'))/ length(
    StudentsPerformance$G1)
137
138
139 G1.categorized <- data.frame(group = c("Very Low", "Low", "High", "Very High"),
140     value = c(freq_vlow, freq_low, freq_high, freq_vhigh))
141
142
143
144 G1.pie <- ggplot(G1.categorized, aes(x="", y = value, fill = group)) +
145     geom_bar(stat = "identity", alpha = 0.7) + coord_polar("y")
146
147
148 G1.pie + scale_fill_manual(values = c("thistle1", "thistle2", "thistle3", "thistle4")) +
149     geom_text(aes(label = paste0(round(value*100), "%"), position = position_stack(vjust =
150     0.5)) +
151     labs(title="Pie-chart of G1", x = 'Frequency', y = 'G1')
152
153 #
154
155 #h.
156
157 boxplot.stats(StudentsPerformance$G1)
158
159 G1.quant <- quantile(StudentsPerformance$G1)
160 G1.iqr <- IQR(StudentsPerformance$G1)

```

```

161 #————
162 #————
163
164 #QUESTION 2
165
166 #Categorical Variable chosen : sex
167
168 StudentsPerformance$sex
169
170 #a.
171 female.freq <- length((((StudentsPerformance %>% filter(sex == 'F'))$sex)) / length(
  StudentsPerformance$sex)
172 male.freq <- length((((StudentsPerformance %>% filter(sex == 'M'))$sex)) / length(
  StudentsPerformance$sex)
173 #————
174
175 #StudentsPerformance.sel <- subset(StudentsPerformance, sex == "F")
176
177
178 #b.
179
180 #freq <-data.frame(female.freq, male.freq)
181 sex.barplot <- ggplot(StudentsPerformance, aes(x = " ", color = sex, fill = sex)) +
182   geom_bar(aes(y = (..count..)/sum(..count..)), alpha = 0.5, width = 0.5) + labs(title="
  Stacked Barplot of sex", y = 'Frequency')
183
184 sex.barplot + scale_color_manual(values = c("thistle1", "gray67")) + xlab("sex") +
185   scale_fill_manual(values = c("thistle1", "gray67")) + scale_y_continuous(labels = scales::
  percent) +
186   geom_text(aes(y = ((..count..)/sum(..count..)), label = scales::percent((..count..)/sum(..
  count..))),
187     stat = "count", hjust = 0.5, size = 4.5, color = 'black', vjust = 1.4, position
     = position_dodge(width = 0.3))
188 #————
189
190 #c.
191 categorizedsex.barplot <- ggplot(StudentsPerformance, aes(x = sex, color = sex, fill = sex))
  +
192   geom_bar(aes(y = (..count..)/sum(..count..)), alpha = 0.7) + labs(title="Barplot of sex",
  y = 'Frequency')
193
194 categorizedsex.barplot + scale_color_manual(values = c("thistle1", "gray67")) +
195   scale_fill_manual(values = c("thistle1", "gray67")) + coord_flip() + scale_x_discrete(
  limits=c("M", "F"))+
196   geom_text(aes(y = ((..count..)/sum(..count..)), label = round(((..count..)/sum(..count..))
  , 3)),
197     stat = "count", vjust = -0.25, size = 4, color = 'black') + theme(axis.text.x=
  element_blank())
198 #————
199
200 #d.
201 sex.df <- data.frame(sex = c("F", "M"), frequency = c(female.freq, male.freq))
202 sex.violinplot <- ggplot(StudentsPerformance, aes(x = sex, y = age, color = sex, fill = sex)
  ) +
203   geom_violin(trim=FALSE, alpha = 0.7) + labs(title="Violinplot")
204
205 sex.violinplot+ scale_color_manual(values = c("thistle1", "gray67")) + xlab("sex") +
206   scale_fill_manual(values = c("thistle1", "gray67"))
207 #————
208 #————
209
210 #QUESTION 3

```

```

211
212 #Numerical Variable chosen : goout and absences -> it actually depends
213
214
215 #b.
216 goout_absences.scatterplot <- ggplot(StudentsPerformance, aes(x = goout, y = absences)) +
217   geom_point(color = "thistle2", size = 2)
218
219 goout_absences.scatterplot + labs(title="Scatterplot of Going out and Absences")
220 #——
221
222 #c.
223 goout_absences.correlation <- cor(StudentsPerformance$goout, StudentsPerformance$absences)
224 goout_absences.correlation
225 #——
226
227 #c.
228 ggscatter(StudentsPerformance, x = "goout", y = "absences", shape = 12, add = "reg.line",
229           conf.int = TRUE,
230           color = "thistle2", add.params = list(color = "thistle3", fill = "gray90"), cor.
231           coef = TRUE,
232           cor.coeff.args = list(method = "pearson", label.x = 3, label.sep = "\n")) + theme_
233           minimal() +
234           labs(title="Scatterplot")
235 #——
236
237 #f.
238 goout_absences_romantic.scatterplot <- ggplot(StudentsPerformance,
239           aes(x = goout, y = absences, color = romantic,
240             shape = romantic)) +
241           geom_point(size = 2) + labs(title="Scatterplot")
242 #——
243
244 #g.
245 breaks <- pretty(StudentsPerformance$goout, n = nclass.FD(StudentsPerformance$goout), min.n
246   = 0)
247 goout.hist <- ggplot(StudentsPerformance, aes(x = goout)) + geom_histogram(binwidth = breaks
248   [2]-breaks[1], color = "thistle2", fill = "thistle2", alpha = 0.5) +theme_void()
249
250 breaks <- pretty(StudentsPerformance$absences, n = nclass.FD(StudentsPerformance$absences),
251   min.n = 0)
252 absences.hist <- ggplot(StudentsPerformance, aes(x = absences)) +
253   geom_histogram(binwidth = breaks[2]-breaks[1], color = "thistle2", fill = "thistle2",
254     alpha = 0.5) + coord_flip() + theme_void()
255
256
257 gar.hexbinplot.log <- ggplot(StudentsPerformance, aes(x = goout, y = log2(absences + 1))) +
258   geom_hex(bins = 10, color = "white", alpha = 0.7) + scale_fill_gradient(low = "thistle1",
259     high = "thistle4", trans="log10")
260 #+ geom_smooth(col = 'grey40')
261
262 goout.hist + plot_spacer() + gar.hexbinplot.log + absences.hist +
263   plot_layout(ncol = 2, nrow = 2, widths = c(4, 1), heights = c(1, 4))
264
265 gar.hexbinplot <- ggplot(StudentsPerformance, aes(x = goout, y = absences)) +
266   geom_hex(bins = 10, color = "white", alpha = 0.7) + scale_fill_gradient(low = "thistle1",
267     high = "thistle4") +
268   geom_smooth(method = "loess", col = 'grey40')

```

```

263
264 goout.hist + plot_spacer() + gar.hexbinplot + absences.hist +
265   plot_layout( ncol = 2, nrow = 2, widths = c(4, 1), heights = c(1, 4))
266
267
268 #——
269
270 #h.
271 gar.2ddensity <- ggplot(StudentsPerformance, aes(x = goout, y = G1)) +
272   stat_density2d(aes(fill = ..level..), geom = "polygon", alpha = 0.5) + lims(x = c(0,6), y
    = c(-5, 20))
273
274 gar.2ddensity + scale_fill_gradient(low = "thistle1", high = "thistle4") +
275   labs(title="2D-density plot")
276 #——
277 #—————
278
279 #QUESTION 4
280
281 #a.
282 ggpairs(dplyr::select_if(StudentsPerformance, is.numeric), title = "Correlogram")
283
284
285 #density, without failure
286 ggpairs(StudentsPerformance[, c(4, 7, 10, 12, 13, 14, 15, 16)],
287   upper = list(continuous = wrap("density", colour="thistle")),
288   lower = list(continuous = wrap("points", colour="grey70")))
289
290 #linear relationship
291 ggpairs(StudentsPerformance[, c(4, 7, 10, 11, 12, 13, 14, 15, 16)],
292   upper = list(continuous = wrap("smooth", colour="thistle")),
293   lower = list(continuous = wrap("points", colour="grey70")))
294 #barplot
295 ggpairs(StudentsPerformance[, c(4, 7, 10, 11, 12, 13, 14, 15, 16)],
296   upper = list(continuous = wrap("barDiag", colour="thistle", fill = "thistle", alpha
    = 0.5)),
297   lower = list(continuous = wrap("points", colour="grey70")))
298 #boxplot
299 ggpairs(StudentsPerformance[, c(4, 7, 10, 11, 12, 13, 14, 15, 16)],
300   upper = list(continuous = wrap("box_no_facet", colour="thistle", fill = "thistle3",
    alpha = 0.5)),
301   lower = list(continuous = wrap("points", colour="grey70")))
302
303 #——
304 #b.
305
306 col <- colorRampPalette(c("grey80", "white", "thistle1", "thistle2"))
307 StudentsPerformance.corr <- rcorr(as.matrix(dplyr::select_if(StudentsPerformance, is.numeric
    )))
308 StudentsPerformance.corr.p <- StudentsPerformance.corr$P
309 StudentsPerformance.corr.p[is.na(StudentsPerformance.corr.p)] <- 1
310
311 M <- cor(dplyr::select_if(StudentsPerformance, is.numeric))
312
313 corplot(M, method = "color", col = col(200), type = "upper", order = "hclust", addCoef.col
    = "black",
314   tl.col = "thistle4", tl.srt = 45, p.mat = StudentsPerformance.corr.p, sig.level =
    0.05, diag = FALSE)
315 #——
316
317 #c.
318 cols <- c("thistle2", "grey50")

```

```

319 with(StudentsPerformance, scatterplot3d(age, health, failures, main="3D scatterplot",
320     pch = 16, color = cols[as.numeric(StudentsPerformance$sex)]))
321
322 legend("right", legend = levels(StudentsPerformance$sex),
323     col = c("thistle2", "grey50"), pch = 16)
324
325 #-----
326
327 #Question 5
328
329 #Chosen : sex and romantic
330 #a.
331 table <- addmargins(table(StudentsPerformance$romantic, StudentsPerformance$sex), c(1,2))
332 print.table(table)
333 #-----
334
335
336
337 #b.
338
339 romantic_sex.groupedbarplot <- ggplot(StudentsPerformance, aes(x = romantic, color = sex,
340     fill = sex)) +
341     geom_bar(position = "dodge", alpha = 0.5) + labs(title="Grouped barplot", x="romantic")
342
343 romantic_sex.groupedbarplot + scale_color_manual(values = c("thistle1", "gray67")) +
344     scale_fill_manual(values = c("thistle1", "gray67")) +
345     geom_text(aes(y = ..count.., label = ..count..), stat = "count", vjust = -0.25, size = 4,
346         color = 'black', position = position_dodge(width = 1))
347
348 romantic_sex.groupedbarplot <- ggplot(StudentsPerformance, aes(x = romantic, color = sex,
349     fill = sex)) +
350     geom_bar(alpha = 0.5, width = 0.5) + labs(title="Segmented barplot", x="romantic")
351
352 romantic_sex.groupedbarplot + scale_color_manual(values = c("thistle2", "gray68")) +
353     scale_fill_manual(values = c("thistle1", "gray67")) +
354     annotate("text", x = 1, label = "134", y = 70, size = 4, angle = 0, color = 'gray29') +
355     annotate("text", x = 1, label = "129", y = 200, size = 4, angle = 0, color = 'gray29') +
356     annotate("text", x = 2, label = "53", y = 25, size = 4, angle = 0, color = 'gray29') +
357     annotate("text", x = 2, label = "79", y = 90, size = 4, angle = 0, color = 'gray29')
358
359 romantic_sex.mosaicplot <- ggplot(StudentsPerformance) +
360     geom_mosaic(aes(x = product(romantic), fill = sex), alpha = 0.5) + labs(title="Mosaicplot"
361     , y = "Frequenct") +
362     scale_y_continuous(labels = scales::percent) +
363     annotate("text", x = 0.33, label = "49%", y = .25, size = 4, angle = 0, color = 'gray29')
364     +
365     annotate("text", x = 0.33, label = "51%", y = .75, size = 4, angle = 0, color = 'gray29')
366     +
367     annotate("text", x = 0.83, label = "62%", y = .3, size = 4, angle = 0, color = 'gray29')
368     +
369     annotate("text", x = 0.83, label = "38%", y = .8, size = 4, angle = 0, color = 'gray29')
370
371 romantic_sex.mosaicplot + scale_fill_manual(values = c("thistle1", "gray67"))
372
373 #-----
374 #-----
375
376 #Question 6

```

```

374 #Chosen Numerical Variable: age
375
376 breaks <- pretty(StudentsPerformance$goout, n = nclass.FD(StudentsPerformance$goout), min.n
    = 0)
377 bwidth <- breaks[2] - breaks[1]
378
379
380 goout_hist <- ggplot(StudentsPerformance, aes(x = goout)) +
381   geom_histogram(binwidth = bwidth, alpha = 0.4, color="thistle1", fill="thistle1") +
382   labs(title = "Histogram for goout", x = "Score", y="Density")
383 goout_hist
384
385 #a.
386 CI.calculate <- function(data.sampled, alpha = 0.05){
387
388   sample.len <- length(data.sampled)
389
390   mu <- mean(data.sampled)
391   s <- sd(data.sampled)
392   SE <- s/sqrt(sample.len)
393
394   if(sample.len > 30){
395     print("Using Z-distribution")
396     Zstar <- abs(qnorm(alpha/2))
397     error.margin <- Zstar * SE}
398
399   else{
400     print("Using t-distribution")
401     tstar <- abs(qt(alpha/2, df = sample.len - 1))
402     error.margin <- tstar * SE }
403   confidence.interval <- c(mu - error.margin, mu + error.margin)
404   return(confidence.interval)
405 }
406
407
408 goout.sampled.t <- sample(StudentsPerformance$goout, 25)
409 confidence.interval.t <- CI.calculate(goout.sampled.t)
410 print(paste("Confidence Interval(using t-test) : (", round(confidence.interval.t[1], 3), ",",
    round(confidence.interval.t[2], 3), ")"))
411
412
413 goout.sampled <- sample(StudentsPerformance$goout, 200)
414 confidence.interval <- CI.calculate(age.sampled)
415 print(paste("Confidence Interval(using z-test) : (", round(confidence.interval[1], 3), ",",
    round(confidence.interval[2], 3), ")"))
416 #————
417
418 #c.
419
420 goout_hist <- ggplot(StudentsPerformance, aes(x = goout)) +
421   geom_histogram(binwidth = bwidth, alpha = 0.2, color="thistle1", fill="thistle1") +
422   labs(title = "Histogram", x = "goout") +
423   geom_vline(xintercept = mean(StudentsPerformance$goout), color = "thistle3", linetype="21",
    size = 0.8) +
424   geom_vline(xintercept = mean(goout.sampled), color = "thistle2", linetype="dotted", size
    = 0.8) +
425   annotate("text", x = mean(StudentsPerformance$goout) + 0.03, label = "mean", y = 90, size
    = 5, angle = 90, color = "thistle3") +
426   annotate("text", x = mean(goout.sampled) - 0.07, label = "sample mean", y = 70, size = 5,
    angle = 90, color = "thistle2")
427
428 goout_hist

```



```

429
430
431 goout.hist <- ggplot(StudentsPerformance, aes(x = goout)) +
432   geom_histogram(binwidth = bwidth, alpha = 0.2, color="thistle1", fill="thistle1") +
433   labs(title = "Histogram", x = "goout") +
434   geom_vline(xintercept = mean(StudentsPerformance$goout), color = "thistle3", linetype="21",
435             size = 0.8) +
436   geom_vline(xintercept = mean(goout.sampled), color = "thistle2", linetype="21", size =
437             0.8) +
438   geom_vline(xintercept = round(confidence.interval[1], 3), color = "grey70", linetype="
439             dotted", size = 1) +
440   geom_vline(xintercept = round(confidence.interval[2], 3), color = "grey70", linetype="
441             dotted", size = 1) +
442   annotate("text", x = mean(StudentsPerformance$goout) + 0.03, label = "mean", y = 90, size
443           = 5, angle = 90, color = "thistle3") +
444   annotate("text", x = mean(goout.sampled) - 0.07, label = "sample mean", y = 70, size = 5,
445           angle = 90, color = "thistle2")
446
447 goout.hist
448
449
450 breaks <- pretty(goout.sampled, n = nclass.FD(goout.sampled), min.n = 0)
451 bwidth <- breaks[2] - breaks[1]
452
453 goout.df <- data.frame(goout.sampled)
454 sampled.goout.hist <- ggplot(goout.df, aes(x = goout.sampled)) +
455   geom_histogram(binwidth = bwidth, alpha = 0.3, color="thistle1", fill="thistle1") +
456   labs(title = "Histogram for Sampled goout", x = "goout") +
457   geom_vline(xintercept = mean(goout.sampled), color = "thistle3", linetype="dotted", size
458             = 0.5) +
459   geom_vline(xintercept = confidence.interval[1], color = "thistle2", linetype="22", size =
460             1) +
461   geom_vline(xintercept = confidence.interval[2], color = "thistle2", linetype="22", size =
462             1) +
463   annotate("text", x = mean(goout.sampled) - 0.07, label = "sample mean", y = 50, size = 4
464           , angle = 90, color = "thistle3")
465
466 sampled.goout.hist
467
468
469 #d.
470 Hypothesis.test <- function(data.sampled, null.value, alpha = 0.05){
471
472   sample.len <- length(data.sampled)
473   print(paste("Null Hypothesis: mean = ", null.value))
474   print(paste("Alternative Hypothesis: mean /= ", null.value))
475
476   x_bar <- mean(data.sampled)
477   s <- sd(data.sampled)
478   SE <- s/sqrt(sample.len)
479   score <- abs((x_bar - null.value)) / SE
480
481   if(sample.len > 30){
482     print("Using Z-distribution")
483     pvalue <- 2*pnorm(score, lower.tail = FALSE)}
484   else{

```

```

481   print("Using t-distribution")
482   pvalue <- 2*pt(score, df = sample.len - 1, lower.tail = FALSE)}
483
484
485   print(paste("p-value =", pvalue))
486
487   if (pvalue < alpha)
488     print("Reject Null Hypothesis.")
489   else
490     print("Fail to Reject Null Hypothesis.")
491 }
492 mean(goout.sampled)
493
494 Hypothesis.test(goout.sampled.t, null.value = 2.8)
495
496 Hypothesis.test(goout.sampled, null.value = 2.8)
497 #-----
498
499 #f. and #g.
500 TypeIIerr <- function(data.sampled, null.value, alpha = 0.05){
501   sample.len <- length(data.sampled)
502   mean.actual <- mean(StudentsPerformance$goout)
503   s <- sd(data.sampled)
504   SE <- s/sqrt(sample.len)
505   ME <- abs(qnorm((alpha/2))) * SE
506   errorTypeII <- pnorm(abs(null.value + ME - mean.actual)/SE, lower.tail = F) +
507     pnorm(abs(null.value - ME - mean.actual)/SE, lower.tail = F)
508
509   print(paste("TypeII error = %", 100*round(errorTypeII,3)))
510   print(paste("Power = %", 100*round(1-errorTypeII,3)))
511 }
512
513 TypeIIerr(goout.sampled, null.value = 2.8)
514
515 power.t.test(n = 200, delta = mean(StudentsPerformance$goout) - 2.8, sd = sd(goout.sampled),
516   type="one.sample", alternative="two.sided")
517
518
519
520
521 differences <- seq(from = 0,to = 1.5,by = 0.1)
522 power.effect <- sapply(differences, function(d){power.t.test(n = 200, delta = d, sd = sd(
523   goout.sampled), type="one.sample")}$power)
524
525 df <- data.frame(differences, power.effect)
526
527 ggplot(data = df, aes(x = differences, y = power.effect)) + ylim(c(0, 1.2)) +
528   geom_line(linetype="dotted", color="thistle2", size=1)+ ylab("Power") + xlab("Effect size
529   ") +
530   geom_point(color="thistle3", size = 2)
531
532 #-----
533 #-----
534
535 #Question 7
536
537 #a. b)
538
539

```

