

# **Statistical Inference: Project Phase I**

Narjes Noorzad - 810196626

## 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

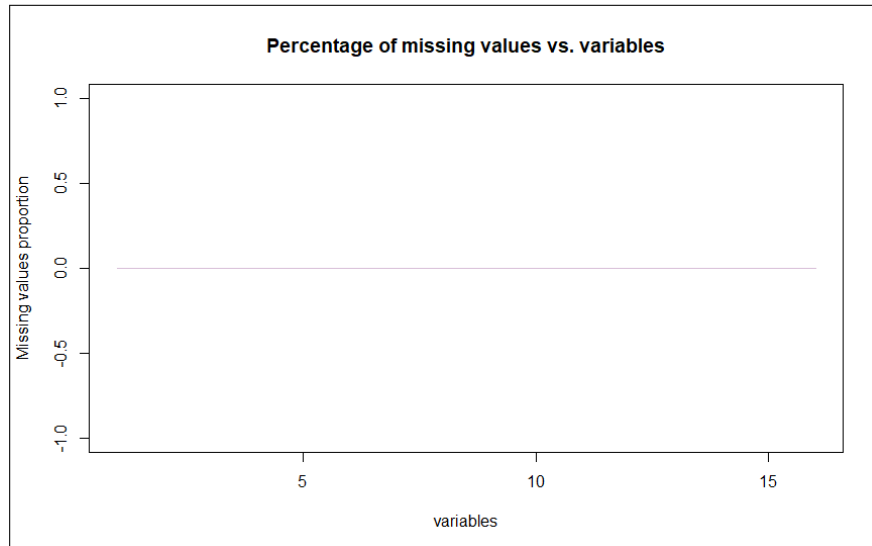


Figure 4: Line plot of missing value proportion

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*.

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

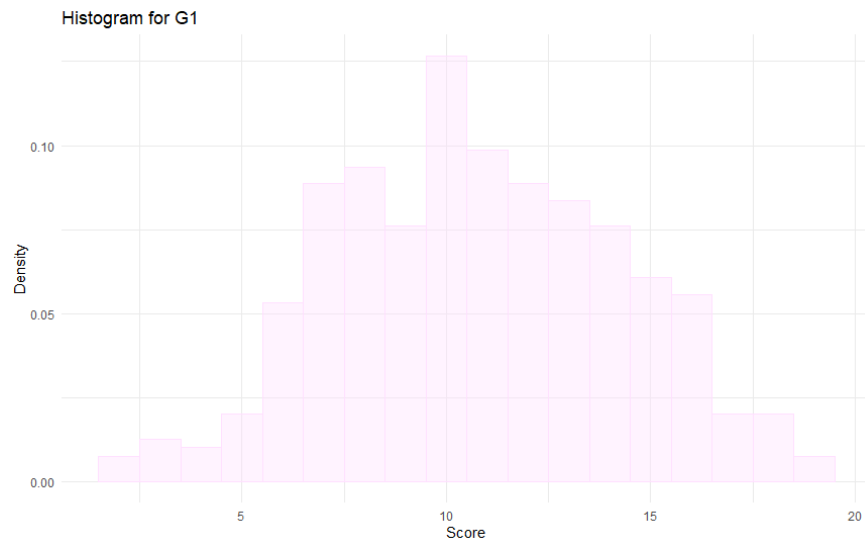


Figure 5: Histogram of *G1*

*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).

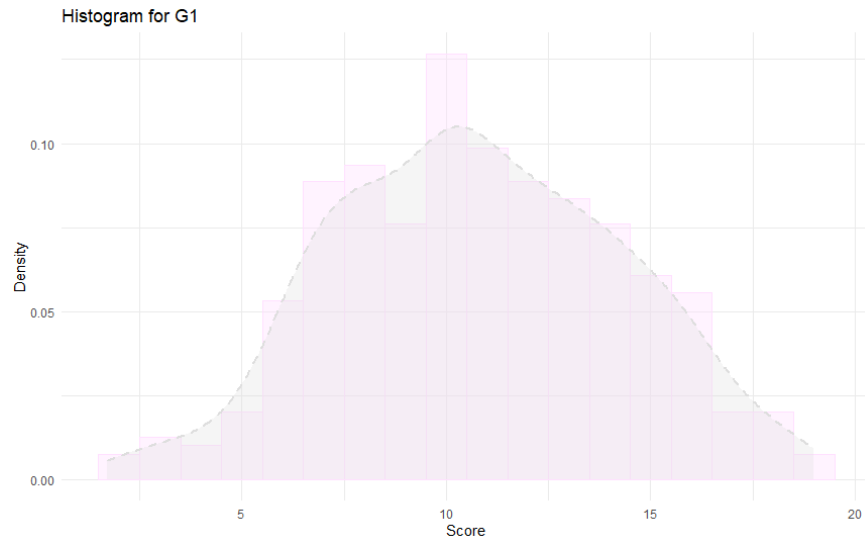


Figure 6: Histogram of G1 overlaid with density plot

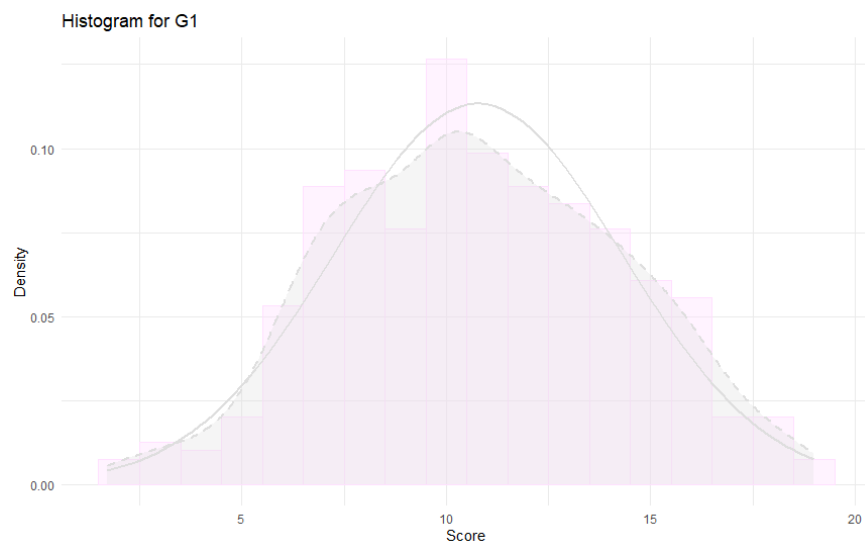


Figure 7: Histogram of G1 overlaid with fitted density plot and MLE density plot

**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   7.000   7.000   7.000
1st Qu.: 98.5   MS: 46  M:187 1st Qu.:16.0   8.000   8.000   8.000
Median :197.0                                9.000   9.000   9.000
Mean   :197.0                                9.709   9.709   9.709
3rd Qu.:295.5                                10.000  10.000  10.000
Max.   :394.0                                12.000  12.000  12.000