```

540 StudentsPerformance.sampled <- sample_n(StudentsPerformance, 25)
541
542 Hypothesis.test <- function(data.sampled.var1, data.sampled.var2, null.value = 0, alpha =
    0.05, paired = FALSE){
543
544   sample.len <- length(data.sampled.var1)
545   print(paste("Null Hypothesis: diff mean = ", null.value))
546   print(paste("Alternative Hypothesis: diff mean != ", null.value))
547
548   x_bar <- mean(data.sampled.var1) - mean(data.sampled.var2)
549   s1 <- sd(data.sampled.var1)
550   s2 <- sd(data.sampled.var2)
551   if (paired)
552     SE <- sd(data.sampled.var1 - data.sampled.var2) / sqrt(sample.len)
553   else
554     SE <- sqrt((s1^2/sample.len) + (s2^2/sample.len))
555   score <- abs((x_bar - null.value)) / SE
556
557   if(sample.len > 30){
558     print("Using Z-distribution")
559     pvalue <- 2*pnorm(abs(score), lower.tail = FALSE)}
560
561   else{
562     print("Using t-distribution")
563     pvalue <- 2*pt(score, df = sample.len - 1, lower.tail = FALSE)}
564
565
566   print(paste("p-value =", pvalue))
567
568   if (pvalue < alpha)
569     print("Reject Null Hypothesis.")
570   else
571     print("Fail to Reject Null Hypothesis.")
572 }
573
574
575
576 Hypothesis.test(StudentsPerformance.sampled$health, StudentsPerformance.sampled$goout,
    paired = TRUE)
577
578 t.test(StudentsPerformance.sampled$health, StudentsPerformance.sampled$goout, paired = TRUE
    )
579
580 #————
581
582 #b.
583
584
585 idx.sampled <- sample(StudentsPerformance$X, 200)
586 health.sampled <- StudentsPerformance$health[idx.sampled[1:100]]
587 goout.sampled <- StudentsPerformance$goout[idx.sampled[1:100]]
588
589 Hypothesis.test(health.sampled, goout.sampled)
590
591 t.test(health.sampled, goout.sampled)
592
593 #————
594 #—————
595
596 #Question 8
597
598

```

```

599 absences_box <- ggplot(StudentsPerformance, aes(x = absences)) +
600   geom_boxplot(outlier.colour="thistle2", color ="gray77", fill ="gray77", alpha = 0.5,
601     outlier.size = 2) +
602   labs(title="Boxplot with outliers")
603 absences_box
604 boxplot.stats(StudentsPerformance$absences)
605
606
607
608 #a.
609
610 quantile(StudentsPerformance$absences, c(0.025, 0.975))
611
612
613
614 bs.size <- 1000
615 rep.size <- 1000
616
617 absences.sample <- replicate(1, sample(StudentsPerformance$absences, size = 200, replace =
618   FALSE))
619 absences.replicated <- replicate(rep.size, sample(absences.sample, size = 100, replace =
620   FALSE))
621
622 means <- apply(X = absences.replicated, MARGIN = 2, FUN = mean, na.rm = TRUE)
623
624
625 means <- sort(means)
626 margin <- 0.025 * bs.size
627
628 print(paste("Confidence Interval: (", round(means[c(margin)], 3),",",round(means[c(bs.size -
629   margin)],3),")"))
630 #————
631
632 #b.
633 bs.size <- 1000
634 rep.size <- 1000
635
636 absences.sample <- replicate(1, sample(StudentsPerformance$absences, size = 20, replace =
637   FALSE))
638 absences.bootstrapped <- replicate(rep.size, sample(absences.sample, size = 1000, replace =
639   TRUE))
640
641 means <- apply(X = absences.bootstrapped, MARGIN = 2, FUN = mean, na.rm = TRUE)
642
643
644 means <- sort(means)
645 margin <- 0.025 * bs.size
646
647 print(paste("Confidence Interval: (", round(means[c(margin)], 3),",",round(means[c(bs.size -
648   margin)],3),")"))
649 #————
650
651 #c.
652 absences.qq <- ggplot(StudentsPerformance, aes(sample = absences, color = "", alpha = 0.7))
653   + geom_qq() +
654   geom_qq_line() + labs(title="QQ-plot ")
655 absences.qq + theme(legend.position="none") + scale_color_manual(values=c("thistle2"))
656
657
658 m.absences.qq <- ggplot(data.frame(mean = means), aes(sample = means, color = "", alpha =
659   0.7)) + geom_qq() +
660   geom_qq_line() + labs(title="QQ-plot ")

```

```

652
653 m.absences.qq + theme(legend.position="none") + scale_color_manual(values=c("thistle2"))
654
655 #————
656 #
657
658
659 #Question 9
660
661
662 StudentsPerformance$Gsum <- StudentsPerformance$G1 + StudentsPerformance$G2 +
  StudentsPerformance$G3
663
664
665 f0.Gsum <- ((StudentsPerformance %>% filter(failures == 0))$Gsum)
666 f1.Gsum <- ((StudentsPerformance %>% filter(failures == 1))$Gsum)
667 f2.Gsum <- ((StudentsPerformance %>% filter(failures == 2))$Gsum)
668 f3.Gsum <- ((StudentsPerformance %>% filter(failures == 3))$Gsum)
669
670 sd.df <- data.frame(groups = c("Group0", "Group1", "Group2", "Group3"),
671                        sds = c(sd(f0.Gsum), sd(f1.Gsum), sd(f2.Gsum), sd(f3.Gsum)))
672
673
674
675 aov.Gsum_failures <- aov(Gsum ~ as.factor(failures), data = StudentsPerformance)
676 aov.Gsum_failures
677
678 summary(aov.Gsum_failures)
679
680
681
682 test1 <- lm(Gsum ~ failures, data = StudentsPerformance)
683 summary(test1)
684
685 TukeyHSD(aov.Gsum_failures)
686
687 plot(TukeyHSD(aov.Gsum_failures), las = 1)
688
689
690 sd(Gsum ~ as.factor(failures))
691
692
693
694 box <- ggplot(StudentsPerformance, aes(x = failures, y = Gsum, group = failures)) +
695   geom_boxplot(alpha = 0.5, outlier.size = 2, color = as.factor(failures), fill = as.factor(
696     failures)) +
697   labs(title = "Boxplot")
698
699 box + scale_color_manual(values=c("thistle1", "thistle2", "thistle3", "thistle4")) +
700   scale_fill_manual(values=c("thistle1", "thistle2", "thistle3", "thistle4"))

```

code.R

Statistical Inference: Project Phase II

Sarmad Zandi Goharrizi - 810199181

Question 0

Student's Performance includes various information about a sample of students studying in two different schools.

A sense of responsibility towards one's education and academic future is a notable information which can be mined from each individual's *study time* and their rate of *going out* which has an effect on their *failures* and their *grades*.

This dataset also contains some semi-relevant factors like each student's parent's job as well as their love life.

X	school	sex	age	Fjob	Mjob	goout	internet	romantic	studytime	failures	health	absences	G1	G2	G3
0	GP	F	18	teacher	at_home	4	no	no	2	0	3	6	5.000000	7.529856	9.289229
1	GP	F	17	other	at_home	3	yes	no	2	0	3	4	5.000000	7.192039	9.424835
2	GP	F	15	other	at_home	2	yes	no	2	3	3	10	3.807703	8.000000	7.354029
3	GP	F	15	services	health	2	yes	yes	3	0	5	2	15.000000	16.373208	17.796916
4	GP	F	16	other	other	2	no	no	2	0	5	4	6.000000	12.138542	12.800024
5	GP	M	16	other	services	2	yes	no	2	0	5	10	15.000000	16.804680	18.347259
6	GP	M	16	other	other	4	yes	no	2	0	3	0	12.000000	13.691091	14.187810
7	GP	F	17	teacher	other	4	no	no	2	0	1	6	6.000000	6.794185	9.012740
8	GP	M	15	other	services	2	yes	no	2	0	1	0	16.000000	19.852952	20.000000
9	GP	M	15	other	other	1	yes	no	2	0	5	0	14.000000	17.180466	18.073614
10	GP	F	15	health	teacher	3	yes	no	2	0	2	0	10.000000	9.609179	11.950918

Figure 1: Head of the dataset

Question 1

Chosen Categorical Variables : *sex* and *Mjob*

a.

In this part we intend to compare the proportion of mothers who are *teachers* between *Male* and *Female* students.

Conditions for inference for comparing two independent proportions :

- Independence :
 - random sample/assignment
 - if sampling without replacement, $n < 10\%$ of population
- Sample size / skew : samples should meet the success-failure condition (at least 10 *successes* and 10 *failures*) :
 - $n_1\hat{p}_1 \geq 10 \rightarrow n_1\hat{p}_1 = 200 \times 0.053 = 10.6 \geq 10$
 - $n_1(1 - \hat{p}_1) \geq 10 \rightarrow n_1(1 - \hat{p}_1) = 200 \times 0.947 = 189.4 \geq 10$
 - $n_2\hat{p}_2 \geq 10 \rightarrow n_2\hat{p}_2 = 200 \times 0.126 = 25.2 \geq 10$
 - $n_2(1 - \hat{p}_2) \geq 10 \rightarrow n_2(1 - \hat{p}_2) = 200 \times 0.874 = 174.8 \geq 10$

All is met.

Confidence Interval : *point estimate* \pm *margin of error* $\rightarrow \hat{p}_1 - \hat{p}_2 \pm z^* SE_{\hat{p}_1 - \hat{p}_2}$

$$SE_{\hat{p}_1 - \hat{p}_2} = \sqrt{\frac{\hat{p}_1(1 - \hat{p}_1)}{n_1} + \frac{\hat{p}_2(1 - \hat{p}_2)}{n_2}} = 0.047$$

Confidence Interval : (0.0655, 0.0811)

If we take repeated samples from this population, and make a confidence interval using each sample, we expect about 95% of the resulting confidence intervals to contain $\hat{p}_1 - \hat{p}_2$.

We are 95% confident that the difference of population proportion of *gMale* and *Female* students whose mother's job is *teachers* is between 0.0655 and 0.0811.

Other confidence intervals can be computed accordingly : mothers who are *at-home* between *Male* and *Female* students.

Confidence Interval : (-0.0632, -0.0567)

We are 95% confident that the difference of population proportion of *gMale* and *Female* students whose mother's job is *being home* is between -0.0632 and -0.0567.

mothers who are *health* between *Male* and *Female* students.

Confidence Interval : (-0.0239, -0.016)

We are 95% confident that the difference of population proportion of *gMale* and *Female* students whose mother's job is *health* is between -0.0239 and -0.016.

mothers who are *services* between *Male* and *Female* students.

Confidence Interval : $(-0.040, -0.019)$

We are 95% confident that the difference of population proportion of *gMale* and *Female* students whose mother's job is *services* is between -0.040 and -0.019 .

b.

To test the independence, I tested my hypothesis using 2 different methods :

H_0 : *Mother's job is independent from sex of the student. (Mother's job does not vary with the sex of the child)*

H_A : *Mother's job is dependent to sex of the student. (Mother's job varies with the sex of the child)*

Method 1 : Pooling :

$$p_{\hat{pool}} = \frac{\# \text{ success}}{\# \text{ total}} = 0.085$$

$$SE_{\hat{p}_1 - \hat{p}_2} = \sqrt{\frac{p_{\hat{pool}}(1 - p_{\hat{pool}})}{n_1} + \frac{p_{\hat{pool}}(1 - p_{\hat{pool}})}{n_2}} = 0.039$$

Conditions for inference for comparing two independent proportions (pooling) :

- **Independence** :
 - random sample/assignment
 - if sampling without replacement, $n < 10\%$ of population
- **Sample size / skew** : samples should meet the success-failure condition (at least 10 *successes* and 10 *failures*) :
 - $n_1 p_{\hat{pool}} \geq 10 \rightarrow n_1 p_{\hat{pool}} = 200 \times 0.085 = 17 \geq 10$
 - $n_1(1 - p_{\hat{pool}}) \geq 10 \rightarrow n_1(1 - p_{\hat{pool}}) = 200 \times 0.915 = 183 \geq 10$
 - $n_2 p_{\hat{pool}} \geq 10 \rightarrow n_2 p_{\hat{pool}} = 200 \times 0.085 = 17 \geq 10$
 - $n_2(1 - p_{\hat{pool}}) \geq 10 \rightarrow n_2(1 - p_{\hat{pool}}) = 200 \times 0.915 = 183 \geq 10$

All is met.

Due to the fact that p-value (0.105) is larger than 0.05 ,we fail to reject the null hypothesis. \rightarrow There are evidence that *Mother's job (teaching specifically) does not vary with the sex of the child* .

Method 2 : χ^2 test :

$$\text{Expected count} = \frac{\text{row total} \times \text{column total}}{\text{table total}}$$

$$\text{test statistic} : \frac{\text{point estimate} - \text{null value}}{\text{SE of point estimate}}$$

$$\chi^2 = \sum_{i=1}^k \frac{(O - E)^2}{E} \quad \text{and} \quad df = (R - 1)(C - 1)$$

Conditions for χ^2 test :

- Independence :
 - random sample/assignment
 - if sampling without replacement, $n < 10\%$ of population
 - each case only contributes to one cell in the table other (non-paired)
- Sample size : Each particular scenario (i.e.cell) must have atleast 5 expected cases. $\rightarrow \times$

Fjob					
sex	at_home	health	other	services	teacher
F	5	6	61	33	6
M	4	5	55	16	9

Figure 2: dataset table

Which will give out the following warning :

```
Chi-squared approximation may be incorrect
Pearson's Chi-squared test

data:  sp.sampled.table
X-squared = 4.6465, df = 4, p-value = 0.3255
```

Figure 3: χ^2 test

For our last condition to meet, we have to merge two columns, *at-home* and *health* :

		other	services	teacher
F	11	61	33	6
M	9	55	16	9

Figure 4: dataset table

So our χ^2 test won't give any warnings :

```
Pearson's Chi-squared test  
data:  sp.sampled.table.bind  
x-squared = 4.6445, df = 3, p-value = 0.1998
```

Figure 5: χ^2 test

Either way, due to the fact that p-value is larger than 0.05 , we fail reject the null hypothesis. \rightarrow There are evidence that *Mother's job does not vary with the sex of the child* .

Note :It's important to mention that in hypothesis testing in categorical variables, CI approach and p-value approach might not always give out the same result.

Question 2

Chosen Categorical Variable : *romantic*

$$H_0 : p = 0.5$$

$$H_A : p < 0.5$$

Conditions for inference for comparing two independent proportions :

- Independence :
 - random sample/assignment
 - if sampling without replacement, $n < 10\%$ of population
- Sample size / skew : samples should meet the success-failure condition (at least 10 *successes* and 10 *failures*) :
 - $n\hat{p} \geq 10 \rightarrow n\hat{p} = 15 \times 0.53 = 5.3 \not\geq 10$
 - $n(1 - \hat{p}) \geq 10 \rightarrow n(1 - \hat{p}) = 400 \times 0.47 = 4.7 \not\geq 10$

Due to the fact that our conditions did not meet, we will use simulation.

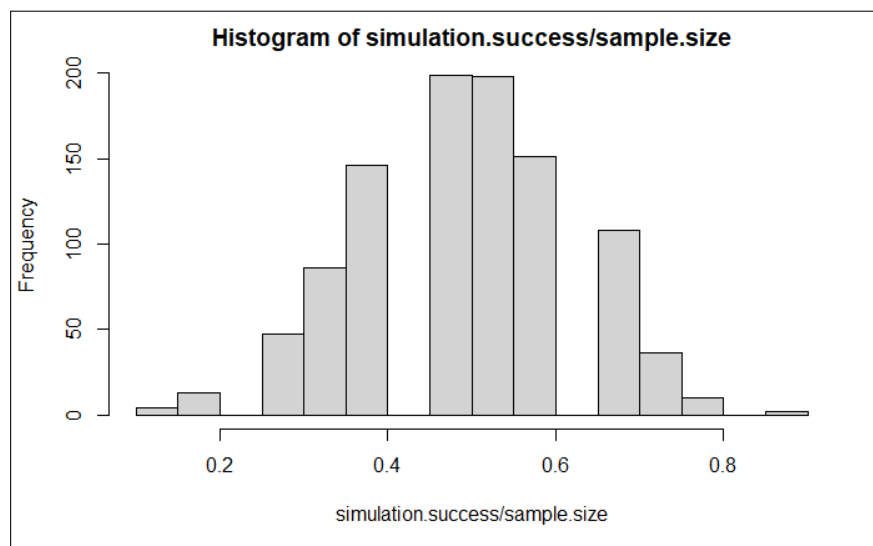


Figure 6: Histogram

Since, the p-value (0.505) is larger than 0.05, we fail to reject the null hypothesis and declare that there is not convincing evidence to accept the alternative hypothesis.

This means that each person is 50% likely to be in a romantic relationship.

Question 3

Chosen Categorical Variable : *Mjob*

sample.original				
at_home	health	other	services	teacher
59	34	141	103	58

Figure 7: Mjob

sample.original				
at_home	health	other	services	teacher
0.1494	0.0861	0.3570	0.2608	0.1468

Figure 8: Mjob - probability distribution

a.

Conditions for χ^2 test :

- Independence :
 - random sample/assignment
 - if sampling without replacement, $n < 10\%$ of population
 - each case only contributes to one cell in the table other (non-paired)
- Sample size : Each particular scenario (i.e.cell) must have atleast 5 expected cases.

H_0 : Samples are randomly chosen and there is nothing going on

H_0 : Samples are not randomly chosen and there is something going on

Randomly selected sample :

sample.unbiased				
at_home	health	other	services	teacher
16	9	39	18	18

Figure 9: 100 samples - randomly

χ^2 test :

Chi-squared test for given probabilities				
data: unbiased.table				
x-squared = 3.6496, df = 4, p-value = 0.4555				

Figure 10: χ^2 test - randomly

Due to the fact that p-value (0.455) is larger than 0.05, we fail to reject the null hypothesis. There is convincing evidence to accept the null hypothesis.

Randomly selected sample with 0.6 bias through teachers :

sample.biased				
at_home	health	other	services	teacher
10	12	33	21	24

Figure 11: 100 samples - biased

χ^2 test :

Chi-squared test for given probabilities	
data:	biased.table
x-squared =	10.071, df = 4, p-value = 0.03924

Figure 12: χ^2 test - biased

Due to the fact that p-value (0.0392) is smaller than 0.05, we fail to reject the null hypothesis, there is convincing evidence that the samples are randomly chosen. (!)

b.

Chosen Categorical Variable : *Fjob*

H_0 : Mother's job and father's job are 2 independent variables

H_0 : Mother's job and father's job are dependent variables

$$\text{Expected count} = \frac{\text{row total} \times \text{column total}}{\text{table total}}$$

$$\text{test statistic} : \frac{\text{point estimate} - \text{null value}}{\text{SE of point estimate}}$$

$$\chi^2 = \sum_{i=1}^k \frac{(O - E)^2}{E} \quad \text{and} \quad df = (R - 1)(C - 1)$$

Conditions for χ^2 test :

- Independence :
 - random sample/assignment
 - if sampling without replacement, $n < 10\%$ of population
 - each case only contributes to one cell in the table other (non-paired)
- Sample size : Each particular scenario (i.e.cell) must have atleast 5 expected cases. $\rightarrow \times$

Mjob	Fjob				
	at_home	health	other	services	teacher
at_home	2	2	21	7	0
health	0	4	8	3	0
other	4	0	51	8	4
services	2	2	24	18	5
teacher	1	3	12	13	6

Figure 13: table

Our last condition is not met, so we get a warning :

```
Chi-squared approximation may be incorrect
Pearson's Chi-squared test

data: Mjob.Fjob
X-squared = 44.517, df = 16, p-value = 0.0001645
```

Figure 14: χ^2 test

We combine *at-home* , *health*, *services* and *teacher* of *Mjob* and compare it to *other* :

	[,1]	[,2]
at_home	11	21
health	7	8
other	16	51
services	27	24
teacher	23	12

Figure 15: combined table

```
Pearson's Chi-squared test

data: Mjob.Fjob.bind
X-squared = 20.514, df = 4, p-value = 0.0003952
```

Figure 16: χ^2 test

Both p-values indicate that we should reject the null hypothesis meaning that parent's job are dependent to each other.

Question 4

Chosen Variables : $G1$ - *failure* and *studytime*

a.

In phase 1, we used *pearson correlation* and *Correlogram* for our predictions.

(In order not to confuse the report of this phase with the previous phase, the question related to phase 1 of the project is also placed in the zip file of this project, although its abstract is also described here.)

Quoting phase 1 : ‘Judging by Figure 34, $G1$ and $G2$ and $G3$ have positive linear associations with each other and with *studytime* as expected. *Failure* and *goout* both have a negative linear associations with $G1$, $G2$ and $G3$.’

From all the variables mentioned, I chose one of the grades ($G1$) as my response variable and from failures, goout and studytime I chose 2 of them that had the most correlation with $G1$ (absolute value of them are aimed).

(Note : Although $G2$ and $G3$ had a very high correlation with $G1$, I didn't pick them, because all three of these variables are scores in different classes and it is better to use other variables to better understand each person and do not estimate their $G1$ score only based on their other scores. (Each one of them can be a great response variable) Although in the end I built the model based on these two variables, because I do not know exactly what was the exact aim of this question, to choose only based on scores or not, simply because I myself thought a better model should be based on a student's other characteristics, I explained more about this.)

$$\text{cor}(G1, \text{failure}) = -0.463$$

$$\text{cor}(G1, \text{studytime}) = 0.176$$

$$\text{cor}(G1, \text{goout}) = -0.161$$

Using those codes here, judging by the results, we can say *failures* is the more significant predictor :

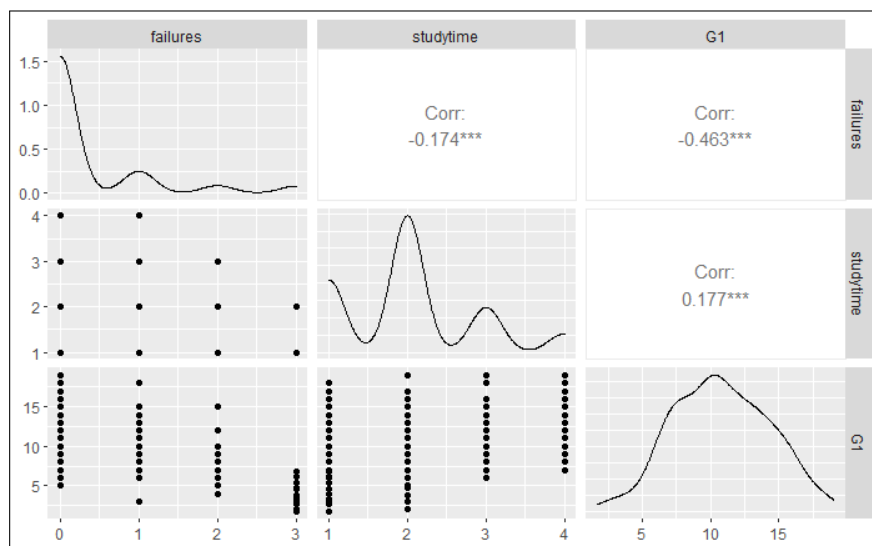


Figure 17: Correlogram

b.

Conditions for linear regression :

- **Residuals vs Fitted** : Used to check the linear relationship assumptions. A horizontal line, without distinct patterns is an indication for a linear relationship, what is good.
- **Normal Q-Q** : Used to examine whether the residuals are normally distributed. It's good if residuals points follow the straight dashed line.
- **Scale-Location** : (or Spread-Location). Used to check the homogeneity of variance of the residuals (homoscedasticity). Horizontal line with equally spread points is a good indication of homoscedasticity.
- **Residuals vs Leverage** : Used to identify influential cases, that is extreme values that might influence the regression results when included or excluded from the analysis.

failures :

```

Call:
lm(formula = G1 ~ failures, data = StudentsPerformance)

Residuals:
    Min       1Q   Median       3Q      Max
-6.5154 -2.5154 -0.5154  2.4846  8.6768

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  11.5154    0.1724   66.79  <2e-16 ***
failures     -2.1922    0.2117  -10.36  <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 3.125 on 393 degrees of freedom
Multiple R-squared:  0.2144,    Adjusted R-squared:  0.2124
F-statistic: 107.2 on 1 and 393 DF,  p-value: < 2.2e-16

```

Figure 18: LM model

$$R^2 = 0.214$$

$$p - value < 2.2e - 16$$

According to R^2 , 0.214 of the variability of the model is explained by failures.

According to the $p - value$, by modeling $G_1 \sim failures$, we can reject the null hypothesis that suggests there is no relationship between these two variables (slope is zero.)

```

Call:
lm(formula = G1 ~ failures, data = StudentsPerformance)

Coefficients:
(Intercept)      failures
   11.515         -2.192

```

Figure 19: LM model

$$G_1 = 11.515 - 2.192 \times failures$$

Intercept: When $failures = 0$, G_1 is expected to equal the intercept (11.515). (Maybe meaningless in context of the data, and only serve to adjust the height of the line.)

In our case when the student has not failed at all , their G_1 score is nearly 11 .

Slope: For each unit increase in *failures*, G_1 is expected to be 2.192 lower on average.

We also need to check whether conditions for using linear regression are met :

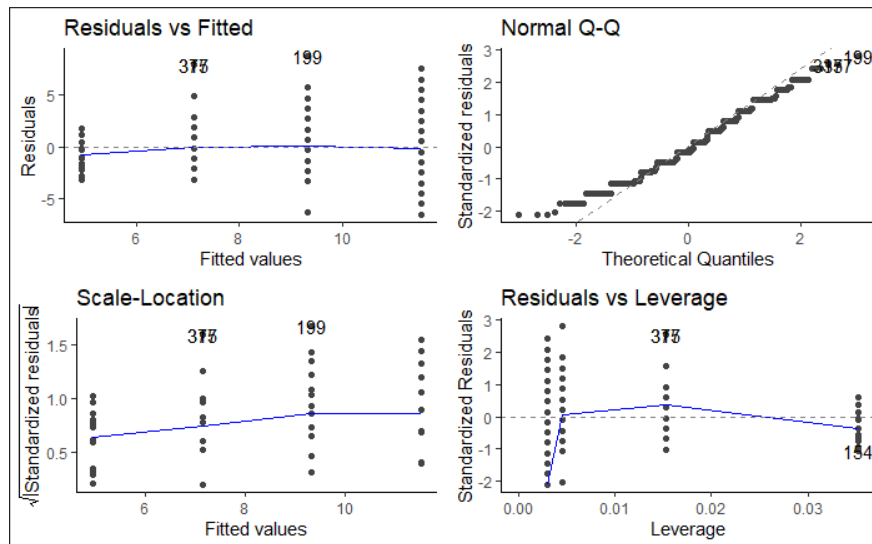


Figure 20: LM conditions - all met

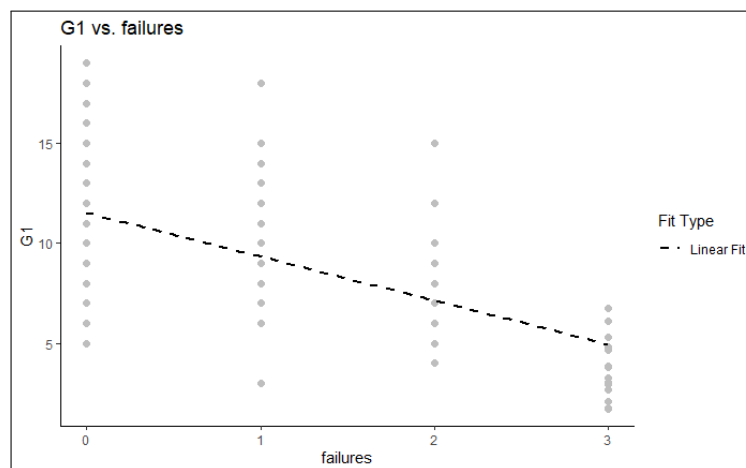


Figure 21: Scatter plot

studytime :

```

Call:
lm(formula = G1 ~ studytime, data = StudentsPerformance)

Residuals:
    Min       1Q   Median       3Q      Max
-8.6592 -2.7566 -0.0149  2.3726  8.2434

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)   9.2732     0.4585  20.224  < 2e-16 ***
studytime     0.7417     0.2083   3.561 0.000415 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 3.47 on 393 degrees of freedom
Multiple R-squared:  0.03125,    Adjusted R-squared:  0.02879
F-statistic: 12.68 on 1 and 393 DF,  p-value: 0.0004154

```

Figure 22: LM model

$$R^2 = 0.031$$

$$p - \text{value} = 0.00041$$

According to R^2 , only 0.03 of the variability of the model is explained by studytime (which is a lot smaller than failures).

According to the $p - \text{value}$, by modeling $G_1 \sim \text{studytime}$, we can reject the null hypothesis that suggests there is no relationship between these two variables (slope is zero.)

```

Call:
lm(formula = G1 ~ studytime, data = StudentsPerformance)

Coefficients:
(Intercept)    studytime
   9.2732         0.7417

```

Figure 23: LM model

$$G_1 = 9.2732 + 0.7417 \times \text{studytime}$$

intercept: When $\text{studytime} = 0$, G_1 is expected to equal the intercept (9.2732). Maybe meaningless in context of the data, and only serve to adjust the height of the line. In our case when the student does not study at all, their G_1 score is nearly 9.

slope: For each unit increase in studytime , G_1 is expected to be 0.7417 higher on average.

We also need to check whether conditions for using linear regression are met :

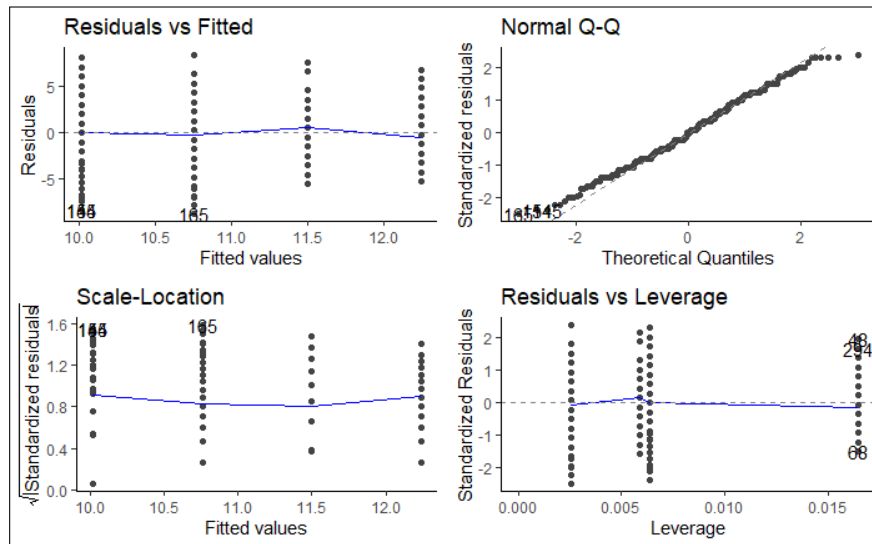


Figure 24: LM conditions - all met

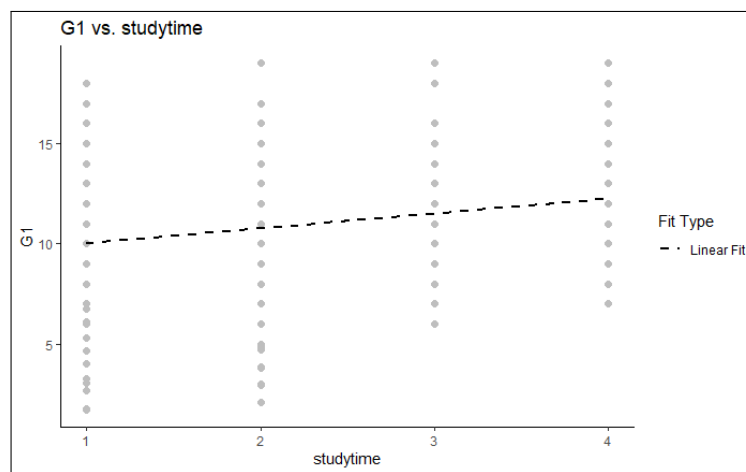


Figure 25: Scatter plot

c.

Judging by above figures, in order to pick the the more significant predictor we can use both R^2_{adj} and $p - value$:

	Adj. R-squared	p-value
failures	0.2124	2.2e-16
studytime	0.02879	0.00041

The more significant predictor is the one with the lowest $p - value$ and highest R^2_{adj} . Both of these point to failures being the best one.

Chosen Variables : $G1$ - $G2$ and $G3$

$$\text{cor}(G_1, G2) = 0.85$$

$$\text{cor}(G_1, G3) = 0.80$$

Using those codes here, judging by the results, we can say $G2$ is the more significant predictor :

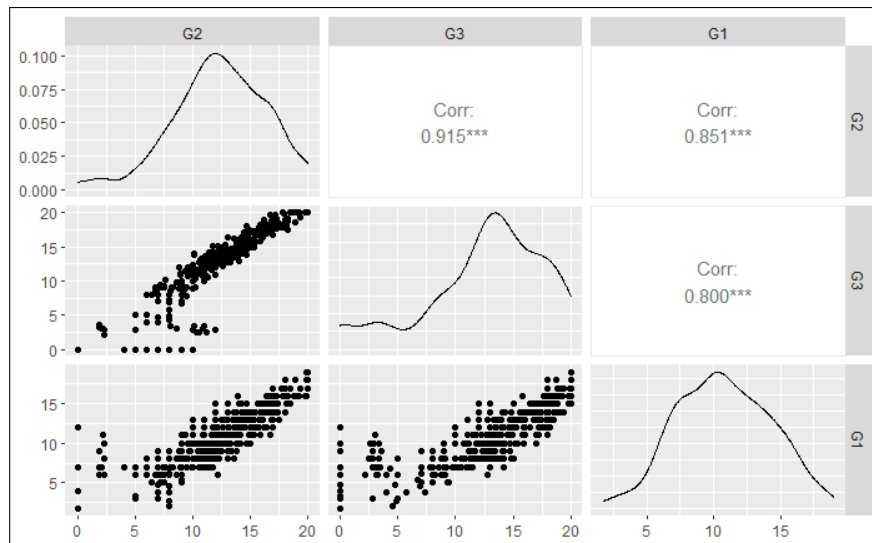


Figure 26: Correlogram

b.

Conditions for linear regression :

- **Residuals vs Fitted** : Used to check the linear relationship assumptions. A horizontal line, without distinct patterns is an indication for a linear relationship, what is good.
- **Normal Q-Q** : Used to examine whether the residuals are normally distributed. It's good if residuals points follow the straight dashed line.
- **Scale-Location** : (or Spread-Location). Used to check the homogeneity of variance of the residuals (homoscedasticity). Horizontal line with equally spread points is a good indication of homoscedasticity.
- **Residuals vs Leverage** : Used to identify influential cases, that is extreme values that might influence the regression results when included or excluded from the analysis.

G2 :

```

Call:
lm(formula = G1 ~ G2, data = StudentsPerformance)

Residuals:
    Min       1Q   Median       3Q      Max
-5.5525 -1.1545 -0.0471  1.0380 10.2153

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  1.78475    0.29527   6.045 3.49e-09 ***
G2           0.73313    0.02283  32.115 < 2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.852 on 393 degrees of freedom
Multiple R-squared:  0.7241,    Adjusted R-squared:  0.7234
F-statistic: 1031 on 1 and 393 DF,  p-value: < 2.2e-16

```

Figure 27: LM model

$$R^2 = 0.724$$

$$p - \text{value} < 2.2e - 16$$

According to R^2 , 0.724 of the variability of the model is explained by failures (which is pretty good).

According to the $p - \text{value}$, by modeling $G1 \sim G2$, we can reject the null hypothesis that suggests there is no relationship between these two variables (slope is zero.)

$$G1 = 1.7847 + 0.7331 \times G2$$

Intercept: When $G2 = 0$, $G1$ is expected to equal the intercept (1.7847). (Maybe meaningless in context of the data, and only serve to adjust the height of the line.)

In our case when the student has not failed at all , their G_1 score is nearly 1.78 .

Slope: For each unit increase in $G2$, $G1$ is expected to be 0.7331 higher on average.

We also need to check whether conditions for using linear regression are met :

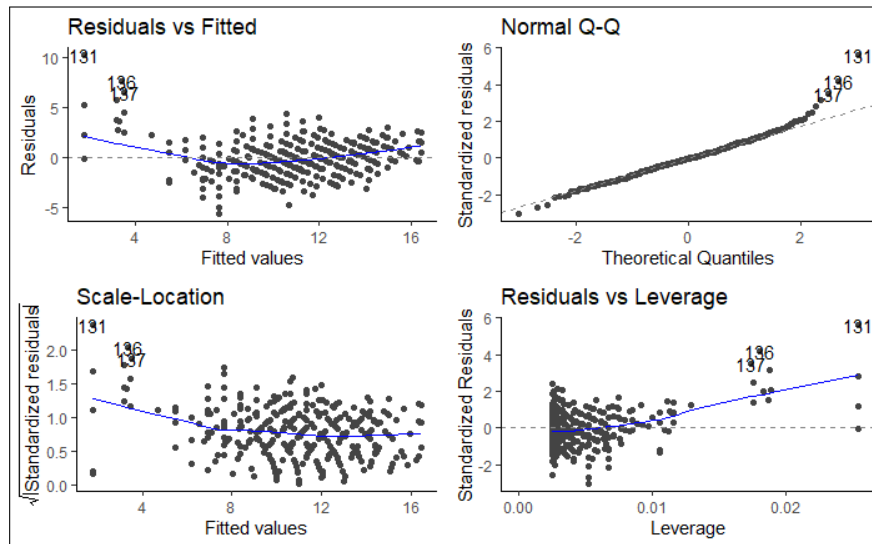


Figure 28: LM conditions - all are hardly met

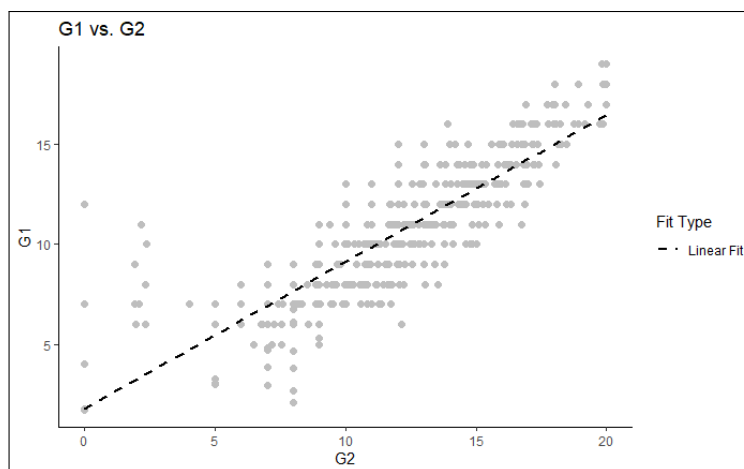


Figure 29: Scatter plot

G3 :

```
Call:
lm(formula = G1 ~ G3, data = StudentsPerformance)

Residuals:
    Min       1Q   Median       3Q      Max
-4.870 -1.623 -0.080  1.338  8.140

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)   3.8602     0.2825   13.66 <2e-16 ***
G3             0.5476     0.0207   26.45 <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 2.114 on 393 degrees of freedom
Multiple R-squared:  0.6403,    Adjusted R-squared:  0.6394
F-statistic: 699.6 on 1 and 393 DF,  p-value: < 2.2e-16
```

Figure 30: LM model

$$R^2 = 0.64$$

$$p - \text{value} < 2.2e - 16$$

According to R^2 , 0.64 of the variability of the model is explained by failures (which is pretty good).

According to the $p - \text{value}$, by modeling $G1 \sim G3$, we can reject the null hypothesis that suggests there is no relationship between these two variables (slope is zero.)

$$G1 = 3.860 + 0.5476 \times G3$$

Intercept: When $G3 = 0$, $G1$ is expected to equal the intercept (3.860). (Maybe meaningless in context of the data, and only serve to adjust the height of the line.)

In our case when the student has not failed at all, their $G3$ score is nearly 1.78.

Slope: For each unit increase in $G3$, $G1$ is expected to be 0.5476 higher on average.

We also need to check whether conditions for using linear regression are met :

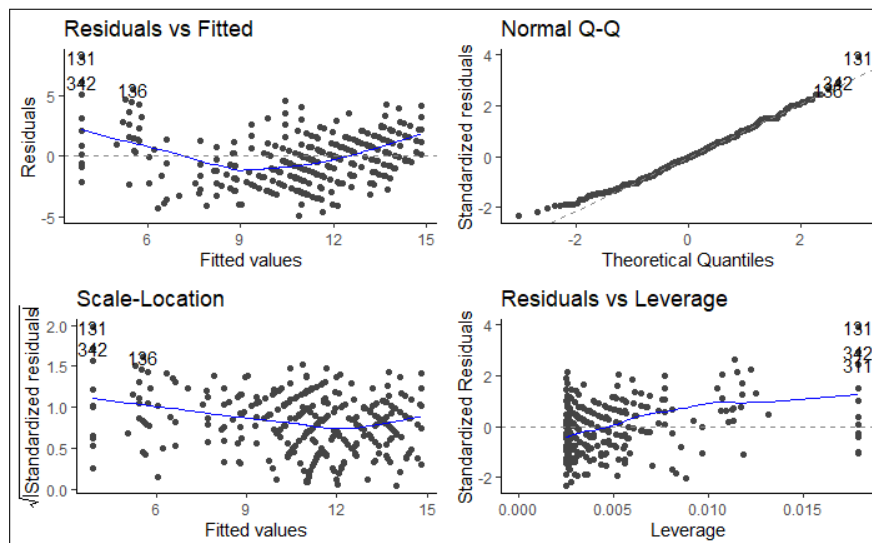


Figure 31: LM conditions - all are hardly met

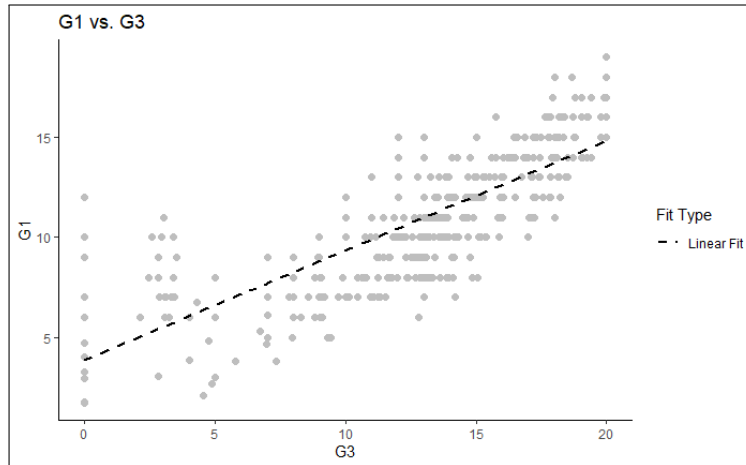


Figure 32: Scatter plot

	Adj. R-squared	p-value
G2	0.724	2.2e-16
G3	0.64	2.2e-16

The more significant predictor is the one with the lowest p -value and highest R_{adj}^2 . Although there is no difference in p -value, according to Adj. R-squared, G2 is the best one.

(Note : Between G2 and failures, G2 has a better R_{adj}^2 , but it did not meet the conditions very well. But R_{adj}^2 is more important so if we have to choose one variable, we choose G2)

d.

From this part forward , i will compare both of the models i made till now :

Adj. R-squared :

As was also mentioned in part c., Comparing failure vs. studytime using R_{adj}^2 will result in $G_1 \sim failures$ to be the better model.

As was also mentioned in part c., Comparing G2 vs. G3 using R_{adj}^2 will result in $G_1 \sim G2$ to be the better model.

As was also mentioned in part c., Comparing failure vs. G2 using R_{adj}^2 will result in $G_1 \sim G2$ to be the better model.

ANOVA table :

In order to compare my models, we first consider a base model, for example :

$$G1 \sim sex$$

Analysis of variance Table						
Response: G1						
	Df	Sum Sq	Mean Sq	F value	Pr(>F)	
sex	1	28.1	28.108	2.2745	0.1323	
Residuals	393	4856.6	12.358			

Figure 33: anove

failure vs. studytime :

Analysis of variance Table						
Response: G1						
	Df	Sum Sq	Mean Sq	F value	Pr(>F)	
sex	1	28.1	28.11	2.9056	0.08906	.
failures	1	1064.5	1064.54	110.0431	< 2e-16	***
Residuals	392	3792.1	9.67			

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1						

Figure 34: anove

Analysis of variance Table						
Response: G1						
	Df	Sum Sq	Mean Sq	F value	Pr(>F)	
sex	1	28.1	28.108	2.3741	0.1242	
studytime	1	215.6	215.641	18.2140	2.479e-05	***
Residuals	392	4641.0	11.839			

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1						

Figure 35: anove

$$R^2 = \frac{SS_{reg}}{SS_{total}}$$

	Base	+ failures	+ studytime
R2	0.01	0.22	0.05

Figure 36: computed R2 base on anova

Comparing failure vs. studytime using R^2 will result in $G_1 \sim failures$ to be the better model. G_2 vs. G_3 :

Analysis of variance Table						
Response: G1						
	Df	Sum Sq	Mean Sq	F value	Pr(>F)	
sex	1	28.1	28.1	8.1791	0.004464	**
G2	1	3509.5	3509.5	1021.2354	< 2.2e-16	***
Residuals	392	1347.1	3.4			

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1						

Figure 37: anove

Analysis of variance Table						
Response: G1						
	Df	Sum Sq	Mean Sq	F value	Pr(>F)	
sex	1	28.11	28.11	6.2773	0.01263	*
G3	1	3101.39	3101.39	692.6334	< 2e-16	***
Residuals	392	1755.25	4.48			

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1						

Figure 38: anove

$$R^2 = \frac{SS_{reg}}{SS_{total}}$$

Analysis of Variance Table						
Response: G1						
	Df	Sum Sq	Mean Sq	F value	Pr(>F)	
sex	1	28.11	28.11	6.2773	0.01263	*
G3	1	3101.39	3101.39	692.6334	< 2e-16	***
Residuals	392	1755.25	4.48			
--- Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1						

Figure 39: computed R^2 base on anova

Comparing G2 vs. G3 using R^2 will result in $G_1 \sim G_2$ to be the better model.

Due to the fact that $n - 1$ and $n - k - 1$ are approximately the same, R^2 and Adjusted R^2 doesn't have a noticeable difference.

e.

When there are many possible predictors, we need some strategy for selecting the best predictors to use in a regression model.

- Adjusted R^2 : Under this criterion, the best model is the one with the highest value of Adjusted R^2 .
- Cross-validation (explained in the next question) :Under this criterion, the best model is the one with the smallest value of MSE.
- Corrected Akaike's Information Criterion :Under this criterion, the best model is the one with the smallest value of AIC.
- Schwarz's Bayesian Information Criterion :Under this criterion, the best model is the one with the smallest value of BIC.

While R^2 is widely used, and has been around longer than the other measures, its tendency to select too many predictor variables makes it less suitable for forecasting.

Many statisticians like to use the BIC because it has the feature that if there is a true underlying model, the BIC will select that model given enough data. However, in reality, there is rarely, if ever, a true underlying model, and even if there was a true underlying model, selecting that model will not necessarily give the best forecasts (because the parameter estimates may not be accurate).

f.

H_0 : The explanatory variable is not a significant predictor of the response variable, i.e. no relationship $\rightarrow \beta = 0$

H_A : The explanatory variable is a significant predictor of the response variable, i.e. relationship $\rightarrow \beta \neq 0$

(a)

	p-value	significant
studytime	0.385	×
failure	2.2e-16	✓
G2	2.2e-16	✓
G3	2.2e-16	✓

Comparing failure vs. studytime : failure is significant predictor of the response variable.

Comparing G2 vs. G3 : Both are significant predictor of the response variable.

Comparing failure vs. G2 : Both are is significant predictor of the response variable.

(b)

$$\text{failure } CI : (-2.448, -1.769)$$

We are 95% confident that for each additional point on failure, G1 is expected on average to be lower by 1.769 to 2.448 points.

$$\text{studytime } CI : (0.019, 0.845)$$

We are 95% confident that for each additional point on studytime, G1 is expected on average to be higher by 0.019 to 0.845 points.

$$G2 \text{ } CI : (0.613, 0.698)$$

We are 95% confident that for each additional point on G2, G1 is expected on average to be higher by 0.613 to 0.846985 points.

$$G3 \text{ } CI : (0.530, 0.565)$$

We are 95% confident that for each additional point on G3, G1 is expected on average to be higher by 0.530 to 0.565 points.

(c)

	Actual	Predicted studytime	Predicted failues	Predicted G2	Predicted G3
209	9	10.2	11.5	10.2	10.9
244	13	10.2	11.5	11.9	12.2
6	15	10.6	11.5	13.8	13.9
358	12	10.6	11.5	12.0	11.6
178	6	10.6	11.5	7.3	8.7
44	8	10.2	11.5	9.4	11.5
201	16	10.6	11.5	14.6	14.3
82	11	11.0	11.5	10.6	11.4
295	14	11.0	11.5	12.4	13.2
30	10	10.6	11.5	11.8	11.5

Figure 40: Predicted

(d)

	Actual	Predicted studytime	Predicted failues	Predicted G2	Predicted G3
209	0	1.2	2.5	1.2	1.9
244	0	2.8	1.5	1.1	0.8
6	0	4.4	3.5	1.2	1.1
358	0	1.4	0.5	0.0	0.4
178	0	4.6	5.5	1.3	2.7
44	0	2.2	3.5	1.4	3.5
201	0	5.4	4.5	1.4	1.7
82	0	0.0	0.5	0.4	0.4
295	0	3.0	2.5	1.6	0.8
30	0	0.6	1.5	1.8	1.5

Figure 41: Prediction error

In order to compute success rate, I accept an 0.1 error which is 2 (data range \times accepted error = $20 \times 0.1 = 2$).

If the predicted result is ± 2 of my actual value, I accept it.

	Predicted studytime	Predicted failues	Predicted G2	Predicted G3
1	40 %	40 %	100 %	80 %

Figure 42: Success rate

Comparing failure vs. studytime : no difference !

Comparing G2 vs. G3 : G2 is the best predictor of the response variable.

Comparing failure vs. G2 : G2 is the best predictor of the response variable.

G2 is by far the best predictor.

Using *min-max calculation* :

$$MinMaxAccuracy = mean\left(\frac{\min(actual, predicted)}{\max(actual, predicted)}\right)$$

	Predicted studytime	Predicted failues	Predicted G2	Predicted G3
1	79.99 %	79.76 %	90 %	86.93 %

Figure 43: Success rate

Comparing failure vs. studytime :studytime is the best predictor of the response variable.

Comparing G2 vs. G3 : G2 is the best predictor of the response variable.

Comparing failure vs. G2 : G2 is the best predictor of the response variable.

G2 is by far the best predictor.

Using *MAPE* calculation :

$$MeanAbsolutePercentageError = mean\left(\frac{abs(actual - predicted)}{actual}\right)$$

	Predicted studytime	Predicted failues	Predicted G2	Predicted G3	
1	79.99 %	79.76 %	90 %	86.93 %	86.93 %

Figure 44: Success rate

Comparing failure vs. studytime :studytime is the best predictor of the response variable.

Comparing G2 vs. G3 : G2 is the best predictor of the response variable.

Comparing failure vs. G2 : G2 is the best predictor of the response variable.

Question 5

Chosen Categorical Variables : $G1 - G2$, $goout$, sex , $failures$, age and $studytime$

a.

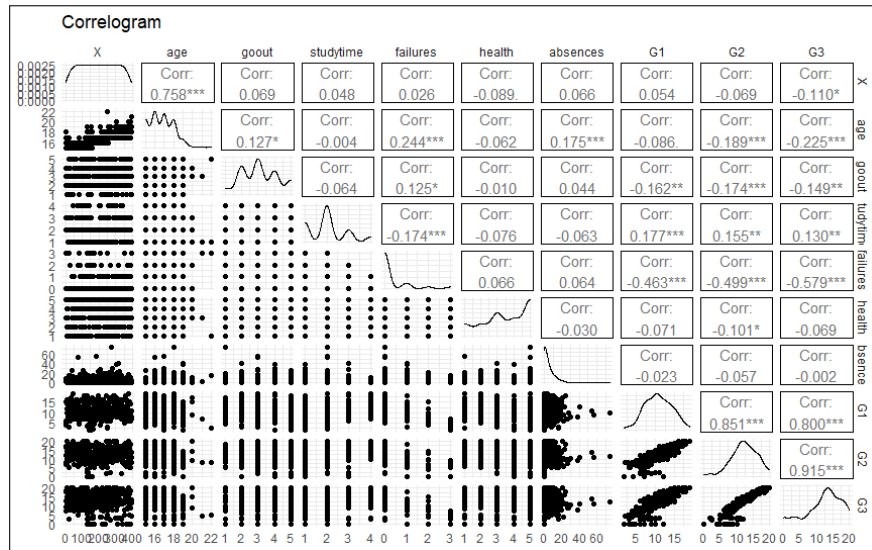


Figure 45: Correlogram

Considering all the analysis and inferences made in the previous question, for this question, from all the choices we have, we will definitely put $G2$ and failure among our options. Adding both $G2$ and $G3$ won't add anything new to the table since they are collinear.

We will add $studytime$, $goout$, age , and sex too, but we have to be careful not to use too many variables; we should pay attention to *occam's razor*; prefer the simplest best model!

The correlation between variables was also explained in the last question, but to summarize, as it was expected, failures and $G2$ have the highest correlation. Surprisingly, age has a higher correlation with $G2$ than $studytime$ (considering their absolute value).

Age and sex are not correlated, just as sex and $G2$. sex has nothing to do with score or age, so it was probably expected.

More explanations can be found on phase 1, so to avoid lengthening the report, we move on to the next part. $Mjob$ and $Fjob$ seemed hardly important to $G1$, so I didn't use them.

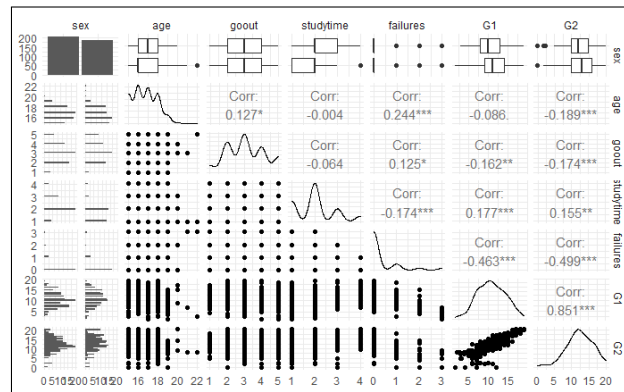


Figure 46: Correlogram

We don't want any collinearity between the variables that we chose and the below figure shows that we did a good job :

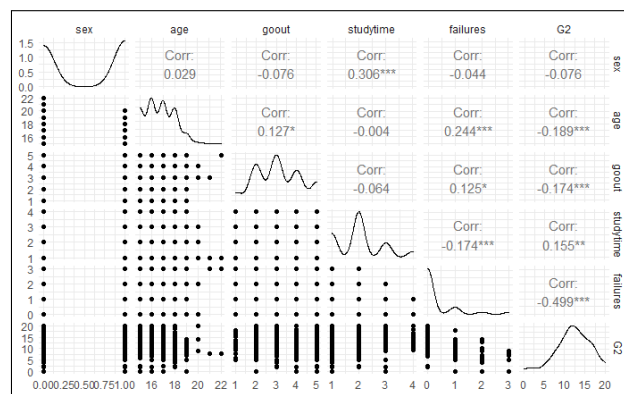


Figure 47: Correlogram - response variable omitted

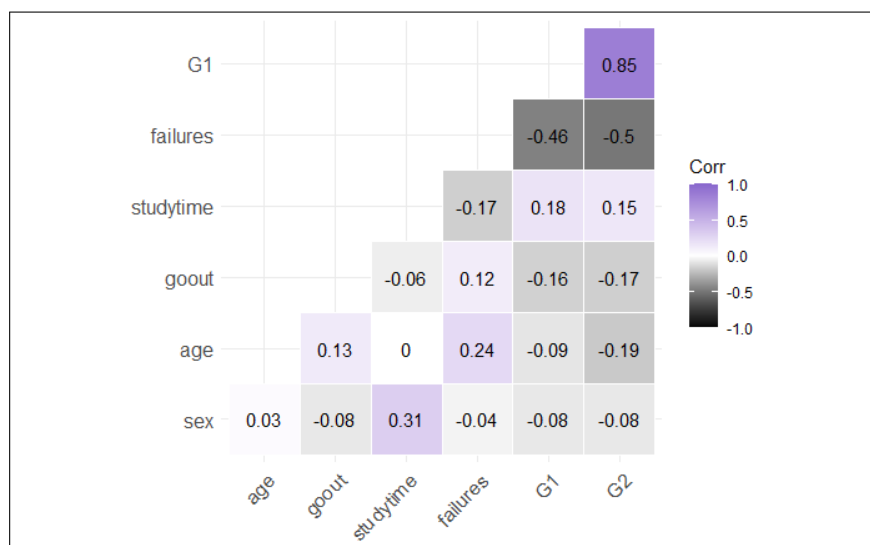


Figure 48: Correlogram

As was said multiple times in the previous question, $G2$ plays a more significant role in prediction.

b.

```
Call:
lm(formula = G1 ~ G2 + goout + failures + studytime + sex + age,
    data = StudentsPerformance)

Residuals:
    Min       1Q   Median       3Q      Max
-5.3720 -1.1898 -0.1367  1.0890 10.6786

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) -1.97651    1.33371   -1.482   0.1392
G2           0.70774    0.02661  26.599 <2e-16 ***
goout        -0.06969    0.08454   -0.824   0.4102
failures     -0.30731    0.14614   -2.103   0.0361 *
studytime     0.20212    0.11755    1.719   0.0863 .
sex          -0.24841    0.19572   -1.269   0.2051
age           0.24627    0.07487    3.289   0.0011 **
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.823 on 388 degrees of freedom
Multiple R-squared:  0.7361,    Adjusted R-squared:  0.732
F-statistic: 180.4 on 6 and 388 DF,  p-value: < 2.2e-16
```

Figure 49: Correlogram

$$G1 = -1.97 + 0.7 \times G2 - 0.06 \times goout + -0.3 \times failures + 0.2 \times studytime - 0.24 \times sex : M + 0.24 \times age$$

c.

R^2 shows what percent of variability in the response variable is explained by the model. In our case nearly 74% of the variability was explained using 6 variables out of the 15 variables available which is pretty good.

d.

Higher R^2 doesn't necessarily guarantee that the model fits the data well, we might face over-fitting if we are not careful.

Adjusted R^2 can be a good indicator of when the model fits the data well, it compares the explanatory power of regression models that contain different numbers of predictors. Adjusted R^2 is around 73% in our fitted model.

The fact that R^2 and Adjusted R^2 are this close is very good which means we don't have overfitting in our model.

Other techniques can help us know whether we have a good fit or not, for example Residuals : A good way to test the quality of the fit of the model is to look at the residuals The idea in here is that the sum of the residuals is approximately zero or as low as possible.

Although all our analyzes so far have promised a good model, the figure below shows that the model is not the best possible model.

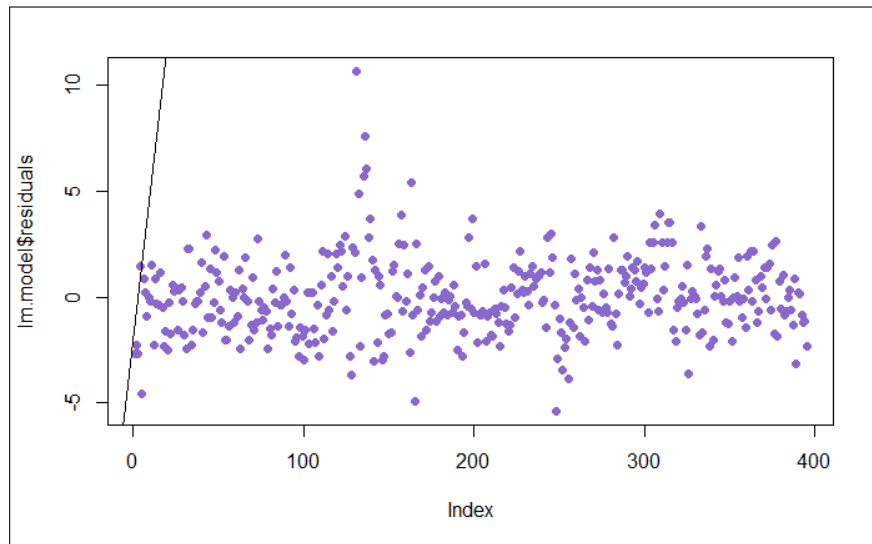


Figure 50: residuals

e.

To develop the best possible model, there are 4 different approaches :

Forward selection : start with an empty model and add one predictor at a time until the parsimonious model is reached.

(a) *p-value* :

Start with single predictor regressions of response vs. each explanatory variable (G2, G3 and failure all had p-values smaller than 2.2×10^{-16} , doesn't matter which one we will choose)

Pick the variable with the lowest significant p-value

Add the remaining variables one at a time to the existing model, and pick the variable with the lowest significant p-value

Repeat until any of the remaining variables do not have a significant p-value

For this part I used `ols_step_forward_p`, details can also be found in my project's file.

Variables Entered:							
+ G2							
+ age							
+ failures							
+ studytime							
+ sex							
Final Model Output							

Model Summary							

R	0.858	RMSE	1.822				
R-Squared	0.736	Coef. Var	16.896				
Adj. R-Squared	0.732	MSE	3.319				
Pred R-Squared	0.722	MAE	1.361				

RMSE: Root Mean Square Error							
MSE: Mean Square Error							
MAE: Mean Absolute Error							
ANOVA							

	Sum of Squares	DF	Mean Square	F	Sig.		
Regression	3593.549	5	718.710	216.526	0.0000		
Residual	1291.200	389	3.319				
Total	4884.749	394					

Parameter Estimates							

model	Beta	Std. Error	Std. Beta	t	Sig.	lower	upper
(Intercept)	-2.139	1.319		-1.622	0.106	-4.731	0.453
G2	0.711	0.026	0.825	26.946	0.000	0.659	0.762
age	0.240	0.075	0.087	3.227	0.001	0.094	0.387
failures	-0.309	0.146	-0.065	-2.119	0.035	-0.597	-0.022
studytime	0.203	0.117	0.048	1.728	0.085	-0.028	0.434
sex	-0.235	0.195	-0.033	-1.206	0.229	-0.618	0.148

Figure 51: Forward - p-value

Final model :

$$G1 \sim G2 + age + failures + sex + studytime$$

(b) *Adjusted R^2* :

Start with single predictor regressions of response vs. each explanatory variable. Pick the model with the highest adjusted R^2 .

Add the remaining variables one at a time to the existing model, and pick the model with the highest adjusted R^2 .

Repeat until the addition of any of the remaining variables does not result in a higher adjusted R^2 .

	best.pred	all.adj.r.squared
1	G2	0.7507277
2	G2 + age	0.7630615
3	G2 + age + failures	0.7639030
4	G2 + age + failures + studytime	0.7635379

Figure 52: Forward - Adjusted R^2

As we can see, adding the fourth variable, studytime, didn't help us with gaining a better fit for our model, so we wrapped this approach up after obtaining this model :

$$G1 \sim G2 + age + failures$$

with Adjusted $R^2 \approx 73.4\%$

Backward elimination : start with a full model (containing all predictors), drop one predictor at a time until the parsimonious model is reached.

(a) *p-value* :

Start with the full model

Drop the variable with the highest p-value and refit a smaller model

Repeat until all variables left in the model are significant

For this part i used *ols_step_backward_p*, details can also be found in my project's file.

Variables Removed:

x goout

x sex

x studytime

Final Model Output

Model Summary

R	0.856	RMSE	1.825
R-Squared	0.733	Coef. Var	16.928
Adj. R-Squared	0.731	MSE	3.332
Pred R-Squared	0.723	MAE	1.363

RMSE: Root Mean Square Error

MSE: Mean Square Error

MAE: Mean Absolute Error

ANOVA

	Sum of Squares	DF	Mean Square	F	Sig.
Regression	3582.019	3	1194.006	358.368	0.0000
Residual	1302.730	391	3.332		
Total	4884.749	394			

Parameter Estimates

model	Beta	Std. Error	Std. Beta	t	Sig.	lower	upper
(Intercept)	-1.994	1.311		-1.521	0.129	-4.573	0.584
G2	0.718	0.026	0.834	27.563	0.000	0.667	0.769
failures	-0.323	0.145	-0.068	-2.225	0.027	-0.608	-0.038
age	0.244	0.075	0.088	3.270	0.001	0.097	0.390

Figure 53: Backward - p-value

Final model :

$$G1 \sim G2 + age + failures$$

Forward and backward p-value approach did not gain the same result.

(b) *Adjusted R^2* :

Start with the full model

Drop one variable at a time and record adjusted R^2 of each smaller model

Pick the model with the highest increase in adjusted R^2

Repeat until none of the models yield an increase in adjusted R^2

	best.pred	all.adj.r.squared
1	G2 + failures + studytime + sex + age	0.7630107
2	G2 + failures + studytime + age	0.7635379
3	G2 + failures + age	0.7639030

Figure 54: Backward - Adjusted R^2

As we can see, omitting either $G2$, failure nor age, didn't help us with gaining a better fit for our model and increasing our Adjusted R^2 , so we wrapped this approach up after obtaining this model :

$$G1 \sim G2 + age + failures$$

with Adjusted $R^2 \approx 73.4\%$

Surprisingly, both forward and backward Adjusted R^2 gained the same result.

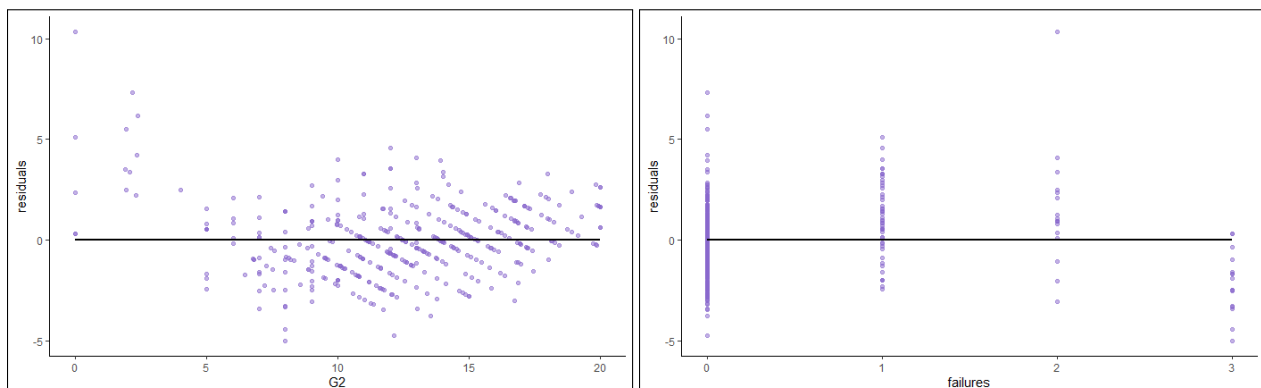
After completing these 4 methods, we see that both Adjusted R^2 approaches and the backward elimination pvalue approach all came to the same result, which is different from the result reached by forward selection pvalue.

According to the criterias mentioned in part d, $G1 \sim G2 + age + failures$ model is better than $G1 \sim G2 + age + failures + studytime + sex$, so we will use the same model in the following parts.

f.

Conditions for linear regression :

- **Linear relationships between x and y :**
Each (numerical) explanatory variable linearly related to the response variable
Check using residuals plots (e vs. x)
Looking for a random scatter around 0
Instead of scatterplot of y vs. x : allows for considering the other variables that are also in the model, and not just the bivariate relationship between a given x and y
- **Nearly normal residuals :**
Some residuals will be positive and some negative
On a residuals plot we look for random scatter of residuals around 0
This translates to a nearly normal distribution of residuals centered at 0
Check using histogram or normal probability plot
- **Constant variability of residuals :**
Residuals should be equally variable for low and high values of the predicted response variable
Check using residuals plots of residuals vs. predicted (e vs. y)
Residuals vs. predicted instead of residuals vs. x because it allows for considering the entire model (with all explanatory variables) at once
Residuals randomly scattered in a band with a constant width around 0 (no fan shape)
Also worthwhile to view absolute value of residuals vs. predicted to identify unusual observations easily



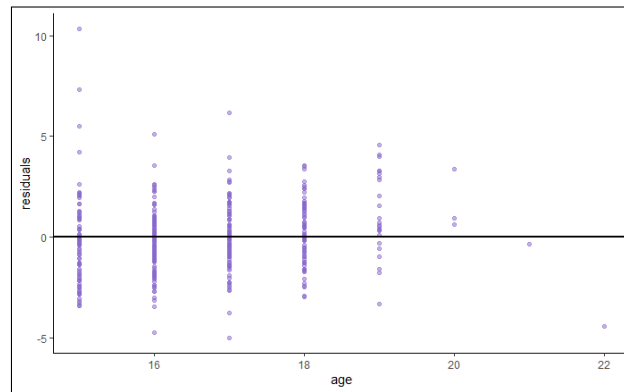


Figure 55: Linear relationship

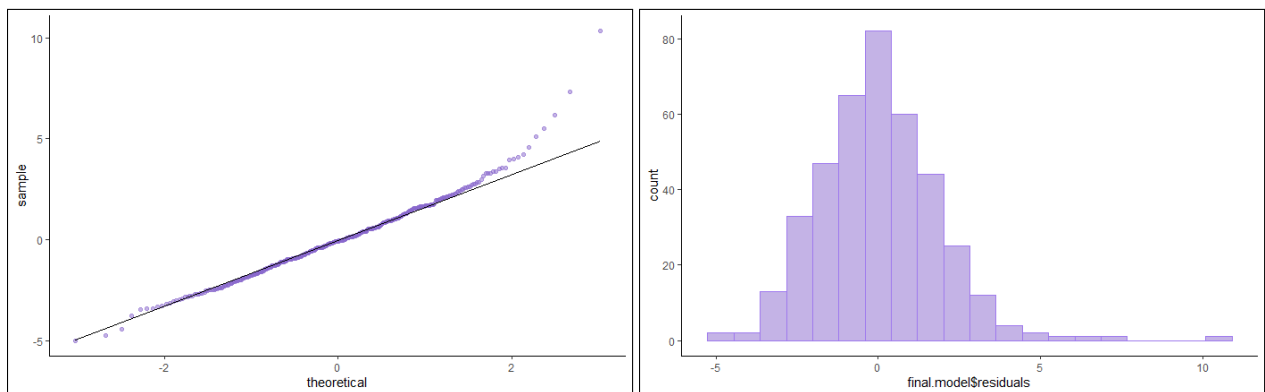


Figure 56: Nearly normal residuals

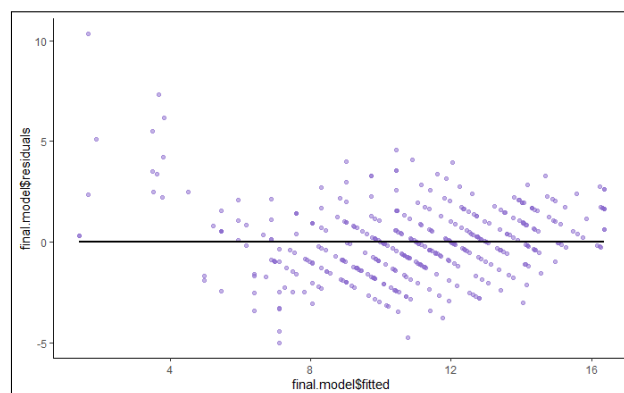


Figure 57: Constant var.

All conditions are met, not fully and perfectly, but they are met.

There are outliers that are effecting our model, if we eliminate them, we will meet them perfectly.

More on ourliers and detecting them in conditions and plots below :

(Note : this conditions were not mentioned in the slides, i checked them too just to be sure, some of them might overlap with the prev. conditions) Conditions for linear regression :

- **Residuals vs Fitted** : Used to check the linear relationship assumptions. A horizontal line, without distinct patterns is an indication for a linear relationship, what is good.
- **Normal Q-Q** : Used to examine whether the residuals are normally distributed. It's good if residuals points follow the straight dashed line.
- **Scale-Location** : (or Spread-Location). Used to check the homogeneity of variance of the residuals (homoscedasticity). Horizontal line with equally spread points is a good indication of homoscedasticity.
- **Residuals vs Leverage** : Used to identify influential cases, that is extreme values that might influence the regression results when included or excluded from the analysis.

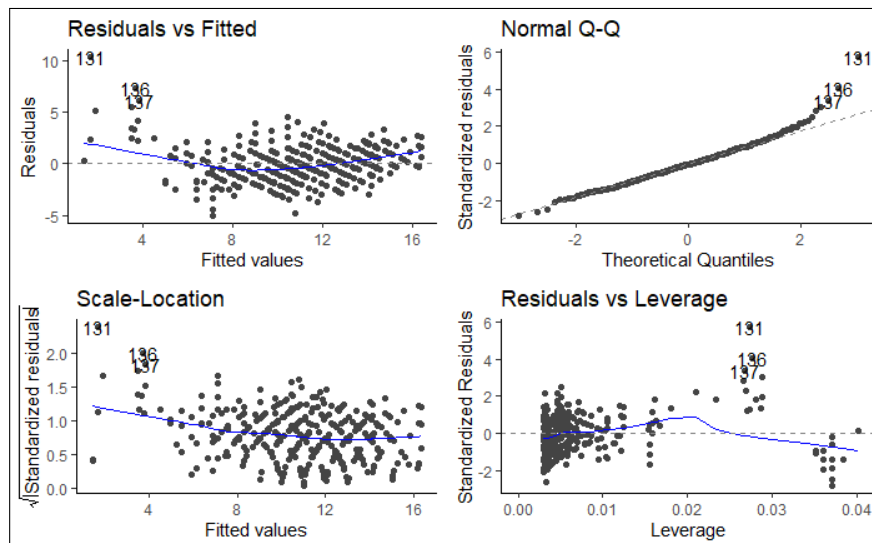


Figure 58: Conditions

All in all, we have a reliable model !

g.

The basic idea, behind cross-validation techniques, consists of dividing the data into two sets:

The training set, used to train (i.e. build) the model; and the testing set (or validation set), used to test (i.e. validate) the model by estimating the prediction error. Cross-validation is also known as a resampling method because it involves fitting the same statistical method multiple times using different subsets of the data.

The k-fold cross-validation method evaluates the model performance on different subset of the training data and then calculate the average prediction error rate. The algorithm is as follow:

- Randomly split the data set into k-subsets (or k-fold) (for example 5 subsets)
- Reserve one subset and train the model on all other subsets
- Test the model on the reserved subset and record the prediction error
- Repeat this process until each of the k subsets has served as the test set.
- Compute the average of the k recorded errors. This is called the cross-validation error serving as the performance metric for the model.

Root Mean Squared Error, which measures the model prediction error. It corresponds to the average difference between the observed known values of the outcome and the predicted value by the model. The lower the RMSE, the better the model.

```
Linear Regression

395 samples
  6 predictor

No pre-processing
Resampling: Cross-Validated (5 fold)
Summary of sample sizes: 317, 315, 316, 316, 316
Resampling results:

   RMSE      Rsquared   MAE
1.85564    0.7268193  1.389432

Tuning parameter 'intercept' was held constant at a value of TRUE
```

Figure 59: Full model

```
Linear Regression

395 samples
  3 predictor

No pre-processing
Resampling: Cross-Validated (5 fold)
Summary of sample sizes: 316, 316, 315, 317, 316
Resampling results:

   RMSE      Rsquared   MAE
1.836361    0.7357864  1.379812

Tuning parameter 'intercept' was held constant at a value of TRUE
```

Figure 60: Best model

Due to the fact that RMSE is lower in the so called 'Best model', we trust what we have done till now.

	RMSE	Rsquared	MAE	Resample
1	1.637994	0.7763565	1.284930	Fold1
2	2.069837	0.6892742	1.364075	Fold2
3	1.712149	0.7383289	1.322173	Fold3
4	1.848170	0.7410741	1.559756	Fold4
5	1.912499	0.7349491	1.378221	Fold5

Figure 61: Different metrics of all 5-fold , best model

Question 6

Chosen Variables : *catG3 - failures, studytime, G2 and sex*

Due to the fact that my dataset lacked a good binary categorical variable, i made G_3 into a binary categorical variable \rightarrow if $G_3 < 10$: *Fail*(0) else *Pass*(1)

a.

```
call:
glm(formula = catG3 ~ failures + studytime + G2 + sex, family = binomial(link = "logit"),
    data = StudentsPerformance)

Deviance Residuals:
    Min       1Q   Median       3Q      Max
-2.81360  0.00018  0.01265  0.13763  2.19874

Coefficients:
            Estimate Std. Error z value Pr(>|z|)
(Intercept) -14.5388    2.4952  -5.827 5.65e-09 ***
failures     -0.7129    0.3903  -1.827  0.0678 .
studytime    -0.2745    0.3209  -0.856  0.3922
G2           1.6466    0.2553   6.449 1.13e-10 ***
sex          -0.4254    0.5741  -0.741  0.4587
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 433.50  on 394  degrees of freedom
Residual deviance: 102.11  on 390  degrees of freedom
AIC: 112.11

Number of Fisher Scoring iterations: 8
```

Figure 62: GLM model

$$\log\left(\frac{p}{1-p}\right) = -14.538 - 0.712 \times \text{failures} - 0.2745 \times \text{studytime} + 1.646 \times G2 + -0.42 \times \text{sex} : M$$

intercept : keeping all other predictors zero, the log odds ratio / odds ratio of catG3 is -14.538 / exp(-14.538) = 4.85e-7

failures : keeping all other predictors constant for a unit increase in failures, the log odds ratio / odds ratio of catG3 will decrease -0.712 / exp(-0.712) = 0.49

studytime : keeping all other predictors constant for a unit increase in studytime, the log odds ratio / odds ratio of catG3 will decrease -0.2745 / exp(-0.2745) = 0.763

G2 : keeping all other predictors constant for a unit increase in G2 , the log odds ratio / odds ratio of catG3 will increase 1.6462 / exp(1.646) = 5.15

sex : keeping all other predictors constant, the log odds ratio / odds ratio of catG3 for reference point (M) is - 0.42 / exp(-0.712) = 0.657 less than F

b.

Odds ratio (OR) is a statistic that quantifies the strength of the association between two events, A and B. The odds ratio is defined as the ratio of the odds of A in the presence of B and vice versa, which, due to symmetry, is equal to the ratio of the odds of B in the presence of A and the odds of B in the absence of A.

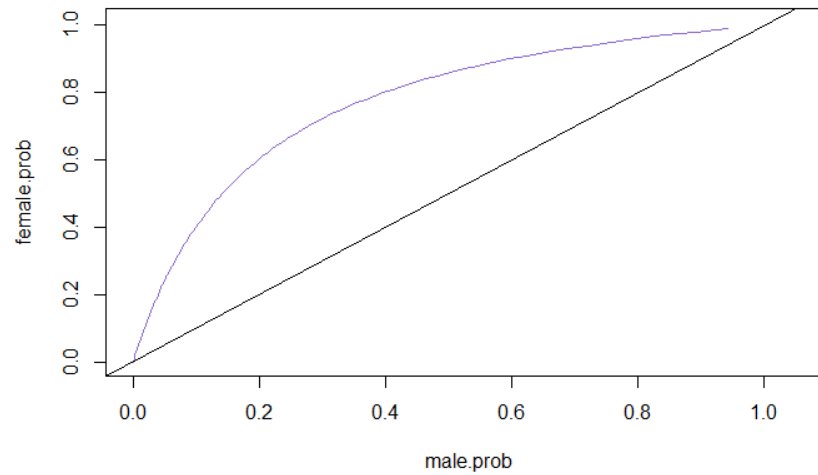


Figure 63: Odds ratio curve for sex - ref : M

This curve indicates the probability of passing G3 ($\text{catG3} = 1$), for male reference point :

$$x : P(\text{catG3}|\text{Male}) \sim y : P(\text{catG3}|\text{Female})$$

c.

ROC stands for Receiver Operating Characteristics, and it is used to evaluate the prediction accuracy of a classifier model. ROC curve is a metric describing the trade-off between the sensitivity (true positive rate, TPR) and specificity (false positive rate, FPR) of a prediction in all probability cutoffs (thresholds).

It can be used for binary and multi-class classification accuracy checking.

To evaluate the ROC in multi-class prediction, we create binary classes by mapping each class against the other classes.

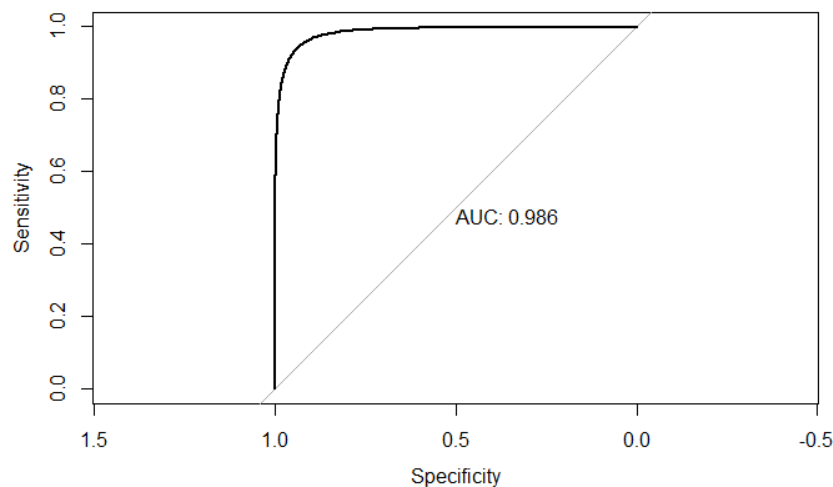


Figure 64: ROC curve - test

The AUC represents the area under the ROC curve. We can evaluate the model the performance by the value of AUC. Higher than 0.5 shows a better model performance. If the curve changes to rectangle it is perfect classifier with AUC value 1.

In our case, AUC is nearly 0.98 which is really good considering all that was mentioned.

d.

The explanatory variable with the lowest p-value in the model, plays the most significant role in the prediction.

e.

According to the summary of our model, *G2* and *failures* are the explanatory variables with the most significant contribution to the model.

```
Call:
glm(formula = catG3 ~ failures + G2, family = binomial(link = "logit"),
    data = StudentsPerformance)

Deviance Residuals:
    Min       1Q   Median       3Q      Max
-2.91186   0.00028   0.01588   0.14760   2.39928

Coefficients:
            Estimate Std. Error z value Pr(>|z|)
(Intercept) -14.9789     2.4607  -6.087 1.15e-09 ***
failures     -0.6456     0.3882  -1.663  0.0963 .
G2           1.6030     0.2453   6.536 6.31e-11 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 433.50  on 394  degrees of freedom
Residual deviance: 104.12  on 392  degrees of freedom
AIC: 110.12

Number of Fisher Scoring iterations: 8
```

Figure 65: GLM model

$$\log\left(\frac{p}{1-p}\right) = -14.978 - 0.645 \times \text{failures} + 1.603 \times G2$$

intercept : keeping all other predictors zero, the log odds ratio / odds ratio of catG3 is $-14.978 / \exp(-14.978) = 4.85e-7$

failures : keeping all other predictors constant for a unit increase in failures, the log odds ratio / odds ratio of catG3 will decrease $-0.645 / \exp(-0.645) = 0.52$

G2 : keeping all other predictors constant for a unit increase in G2 , the log odds ratio / odds ratio of catG3 will increase $1.603 / \exp(1.603) = 4.96$

Produces a table of fit statistics for multiple glm models: AIC, AICc, BIC, p-value, pseudo R-squared (McFadden, Cox and Snell, Nagelkerke).

Smaller values for AIC, AICc, and BIC indicate a better balance of goodness-of-fit of the model and the complexity of the model. The goal is to find a model that adequately explains the data without having too many terms.

BIC tends to choose models with fewer parameters relative to AIC.

Rank <dbl>	Df.res <dbl>	AIC <dbl>	AICc <dbl>	BIC <dbl>	McFadden <dbl>	Cox.and.Snell <dbl>	Nagelkerke <dbl>	p.value <dbl>
5	390	114.1	114.3	138	0.7645	0.5678	0.8523	9.060e-71
3	392	112.1	112.2	128	0.7598	0.5656	0.8489	1.498e-72

Figure 66: GLM model comparison

Model analysis :

```

Confusion Matrix and Statistics

      Reference
Prediction 0  1
0      21  2
1       3 73

      Accuracy : 0.9495
      95% CI : (0.8861, 0.9834)
No Information Rate : 0.7576
P-Value [Acc > NIR] : 3.298e-07

      kappa : 0.8605

McNemar's Test P-value : 1

      Sensitivity : 0.8750
      Specificity : 0.9733
Pos Pred value : 0.9130
Neg Pred value : 0.9605
Prevalence : 0.2424
Detection Rate : 0.2121
Detection Prevalence : 0.2323
Balanced Accuracy : 0.9242

      'Positive' class : 0

```

Figure 67: Confusion matrix and accuracy

f.

catG3 is a binary numerical variable indicating whether you pass the test or not.

A perfect regression model needs to have a low false-positive rate and a low false-negative rate.

In minimizing these factors, we face a dilemma, and we have to decide in which case it is more harmful for us to make mistakes.

It will be costly to have a large false-positive. False-positive might ruin your study plans; failing a course might have some harmful effects on your future. But having false-negative, although still bad, is not as costly as false-positive. False-negative will make you study more, although it might cause depression. :))

Outcome	Utility
True Positive	1
True Negative	1
False positive	-80
False Negative	-10

$$U(p) = TP(p) + TN(p) - 80 \times FP(p) - 10 \times FN(p)$$

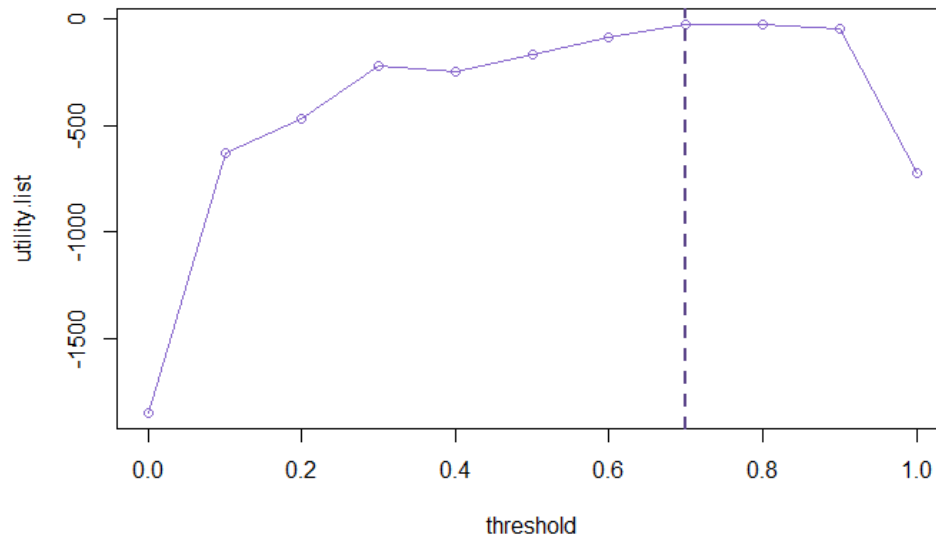


Figure 68: Utility curve

Best threshold : 0.7

Question 7

After converting the sums of Gs to a numeric binary variable :

```
call:
glm(formula = Gsum ~ school + age + Fjob + Mjob + internet +
    romantic + health + failures + goout + studytime + absences +
    sex, family = binomial, data = train)

Deviance Residuals:
    Min       1Q   Median       3Q      Max
-1.7955  -0.5313  -0.3384  -0.1397   2.9123

Coefficients:
            Estimate Std. Error z value Pr(>|z|)
(Intercept)  -6.44436    3.08877  -2.086  0.03694 *
schoolMS      -0.15419    0.55473  -0.278  0.78105
age           0.28826    0.17245   1.672  0.09461 .
Fjobhealth    -0.33615    1.49316  -0.225  0.82188
Fjobother     0.54473    0.79098   0.689  0.49103
Fjobservices  -0.45337    0.82518  -0.549  0.58272
Fjobteacher   0.79923    1.00316   0.797  0.42561
Mjobhealth    -1.90779    1.02396  -1.863  0.06244 .
Mjobother     -0.58612    0.53176  -1.102  0.27036
Mjobservices  -0.45011    0.56543  -0.796  0.42600
Mjobteacher   -0.54123    0.75417  -0.718  0.47297
internetyes   0.25312    0.51771   0.489  0.62490
romanticyes   0.38277    0.39121   0.978  0.32788
health        0.05246    0.13736   0.382  0.70252
failures      1.80570    0.31693   5.697 1.22e-08 ***
goout         0.45450    0.17554   2.589  0.00962 **
studytime     -0.74356    0.28326  -2.625  0.00867 **
absences      -0.10105    0.03580  -2.822  0.00477 **
sexM          -0.78753    0.42097  -1.871  0.06138 .
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 301.62  on 296  degrees of freedom
Residual deviance: 201.84  on 278  degrees of freedom
AIC: 239.84

Number of Fisher Scoring iterations: 6
```

Figure 69: GLM model of all variables

Significant predictors are the ones with the p – value smaller than 0.05 :

	p_value
(Intercept)	1
schoolMS	0
age	0
Fjobhealth	0
Fjobother	0
Fjobservices	0
Fjobteacher	0
Mjobhealth	0
Mjobother	0
Mjobservices	0
Mjobteacher	0
internetyes	0
romanticyes	0
health	0
failures	1
goout	1
studytime	1
absences	1
sexM	0

Figure 70: GLM model of all variables

According to figure 63, the variables that have significant p-value will be selected

$$Gsum \sim failures + goout + studytime + absence$$

Accuracy is 0.86, which is good enough for this model.

86% of the time, we can correctly predict whether a student will be on academic probation or not.

There are several statistics that can help us determine which predictor variables are most important in regression models. These statistics might not agree because the manner in which each one defines "most important" is a bit different :

- P-value : Look for the predictor variable with the lowest p-value
- Standardized regression coefficients : Look for the predictor variable with the largest absolute value for the standardized coefficient.
- Change in R-squared when the variable is added to the model last : Look for the predictor variable that is associated with the greatest increase in R-squared. (explained comprehensively in next question)

The variable with the most effect on academic probation is the variable with the least p-value, which is failures which makes sense.

Confusion Matrix and Statistics

```

      Reference
Prediction 0  1
0      10  1
1      12  75

      Accuracy : 0.8673
      95% CI : (0.7838, 0.9274)
      No Information Rate : 0.7755
      P-Value [Acc > NIR] : 0.015801

      Kappa : 0.5367

      McNemar's Test P-value : 0.005546

      Sensitivity : 0.4545
      Specificity : 0.9868
      Pos Pred Value : 0.9091
      Neg Pred Value : 0.8621
      Prevalence : 0.2245
      Detection Rate : 0.1020
      Detection Prevalence : 0.1122
      Balanced Accuracy : 0.7207

      'Positive' Class : 0
```

Figure 71: Prediction accuracy

R Codes

```

1  ———
2  title: "Statistical Inference"
3  output:
4    pdf_document: default
5    #html_notebook: default
6  ———
7  <h1> Phase 2 </h1>
8  <h2> Dataset : Students Performance </h2>
9
10
11  ### Question 0
12  #### Refreshing the memory:
13  ```{r}
14  set.seed(NULL)
15  StudentsPerformance <- read.csv("StudentsPerformance.csv")
16  head(StudentsPerformance)
17  summary(StudentsPerformance)
18
19  ```
20  ### Question 1
21  #### Chosen variables : *sex* and *Mjob*
22
23  ##### a.
24
25
26
27  First we have to compute the proportions :
28  ```{r warning=FALSE}
29  sample.size <- 200
30  sp.sample <- StudentsPerformance[sample(nrow(StudentsPerformance), sample.size), ]
31  sp.sampled.table <- table(sp.sample[,c("sex", "Fjob")])
32  sp.sampled.table
33
34  F.phat <- sp.sampled.table["F", "teacher"]/sum(sp.sampled.table["F",])
35  F.phat
36
37  M.phat <- sp.sampled.table["M", "teacher"]/sum(sp.sampled.table["M",])
38  M.phat
39  ```
40
41  ```{r}
42  SE <- sqrt( M.phat*(1-M.phat)/sum(sp.sampled.table["F",]) + M.phat*(1-M.phat)/sum(sp.sampled
43    .table["M",]) )
44  SE
45  ```
46  ```{r}
47  (M.phat - F.phat) + c(-1, 1)*pnorm(0.975, lower.tail = F)*SE
48  ```
49
50
51  ##### b.
52  ```{r}
53
54  p.pool <- (sp.sampled.table["F", "teacher"] + sp.sampled.table["M", "teacher"]) / (sum(sp.
55    sampled.table["F",]) + sum(sp.sampled.table["M",]))
56  p.pool
57  SE.pool <- sqrt( p.pool * (1-p.pool)*( 1/ (sum(sp.sampled.table["F",])) + 1/ (sum(sp.sampled
58    .table["M",]))))

```

```

58 SE.pool
59
60 p_value <- pnorm((M.phat - F.phat) / SE.pool, lower.tail = FALSE)
61
62 hypothesis.test <-function(pvalue, alpha = 0.05){
63   if (pvalue < alpha){cat("Due to the fact that p-value (", round(pvalue, 3) , ") is smaller
64     than", alpha, ", we reject the null hypothesis.")}
65   else {cat("Due to the fact that p-value (", round(pvalue, 3) , ") is larger than " ,alpha
66     , ",we fail to reject the null hypothesis.")}
67   }
68
69 hypothesis.test(p_value)
70
71 " "
72
73 " "{r}
74
75 sp.sampled.table
76 sp.sampled.table.bind <- cbind(sp.sampled.table[,1] + sp.sampled.table[, 2], sp.sampled.
77   table[, 3:5] )
78 sp.sampled.table.bind
79
80 chisq.test(sp.sampled.table , rescale.p = T)
81 chisq.test(sp.sampled.table.bind , rescale.p = T)
82
83 " "
84
85 ### Question 2
86 ##### Chosen varibales : *romantic*
87
88 " "{r}
89 sample.size <- 15
90 romantic.sample <- StudentsPerformance[sample(nrow(StudentsPerformance), sample.size), ]$
91   romantic
92 p.hat <- length(which(romantic.sample == 'yes' ))/sample.size
93 p.hat
94
95 simulation <- data.frame(t(replicate(n = 1000, sample(levels(as.factor(StudentsPerformance$
96   romantic))), size = sample.size, replace = TRUE)))
97
98 simulation.success <- apply(simulation, 1, function(x) length(which(x == 'yes'))
99 p_value <- length(which(simulation.success >= 8))/1000
100 hypothesis.test(p_value)
101
102 hist(simulation.success/sample.size)
103
104 " "
105
106
107
108 ### Question 3
109 ##### Chosen varibales : *Mjob*
110
111
112 ##### a.
113 " "{r}
114
115 sample.original <- StudentsPerformance$Mjob
116 round(table(sample.original) / length(StudentsPerformance$Mjob), 4)
117
118 sample.size <- 100
119
120

```

```

115 sample.unbiased <- sample(StudentsPerformance$Mjob, sample.size, replace = FALSE)
116 unbiased.table <- table(sample.unbiased)
117 unbiased.table
118
119
120 biased.prob <- ifelse(StudentsPerformance$Mjob == "teacher", 0.6, 0.4)
121 sample.biased <- sample(StudentsPerformance$Mjob, sample.size, prob = biased.prob)
122 biased.table <- table(sample.biased)
123 biased.table
124
125
126 original_prob <- c(prop.table(table(StudentsPerformance$Mjob)))
127 chisq.test(unbiased.table, p = original_prob)
128 chisq.test(biased.table, p = original_prob)
129 '''
130 ##### Chosen variables : *Fjob*
131
132
133 ##### b.
134 '''{r}
135 Mjob.Fjob <- table(sp.sample[,c("Mjob","Fjob")])
136 Mjob.Fjob
137 chisq.test(Mjob.Fjob)
138
139 Mjob.Fjob.bind <- cbind(Mjob.Fjob[,1] + Mjob.Fjob[, 2] + Mjob.Fjob[, 4] + Mjob.Fjob[, 5],
140                        Mjob.Fjob[, 3] )
141 Mjob.Fjob.bind
142
143 chisq.test(Mjob.Fjob.bind, rescale.p = T)
144 '''
145 ##### Question 4
146 ##### Chosen variables : *G1* , *failure* and *studytime*
147
148
149 ##### a.
150
151 '''{r}
152 library(ggplot2)
153 library("ggpubr")
154 library(GGally)
155 cor(StudentsPerformance$failures, StudentsPerformance$G1)
156 cor(StudentsPerformance$studytime, StudentsPerformance$G1)
157 cor(StudentsPerformance$goout, StudentsPerformance$G1)
158
159 ggpairs(StudentsPerformance[, c(11, 10, 14)])
160 '''
161
162 ##### b.
163 ##### a. and b.
164 '''{r}
165 #just failure
166 lm.G1.failure <- lm(G1 ~ failures, data = StudentsPerformance)
167 summary(lm.G1.failure)
168 lm.G1.failure
169 '''
170
171 '''{r}
172 #condition
173 library(ggplot2)
174 library(ggfortify)
175 autoplot(lm.G1.failure)+ theme_classic()

```

```

176
177 ' '
178
179
180 '{r}'
181 #just studytime
182 lm.G1.studytime <- lm(G1 ~ studytime, data = StudentsPerformance)
183 summary(lm.G1.studytime)
184 lm.G1.studytime
185 ' '
186
187 '{r}'
188 #condition
189 library(ggplot2)
190 library(ggfortify)
191 autoplot(lm.G1.studytime)+ theme_classic()
192 ' '
193
194 ##### c.
195 '{r}'
196 G1.failures <- ggplot(StudentsPerformance, aes(x = failures)) + geom_point(aes(y = G1), size
    = 2, colour = "grey") + stat_smooth(aes(x = failures, y = G1, linetype = "Linear Fit"),
    method = "lm", formula = y ~ x, se = F, color = "black")+ scale_linetype_manual(name =
    "Fit Type", values = c(2, 2)) + ggtitle("G1 vs. failures")
197
198 G1.failures + theme_classic()
199
200 G1.studytime <- ggplot(StudentsPerformance, aes(x = studytime)) + geom_point(aes(y = G1),
    size = 2, colour = "grey") + stat_smooth(aes(x = studytime, y = G1, linetype = "Linear
    Fit"), method = "lm", formula = y ~ x, se = F, color = "black")+ scale_linetype_manual(
    name = "Fit Type", values = c(2, 2)) + ggtitle("G1 vs. studytime")
201
202 G1.studytime + theme_classic()
203
204 ' '
205
206
207
208 ##### e.
209 '{r}'
210
211 compute.R.sqr <- function(model){
212   SS.reg <- (anova(model)[[2]])[1] + (anova(model)[[2]])[2]
213   SS.res <- (anova(model)[[2]])[3]
214   R.sqr.f <- SS.reg / (SS.res + SS.reg)
215   return(R.sqr.f)
216 }
217
218
219 base.model <- lm(G1 ~ sex, data = StudentsPerformance)
220 anova(base.model)
221 SS.reg <- (anova(base.model)[[2]])[1]
222 SS.res <- (anova(base.model)[[2]])[2]
223 R.sqr <- SS.reg / (SS.res + SS.reg)
224
225
226 #failure vs. studytime :
227 model.s.f <- lm(G1 ~ sex + failures, data = StudentsPerformance)
228 anova(model.s.f)
229 R.sqr.f <- compute.R.sqr(model.s.f)
230
231 model.s.s <- lm(G1 ~ sex + studytime, data = StudentsPerformance)

```

```

232 anova(model.s.s)
233 R.sqr.s <- compute.R.sqr(model.s.s)
234
235 R.square <- c(R.sqr, R.sqr.f, R.sqr.s)
236 df <- data.frame(R2 = round(R.square, 2))
237 df <- t(df)
238 colnames(df) <- c("Base", " + failures", " + studytime")
239 df
240
241 #G2 vs. G3
242 model.s.2 <- lm(G1 ~ sex + G2, data = StudentsPerformance)
243 anova(model.s.2)
244 R.sqr.2 <- compute.R.sqr(model.s.2)
245
246 model.s.3 <- lm(G1 ~ sex + G3, data = StudentsPerformance)
247 anova(model.s.3)
248 R.sqr.3 <- compute.R.sqr(model.s.3)
249
250 R.square. <- c(R.sqr, R.sqr.2, R.sqr.3)
251 df <- data.frame(R2 = round(R.square., 2))
252 df <- t(df)
253 colnames(df) <- c("Base", " + G2", " + G3")
254 df
255
256 ‘‘‘
257
258 ##### e.
259 ##### a.
260 ‘‘{r}
261 require(caTools)
262 set.seed(101)
263
264 sample.size <- 100
265 sp.sample <- StudentsPerformance[sample(nrow(StudentsPerformance), sample.size), ]
266
267 sample <- sample.split(sp.sample$G1, SplitRatio = 9/10)
268 G1.train <- subset(sp.sample, sample == TRUE)
269 G1.test <- subset(sp.sample, sample == FALSE)
270 ‘‘‘
271
272
273 ‘‘{r}
274 #failures
275 lm.G1.failures <- lm(G1 ~ failures, data = G1.train)
276 summary(lm.G1.failures)
277 p_value <- summary(lm.G1.failures)$coefficients[8]
278 hypothesis.test(p_value)
279 ‘‘‘
280
281
282 ‘‘{r}
283 #studytime
284 lm.G1.studytime <- lm(G1 ~ studytime, data = G1.train)
285 summary(lm.G1.studytime)
286 p_value <- summary(lm.G1.studytime)$coefficients[8]
287 hypothesis.test(p_value)
288 ‘‘‘
289
290 ‘‘{r}
291 #G2
292 lm.G1.G2 <- lm(G1 ~ G2, data = G1.train)
293 summary(lm.G1.G2)

```

```

294 p_value <- summary(lm.G1.G2)$coefficients[8]
295 hypothesis.test(p_value)
296 ' '
297
298
299 '{r}'
300 #G3
301 lm.G1.G3 <- lm(G1 ~ G3, data = G1.train)
302 summary(lm.G1.G3)
303 p_value <- summary(lm.G1.G3)$coefficients[8]
304 hypothesis.test(p_value)
305 ' '
306
307 ##### b.
308 '{r}'
309
310 calculate.CI <- function(model, alpha = 0.05){
311
312   point.est <- summary(model)$coefficients[2]
313   std.error <- summary(model)$coefficients[4]
314
315   round(point.est + c(-1, 1) * pnorm(1 - alpha/2) * std.error, 3)
316 }
317
318 calculate.CI(lm.G1.failures)
319
320 calculate.CI(lm.G1.studytime)
321
322 calculate.CI(lm.G1.G2)
323
324 calculate.CI(lm.G1.G3)
325
326 ' '
327
328 ##### c.
329 '{r}'
330 predicted.s <- round(predict(lm.G1.studytime, G1.test, type = "response"),1)
331 predicted.f <- round(predict(lm.G1.failures, G1.test, type = "response"),1)
332 predicted.2 <- round(predict(lm.G1.G2, G1.test, type = "response"),1)
333 predicted.3 <- round(predict(lm.G1.G3, G1.test, type = "response"),1)
334
335
336
337 pred.actual <- data.frame(G1.test$G1, predicted.s, predicted.f, predicted.2, predicted.3)
338 colnames(pred.actual) <- c("Actual", "Predicted studytime", "Predicted failues", "Predicted
    G2", "Predicted G3")
339
340 ' '
341 ##### d.
342 '{r}'
343 # 0.1 * data_range = error
344 error <- abs(G1.test$G1 - pred.actual)
345 error
346
347 succes.rate.list <- c()
348 for (predictor in 1:length(error)) {
349   error.accepted <- length(which(error[predictor] <= 2))
350   succes.rate <- paste((error.accepted / length(G1.test$G1))*100, "%")
351   succes.rate.list <- c(succes.rate.list, succes.rate)
352 }
353
354

```

```

355 succes.rate <- data.frame(t(succes.rate.list[2:5]))
356 colnames(succes.rate) <- c("Predicted studytime", "Predicted failues", "Predicted G2", "
    Predicted G3")
357 succes.rate
358 ""
359
360
361 ""{r}
362 # Min-Max Accuracy Calculation
363 predictors <- data.frame(predicted.s, predicted.f, predicted.2, predicted.3)
364
365 mm.succes.rate.list <- c()
366 for (p in 1:length(predictors)) {
367   actuals.preds <- data.frame(cbind(actuals = G1.test$G1, predicteds = predictors[p]))
368   min.max.succes.rate <- paste(round((mean(apply(actuals.preds, 1, min) / apply(actuals.
    preds, 1, max)))*100, 2), "%")
369   mm.succes.rate.list <- c(mm.succes.rate.list, min.max.succes.rate)
370 }
371
372 succes.rate <- data.frame((t(mm.succes.rate.list)))
373 colnames(succes.rate) <- c("Predicted studytime", "Predicted failues", "Predicted G2", "
    Predicted G3")
374 succes.rate
375
376
377 ""
378
379
380 ""{r}
381 # MAPE Calculation
382
383
384 mape.succes.rate.list <- c()
385 for (p in 1:length(predictors)) {
386   actuals.preds <- data.frame(cbind(actuals = G1.test$G1, predicteds = predictors[p]))
387   mape.succes.rate <- paste(round((mean(abs((actuals.preds$predicted - actuals.preds$
    actuals))/actuals.preds$actuals) )*100, 2), "%")
388   mape.succes.rate.list <- c(mm.succes.rate.list, min.max.succes.rate)
389 }
390
391
392
393 succes.rate <- data.frame((t(mape.succes.rate.list)))
394 colnames(succes.rate) <- c("Predicted studytime", "Predicted failues", "Predicted G2", "
    Predicted G3")
395 succes.rate
396
397
398
399 ""
400
401
402 ##### extra part
403 ""{r}
404 library(ggplot2)
405 library("ggpubr")
406 library(GGally)
407 cor(StudentsPerformance$G2, StudentsPerformance$G1)
408 cor(StudentsPerformance$G3, StudentsPerformance$G1)
409
410
411 ggpairs(StudentsPerformance[, c(15, 16, 14)])

```



```

412  ``
413  ``{r}
414  #G2
415  lm.G1.G2 <- lm(G1 ~ G2, data = StudentsPerformance)
416  summary(lm.G1.G2)
417  lm.G1.G2
418
419  library(ggplot2)
420  library(ggfortify)
421  autoplot(lm.G1.G2)+ theme_classic()
422
423  G1.G2 <- ggplot(StudentsPerformance, aes(x = G2)) + geom_point(aes(y = G1), size = 2, colour
    = "grey") + stat_smooth(aes(x = G2, y = G1, linetype = "Linear Fit"), method = "lm",
    formula = y ~ x, se = F, color = "black")+ scale_linetype_manual(name = "Fit Type",
    values = c(2, 2)) + ggtitle("G1 vs. G2")
424
425  G1.G2 + theme_classic()
426
427  ``
428  ``{r}
429  #G2
430  lm.G1.G3 <- lm(G1 ~ G3, data = StudentsPerformance)
431  summary(lm.G1.G3)
432  lm.G1.G3
433
434  library(ggplot2)
435  library(ggfortify)
436  autoplot(lm.G1.G3)+ theme_classic()
437
438  G1.G3 <- ggplot(StudentsPerformance, aes(x = G3)) + geom_point(aes(y = G1), size = 2, colour
    = "grey") + stat_smooth(aes(x = G3, y = G1, linetype = "Linear Fit"), method = "lm",
    formula = y ~ x, se = F, color = "black")+ scale_linetype_manual(name = "Fit Type",
    values = c(2, 2)) + ggtitle("G1 vs. G3")
439
440  G1.G3 + theme_classic()
441  ``
442
443  ### Question 5
444  ##### Chosen response variable : *G1*
445  #####Chosen explanatory variables : *G2*, *goout*, *failures*, *studytime*, *sex* , *age*
446
447
448  ##### a.
449
450  ``{r message=FALSE, warning=FALSE}
451
452  library(GGally)
453  p_ <- GGally::print_if_interactive
454  pm <- ggpairs(StudentsPerformance[, c(3, 4, 7, 10, 11, 14, 15)], progress = FALSE) + theme_
    minimal()
455  p_(pm)
456
457  pm <- ggpairs(StudentsPerformance[, c(3, 4, 7, 10, 11, 15)], progress = FALSE) + theme_
    minimal()
458  p_(pm)
459
460
461  StudentsPerformance$sex <- ifelse(StudentsPerformance$sex == "F", 1, 0)
462  library(ggcorrplot)
463  ggcorrplot(cor(StudentsPerformance[, c(3, 4, 7, 10, 11, 14, 15)]) , type = "lower", lab =
    TRUE, outline.color = "white", colors = c("black", "white", "mediumpurple3"))
464

```

```

465
466 ' '
467
468
469 ##### b.
470 '{r}
471 lm.model <- lm(G1 ~ G2 + goout + failures + studytime + sex + age , data =
  StudentsPerformance)
472 summary(lm.model)
473
474 ' '
475
476 '{r}
477 plot(lm.model$residuals , pch = 16, col = "mediumpurple3") + abline(lm.model)
478 ' '
479
480
481 ##### e.
482 '{r}
483 library(olsrr)
484 #forward - p-value
485 forward.selection.p <- ols_step_forward_p(lm.model, details = TRUE, prem = 0.05)
486
487 #backward - p-value
488 backward.elimination.p <- ols_step_backward_p(lm.model, details = TRUE, prem = 0.05)
489
490 ' '
491
492 '{r}
493 #forward - adjusted R-sqrt
494 library(rms)
495
496 best.pred <- c()
497
498 adj.r.squared <- function(formula , dataset , k = 1) {
499   n <- length(StudentsPerformance$G1)
500   r.squared <- lrm(formula = formula , data = dataset)$stat["R2"]
501   adjR2 <- 1 - (((n-1)/(n-k-1)) * (1-r.squared))
502 }
503
504
505 #step 1
506 adj.r.squared.list1 <- c()
507 names <- c("G2" , "goout" , "failures" , "studytime" , "sex" , "age")
508 adj.r.squared.list1 <- c(adj.r.squared(G1 ~ G2, StudentsPerformance),
509   adj.r.squared(G1 ~ goout, StudentsPerformance),
510   adj.r.squared(G1 ~ failures, StudentsPerformance),
511   adj.r.squared(G1 ~ studytime, StudentsPerformance),
512   adj.r.squared(G1 ~ sex, StudentsPerformance),
513   adj.r.squared(G1 ~ age, StudentsPerformance))
514
515
516
517 max.adj.r.squared <- names[which.max(adj.r.squared.list1)]
518 if (max(adj.r.squared.list1 , 0) > 0) { best.pred <- c(best.pred , max.adj.r.squared) }
519 best.pred
520
521
522 #step 2
523 names <- c("G2 + goout" , "G2 + failures" , "G2 + studytime" , "G2 + sex" , "G2 + age")
524 adj.r.squared.list2 <- c(adj.r.squared(G1 ~ goout + G2, StudentsPerformance, k = 2),
525   adj.r.squared(G1 ~ failures + G2, StudentsPerformance, k = 2),

```

```

526         adj.r.square(G1 ~ studytime + G2, StudentsPerformance, k = 2),
527         adj.r.square(G1 ~ sex + G2, StudentsPerformance, k = 2),
528         adj.r.square(G1 ~ age + G2, StudentsPerformance, k = 2))
529
530
531 max.adj.r.squared <- names[which.max(adj.r.squared.list2 - max(adj.r.squared.list1))]
532 if (max(adj.r.squared.list2 - max(adj.r.squared.list1)) > 0) { best.pred <- c(best.pred, max
533   .adj.r.squared) }
534 best.pred
535
536 #step 3
537 names <- c("G2 + age + goout" , "G2 + age + failures" , "G2 + age + studytime" , "G2 + age
538   + sex" )
539 adj.r.squared.list3 <- c(adj.r.square(G1 ~ goout + G2 + age , StudentsPerformance, k = 3),
540   adj.r.square(G1 ~ failures + G2 + age , StudentsPerformance, k = 3),
541   adj.r.square(G1 ~ studytime + G2 + age , StudentsPerformance, k = 3)
542   ,
543   adj.r.square(G1 ~ sex + G2 + age , StudentsPerformance, k = 3))
544 max.adj.r.squared <- names[which.max(adj.r.squared.list3 - max(adj.r.squared.list2))]
545 if (max(adj.r.squared.list3 - max(adj.r.squared.list2)) > 0) { best.pred <- c(best.pred, max
546   .adj.r.squared) }
547 best.pred
548
549 #step 4
550 names <- c("G2 + age + failures + goout" , "G2 + age + failures + studytime" , "G2 + age
551   + failures + sex" )
552 adj.r.squared.list4 <- c(adj.r.square(G1 ~ goout + G2 + age + failures ,
553   StudentsPerformance, k = 4),
554   adj.r.square(G1 ~ studytime + G2 + age + failures ,
555   StudentsPerformance, k = 4),
556   adj.r.square(G1 ~ sex + G2 + age + failures , StudentsPerformance,
557   k = 4))
558 adj.r.squared.list4
559 max.adj.r.squared <- names[which.max(adj.r.squared.list4 - max(adj.r.squared.list3))]
560 adj.r.squared.list4 - max(adj.r.squared.list3)
561
562 if (max(adj.r.squared.list3 - max(adj.r.squared.list2)) > 0) { best.pred <- c(best.pred, max
563   .adj.r.squared) }
564 best.pred
565
566 all.adj.r.squared <- c(max(adj.r.squared.list1), max(adj.r.squared.list2), max(adj.r.squared
567   .list3), max(adj.r.squared.list4))
568
569 model <- data.frame(best.pred, all.adj.r.squared)
570 model
571
572
573 ""
574 ""{r}
575 #backward - adjusted R-sqrt
576 library(rms)
577
578 fullmodel.adj.r.sqr <- adj.r.square(G1 ~ G2 + goout + failures + studytime + sex + age ,
579   StudentsPerformance ,k = 6)
580 best.pred <- c()
581
582

```

```

577 #step 1
578 adj.r.squared.list1 <- c()
579 names <- c("G2 + goout + failures + studytime + sex" , "G2 + goout + failures + studytime
+ age" ,
580           "G2 + goout + failures + sex + age" , "G2 + goout + studytime + sex + age" ,
581           "G2 + failures + studytime + sex + age" , "goout + failures + studytime + sex +
age")
582 adj.r.squared.list1 <- c(adj.r.square(G1 ~ G2 + goout + failures + studytime + sex,
StudentsPerformance, k = 5),
583                          adj.r.square(G1 ~ G2 + goout + failures + studytime + age,
StudentsPerformance, k = 5),
584                          adj.r.square(G1 ~ G2 + goout + failures + sex + age,
StudentsPerformance, k = 5),
585                          adj.r.square(G1 ~ G2 + goout + studytime + sex + age,
StudentsPerformance, k = 5),
586                          adj.r.square(G1 ~ G2 + failures + studytime + sex + age,
StudentsPerformance, k = 5),
587                          adj.r.square(G1 ~ goout + failures + studytime + sex + age,
StudentsPerformance, k = 5))
588
589
590
591 max.adj.r.squared <- names[which.max(adj.r.squared.list1 - fullmodel.adj.r.sqr)]
592 if ( (max(adj.r.squared.list1) - fullmodel.adj.r.sqr) > 0) { best.pred <- c(best.pred, max.
adj.r.squared) }
593 best.pred
594
595
596 #step 2
597 adj.r.squared.list2 <- c()
598 names <- c("G2 + failures + studytime + sex" , "G2 + failures + studytime + age" ,
599           "G2 + failures + sex + age" , "G2 + studytime + sex + age",
600           "failures + studytime + sex + age")
601 adj.r.squared.list2 <- c(adj.r.square(G1 ~ G2 + failures + studytime + sex,
StudentsPerformance, k = 4),
602                          adj.r.square(G1 ~ G2 + failures + studytime + age,
StudentsPerformance, k = 4),
603                          adj.r.square(G1 ~ G2 + failures + sex + age, StudentsPerformance, k
= 4),
604                          adj.r.square(G1 ~ G2 + studytime + sex + age, StudentsPerformance, k
= 4),
605                          adj.r.square(G1 ~ failures + studytime + sex + age,
StudentsPerformance, k = 4))
606
607
608
609 max.adj.r.squared <- names[which.max(adj.r.squared.list2 - max(adj.r.squared.list1))]
610 if ((max(adj.r.squared.list2) - max(adj.r.squared.list1)) > 0) { best.pred <- c(best.pred,
max.adj.r.squared) }
611 best.pred
612
613
614 #step 3
615 adj.r.squared.list3 <- c()
616 names <- c("G2 + failures + studytime" , "G2 + failures + age" ,
617           "G2 + studytime + age", "failures + studytime + age")
618 adj.r.squared.list3 <- c(adj.r.square(G1 ~ G2 + failures + studytime, StudentsPerformance, k
= 3),
619                          adj.r.square(G1 ~ G2 + failures + age, StudentsPerformance, k = 3),
620                          adj.r.square(G1 ~ G2 + studytime + age, StudentsPerformance, k = 3),
621                          adj.r.square(G1 ~ failures + studytime + age, StudentsPerformance, k
= 3))

```

```

622
623
624
625 max.adj.r.squared <- names[which.max(adj.r.squared.list3 - max(adj.r.squared.list2))]
626 if ((max(adj.r.squared.list3) - max(adj.r.squared.list2)) > 0) { best.pred <- c(best.pred,
627   max.adj.r.squared) }
628 best.pred
629
630 #step 4
631 adj.r.squared.list4 <- c()
632 names <- c("G2 + failures" , "G2 + age", "failures + age")
633 adj.r.squared.list4 <- c(adj.r.square(G1 ~ G2 + failures , StudentsPerformance , k = 2) ,
634   adj.r.square(G1 ~ G2 + age , StudentsPerformance , k = 2) ,
635   adj.r.square(G1 ~ failures + age , StudentsPerformance , k = 2))
636
637
638
639 max.adj.r.squared <- names[which.max(adj.r.squared.list4 - max(adj.r.squared.list3))]
640 if ((max(adj.r.squared.list4) - max(adj.r.squared.list3)) > 0) { best.pred <- c(best.pred,
641   max.adj.r.squared) }
642 best.pred
643
644 all.adj.r.squared <- c(max(adj.r.squared.list1) , max(adj.r.squared.list2) , max(adj.r.squared
645   .list3))
646
647 model <- data.frame(best.pred , all.adj.r.squared)
648 model
649 ""
650
651
652 "{r}"
653 final.model <- lm(G1 ~ G2 + failures + age , data = StudentsPerformance)
654 summary(final.model)
655 ""
656
657 ##### f.
658 "{r}"
659 #linearity
660 data <- data.frame(G2 = StudentsPerformance$G2, residuals = final.model$residuals)
661 ggplot(data = data, aes(G2, residuals)) + geom_point(color = "mediumpurple3", alpha = 0.5) +
662   stat_smooth(method = lm, se = F, color = "black") + theme_classic()
663
664 data <- data.frame(failures = StudentsPerformance$failures , residuals = final.model$
665   residuals)
666 ggplot(data = data, aes(failures , residuals)) + geom_point(color = "mediumpurple3", alpha =
667   0.5) + stat_smooth(method = lm, se = F, color = "black") + theme_classic()
668
669 data <- data.frame(age = StudentsPerformance$age, residuals = final.model$residuals)
670 ggplot(data = data, aes(age, residuals)) + geom_point(color = "mediumpurple3", alpha = 0.5) +
671   geom_hline( yintercept = 0, size = 1) + theme_classic()
672
673 #nearly normal
674 ggplot(final.model, aes(sample = final.model$residuals)) + stat_qq(col = "mediumpurple3",
675   alpha = 0.5) + stat_qq_line() + theme_classic()
676
677 ggplot(data = final.model, aes(final.model$residuals)) + geom_histogram(bins = 20, col = "
678   mediumpurple2", fill="mediumpurple3", alpha = 0.5) + theme_classic()
679
680 #cons. var

```

```

675 ks.test(unique(final.model$residuals), "pnorm", mean=0, sd=1)
676 ggplot(data = final.model, aes(final.model$fitted, final.model$residuals)) + geom_point(color
    = "mediumpurple3", alpha = 0.5) + stat_smooth(method = lm, se = F, color = "black") +
    theme_classic()
677
678 ""
679
680
681
682 "{r}"
683 library(ggplot2)
684 library(ggfortify)
685 autoplot(final.model) + theme_classic()
686 ""
687
688 ##### g.
689 "{r}"
690 library(caret)
691 model <- trainControl(method = "cv", number = 5)
692 fullmodel.cv <- train(G1 ~ G2 + goout + failures + studytime + sex + age, data =
    StudentsPerformance, trControl = model, method = "lm")
693
694 bestmodel.cv <- train(G1 ~ G2 + failures + age, data = StudentsPerformance, trControl =
    model, method = "lm")
695
696
697 fullmodel.cv
698 bestmodel.cv
699
700
701 fullmodel.cv$finalModel
702 bestmodel.cv$finalModel
703
704 allfolds <- bestmodel.cv$resample
705 ""
706
707 ### Question 6
708
709 "{r}"
710 StudentsPerformance$catG3 <- ifelse(StudentsPerformance$G3 < 10, 0, 1)
711
712 sample <- sample.split(StudentsPerformance$catG3, SplitRatio = 3/4)
713 train <- subset(StudentsPerformance, sample == TRUE)
714 test <- subset(StudentsPerformance, sample == FALSE)
715
716
717 ""
718
719
720 ##### Chosen response variable : *catG3*
721 ##### Chosen explanatory variables : *failures*, *studytime*, *G2* and *sex*
722
723
724 ##### a.
725
726 "{r}"
727 model.glm <- glm(catG3 ~ failures + studytime + G2 + sex, family = binomial(link='logit'),
    data = train)
728
729 summary(model.glm)
730 ""
731

```

```

732 ##### b.
733 ```{r}
734
735 female.prob <- seq(0, 1.01, 0.01)
736 OR.ratio = abs(summary(model.gml)$coefficients[3])
737
738 pred.y <- function(x) {
739   return ((OR.ratio*x/(1-x)) / (1 + (OR.ratio*x/(1-x))))
740 }
741 male.prob <- sapply(female.prob, pred.y)
742 plot(male.prob, female.prob, type = "l", col = "mediumpurple3", lwd = 1.3) + abline(a=0, b
743   =1)
744
745 ```
746
747 ##### c.
748 ```{r message=FALSE, warning=FALSE}
749 library(pROC)
750 require(ROCR)
751
752 pred <- predict(model.gml, train , type="response")
753 roc(catG3 ~ pred, data = train, plot = TRUE, print.auc = TRUE, smooth = TRUE)
754
755
756 pred.t <- predict(model.gml, test , type="response")
757 roc(catG3 ~ pred.t, data = test, plot = TRUE, print.auc = TRUE, smooth = TRUE)
758
759
760 ```
761
762 ##### e.
763 ```{r}
764 library(rcompanion)
765 better.model.gml <- glm(catG3 ~ failures + G2 , family = binomial(link='logit'), data =
766   StudentsPerformance)
767 summary(better.model.gml)
768
769 compareGLM(model.gml, better.model.gml)
770
771 ```
772 ##### f.
773
774 ```{r}
775 library(caret)
776
777
778 confusion.matrix <- function(threshold){
779   prediction.probability <- predict(better.model.gml, newdata = test, type = "response")
780   pos.neg <- ifelse(prediction.probability > threshold, "1", "0")
781   p.class <- factor(pos.neg, levels = c("0", "1"))
782   cm <- confusionMatrix(p.class, as.factor(test$catG3))
783   return(cm)}
784
785
786 confusion.matrix(0.5)
787
788
789 threshold <- seq(0, 1, by = 0.1)
790 utility.list <- c()
791 for (i in 1:length(threshold)){

```

```

792
793   cm <- confusion.matrix(threshold[i])
794
795   TP <- cm$table[1]
796   FP <- cm$table[2]
797   FN <- cm$table[3]
798   TN <- cm$table[4]
799
800   utility <- TP + TN - 80*FP - 10*FN
801   utility.list <- c(utility.list, utility)
802
803 }
804
805 plot(threshold, utility.list, type = "o", col = "mediumpurple3", lwd = 1.3) + abline(v =
      threshold[which.max(utility.list)], col="mediumpurple4", lwd = 2, lty=2)
806
807
808
809
810 ““
811
812 ### Question 7
813 ““{r}
814
815 G.sum <- StudentsPerformance$G1 + StudentsPerformance$G2 + StudentsPerformance$G3
816 StudentsPerformance$Gsum <- ifelse(G.sum < 25, 1, 0)
817
818
819 sample <- sample.split(StudentsPerformance$Gsum, SplitRatio = 3/4)
820 train <- subset(StudentsPerformance, sample == TRUE)
821 test  <- subset(StudentsPerformance, sample == FALSE)
822
823
824
825 model.gml <- glm(Gsum ~ school + age + Fjob + Mjob + internet + romantic + health +failures
      +goout + studytime + absences + sex , family = binomial, data = train)
826
827 summary(model.gml)
828
829 ““
830
831
832 ““{r}
833
834 p.values <- coef(summary(model.gml))[,4]
835
836 p.value <- ifelse(p.values < 0.05, 1, 0)
837 significant.pvalue <- data.frame(p.value)
838
839
840 ““
841
842
843 ““{r}
844
845 prediction.probability <- predict(model.gml, newdata = test, type = "response")
846 pos.neg <- ifelse(prediction.probability > 0.5, "0", "1")
847 p.class <- factor(pos.neg, levels = c("0", "1"))
848 cm <- confusionMatrix(p.class, as.factor(test$catG3))
849
850 cm
851

```


852 | ' ' ' |

code.Rmd

Forward Selection Method

Candidate Terms:

1. G2
2. goout
3. failures
4. studytime
5. sex
6. age

We are selecting variables based on p value...

Forward Selection: Step 1

+ G2

Model Summary			
R	0.851	RMSE	1.852
R-Squared	0.724	Coef. Var	17.174
Adj. R-Squared	0.723	MSE	3.429
Pred R-Squared	0.719	MAE	1.383

RMSE: Root Mean Square Error
MSE: Mean Square Error
MAE: Mean Absolute Error

ANOVA					
	Sum of Squares	DF	Mean Square	F	Sig.
Regression	3537.013	1	3537.013	1031.393	0.0000
Residual	1347.737	393	3.429		
Total	4884.749	394			

Parameter Estimates							
model	Beta	Std. Error	Std. Beta	t	Sig.	lower	upper
(Intercept)	1.785	0.295		6.045	0.000	1.204	2.365
G2	0.733	0.023	0.851	32.115	0.000	0.688	0.778

Forward Selection: Step 2

+ age

Model Summary			
R	0.854	RMSE	1.834
R-Squared	0.730	Coef. Var	17.013
Adj. R-Squared	0.729	MSE	3.365
Pred R-Squared	0.722	MAE	1.348

RMSE: Root Mean Square Error
MSE: Mean Square Error
MAE: Mean Absolute Error

ANOVA					
	Sum of Squares	DF	Mean Square	F	Sig.
Regression	3565.520	2	1782.760	529.735	0.0000
Residual	1319.229	392	3.365		

Total 4884.749 394

Parameter Estimates							
model	Beta	Std. Error	Std. Beta	t	Sig.	lower	upper
(Intercept)	-1.956	1.318		-1.484	0.139	-4.547	0.636
G2	0.746	0.023	0.866	32.383	0.000	0.701	0.791
age	0.215	0.074	0.078	2.910	0.004	0.070	0.360

Forward Selection: Step 3

+ failures

Model Summary			
R	0.856	RMSE	1.825
R-Squared	0.733	Coef. Var	16.928
Adj. R-Squared	0.731	MSE	3.332
Pred R-Squared	0.723	MAE	1.363

RMSE: Root Mean Square Error
MSE: Mean Square Error
MAE: Mean Absolute Error

ANOVA					
	Sum of Squares	DF	Mean Square	F	Sig.
Regression	3582.019	3	1194.006	358.368	0.0000
Residual	1302.730	391	3.332		
Total	4884.749	394			

Parameter Estimates							
model	Beta	Std. Error	Std. Beta	t	Sig.	lower	upper
(Intercept)	-1.994	1.311		-1.521	0.129	-4.573	
G2	0.718	0.026	0.834	27.563	0.000	0.667	
age	0.244	0.075	0.088	3.270	0.001	0.097	
failures	-0.323	0.145	-0.068	-2.225	0.027	-0.608	-

Forward Selection: Step 4

+ studytime

Model Summary			
R	0.857	RMSE	1.823
R-Squared	0.735	Coef. Var	16.906
Adj. R-Squared	0.732	MSE	3.323
Pred R-Squared	0.723	MAE	1.361

RMSE: Root Mean Square Error
MSE: Mean Square Error
MAE: Mean Absolute Error

ANOVA

	Sum of Squares	DF	Mean Square	F	Sig.
Regression	3588.724	4	897.181	269.98	0.0000
Residual	1296.025	390	3.323		
Total	4884.749	394			

Parameter Estimates

model	Beta	Std. Error	Std. Beta	t	Sig.	lower
(Intercept)	-2.205	1.318		-1.673	0.095	-4.796
G2	0.715	0.026	0.830	27.384	0.000	0.664
age	0.239	0.075	0.087	3.204	0.001	0.092
failures	-0.298	0.146	-0.063	-2.043	0.042	-0.585
studytime	0.159	0.112	0.038	1.420	0.156	-0.061

Forward Selection: Step 5

+ sex

Model Summary

R	0.858	RMSE	1.822
R-Squared	0.736	Coef. Var	16.896
Adj. R-Squared	0.732	MSE	3.319
Pred R-Squared	0.722	MAE	1.361

RMSE: Root Mean Square Error
MSE: Mean Square Error
MAE: Mean Absolute Error

ANOVA

	Sum of Squares	DF	Mean Square	F	Sig.
Regression	3593.549	5	718.710	216.526	0.0000
Residual	1291.200	389	3.319		
Total	4884.749	394			

Parameter Estimates

model	Beta	Std. Error	Std. Beta	t	Sig.	lower
(Intercept)	-2.139	1.319		-1.622	0.106	-4.731

0.762	G2	0.711	0.026	0.825	26.946	0.000	0.659	
0.387	age	0.240	0.075	0.087	3.227	0.001	0.094	
0.022	failures	-0.309	0.146	-0.065	-2.119	0.035	-0.597	-
0.434	studytime	0.203	0.117	0.048	1.728	0.085	-0.028	
0.148	sex	-0.235	0.195	-0.033	-1.206	0.229	-0.618	

--

No more variables to be added.

Variables Entered:

+ G2
+ age
+ failures
+ studytime
+ sex

Final Model Output

Model Summary			
R	0.858	RMSE	1.822
R-Squared	0.736	Coef. Var	16.896
Adj. R-Squared	0.732	MSE	3.319
Pred R-Squared	0.722	MAE	1.361

RMSE: Root Mean Square Error
MSE: Mean Square Error
MAE: Mean Absolute Error

ANOVA					
	Sum of Squares	DF	Mean Square	F	Sig.
Regression	3593.549	5	718.710	216.526	0.0000
Residual	1291.200	389	3.319		
Total	4884.749	394			

Parameter Estimates						
model	Beta	Std. Error	Std. Beta	t	Sig.	lower
(Intercept)	-2.139	1.319		-1.622	0.106	-4.731
G2	0.711	0.026	0.825	26.946	0.000	0.659
age	0.240	0.075	0.087	3.227	0.001	0.094
failures	-0.309	0.146	-0.065	-2.119	0.035	-0.597
studytime	0.203	0.117	0.048	1.728	0.085	-0.028
sex	-0.235	0.195	-0.033	-1.206	0.229	-0.618

Backward Elimination Method

Candidate Terms:

- 1 . G2
- 2 . goout
- 3 . failures
- 4 . studytime
- 5 . sex
- 6 . age

We are eliminating variables based on p value...

x goout

Backward Elimination: Step 1

Variable goout Removed

Model Summary			
R	0.858	RMSE	1.822
R-Squared	0.736	Coef. Var	16.896
Adj. R-Squared	0.732	MSE	3.319
Pred R-Squared	0.722	MAE	1.361

RMSE: Root Mean Square Error
MSE: Mean Square Error
MAE: Mean Absolute Error

ANOVA					
	Sum of Squares	DF	Mean Square	F	Sig.
Regression	3593.549	5	718.710	216.526	0.0000
Residual	1291.200	389	3.319		
Total	4884.749	394			

Parameter Estimates							
model	Beta	Std. Error	Std. Beta	t	Sig	lower	
(Intercept)	-2.139	1.319		-1.622	0.106	-4.731	
G2	0.711	0.026	0.825	26.946	0.000	0.659	
failures	-0.309	0.146	-0.065	-2.119	0.035	-0.597	-
studytime	0.203	0.117	0.048	1.728	0.085	-0.028	
sex	-0.235	0.195	-0.033	-1.206	0.229	-0.618	
age	0.240	0.075	0.087	3.227	0.001	0.094	

x sex

Backward Elimination: Step 2

Variable sex Removed

Model Summary			
---------------	--	--	--

R	0.857	RMSE	1.823
R-Squared	0.735	Coef. Var	16.906
Adj. R-Squared	0.732	MSE	3.323
Pred R-Squared	0.723	MAE	1.361

RMSE: Root Mean Square Error
MSE: Mean Square Error
MAE: Mean Absolute Error

ANOVA

	Sum of Squares	DF	Mean Square	F	Sig.
Regression	3588.724	4	897.181	269.98	0.0000
Residual	1296.025	390	3.323		
Total	4884.749	394			

Parameter Estimates

model	Beta	Std. Error	Std. Beta	t	Sig.	lower
(Intercept)	-2.205	1.318		-1.673	0.095	-4.796
G2	0.715	0.026	0.830	27.384	0.000	0.664
failures	-0.298	0.146	-0.063	-2.043	0.042	-0.585
studytime	0.159	0.112	0.038	1.420	0.156	-0.061
age	0.239	0.075	0.087	3.204	0.001	0.092

x studytime

Backward Elimination: Step 3

Variable studytime Removed

Model Summary

R	0.856	RMSE	1.825
R-Squared	0.733	Coef. Var	16.928
Adj. R-Squared	0.731	MSE	3.332
Pred R-Squared	0.723	MAE	1.363

RMSE: Root Mean Square Error
MSE: Mean Square Error
MAE: Mean Absolute Error

ANOVA

	Sum of Squares	DF	Mean Square	F	Sig.
Regression	3582.019	3	1194.006	358.368	0.0000
Residual	1302.730	391	3.332		
Total	4884.749	394			

Parameter Estimates

model	Beta	Std. Error	Std. Beta	t	Sig.	lower
upper						