  studytime failures health      G1
Min.   :1.000 min. :0.0000 min.   :1.00  7.000
1st Qu.:1.000 1st Qu.:0.0000 1st Qu.:13.00  8.000
Median :2.000 Median :0.0000 Median :4.00  9.000
Mean   :2.035 Mean   :0.3342 Mean   :3.55  9.709
3rd Qu.:12.000 3rd Qu.:10.0000 3rd Qu.:15.00 10.000
Max.   :4.000 Max.   :3.0000 Max.   :5.00  12.000
```

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.)

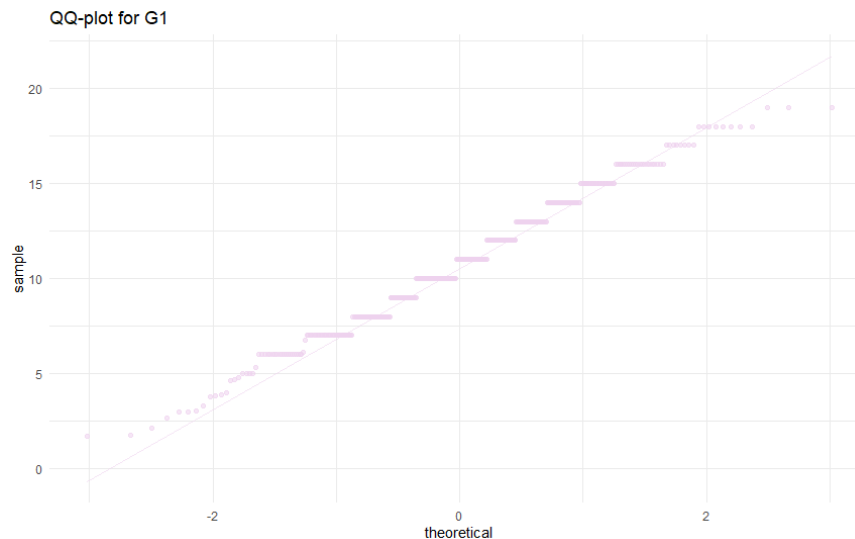


Figure 9: Normal Q-Q plot of G1

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)

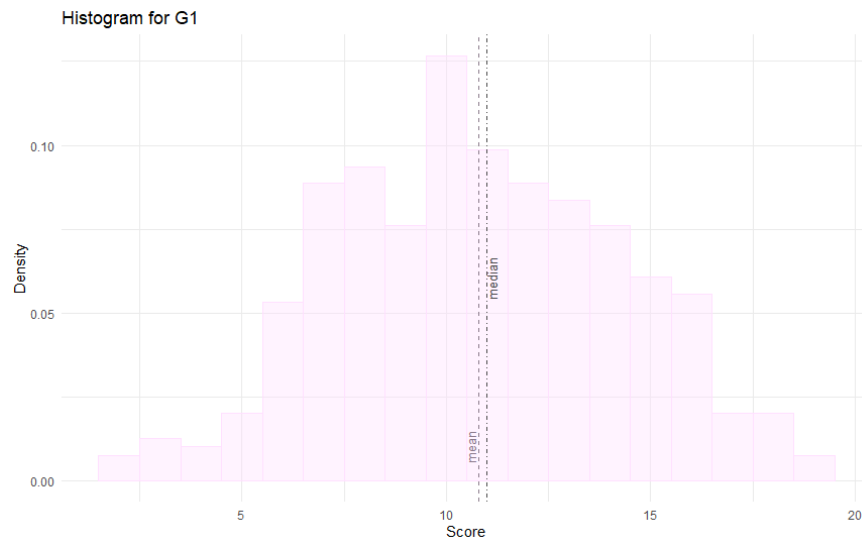


Figure 10: Median and mean marked on histogram of G1

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*, *1<sup>st</sup> quartile*, *3<sup>rd</sup> 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.

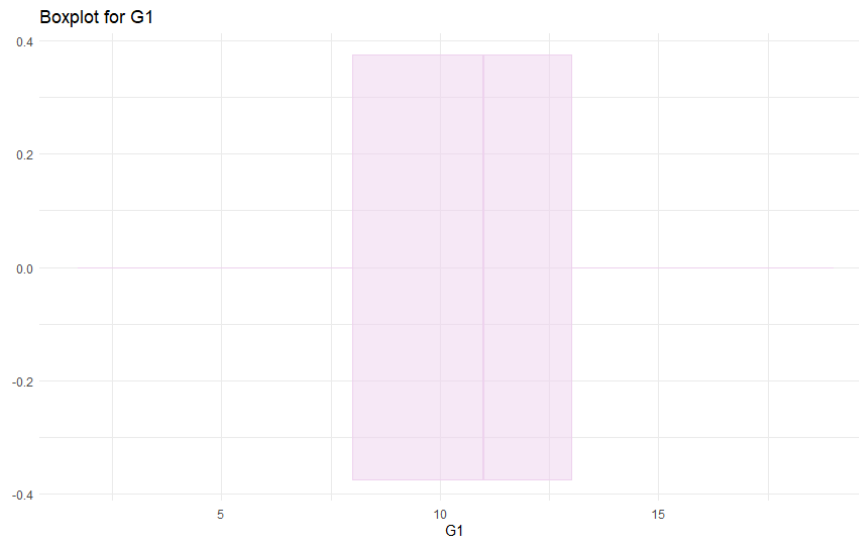


Figure 12: Boxplot of G1 to visualize outliers

```
>
> 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.*

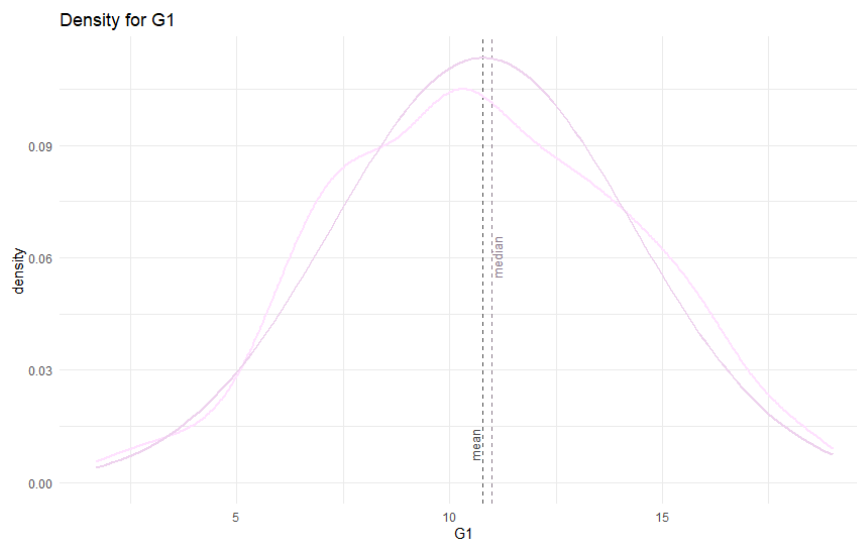


Figure 15: Median and mean marked on density of G1 - darker one is drawn using *dnorm*

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.

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

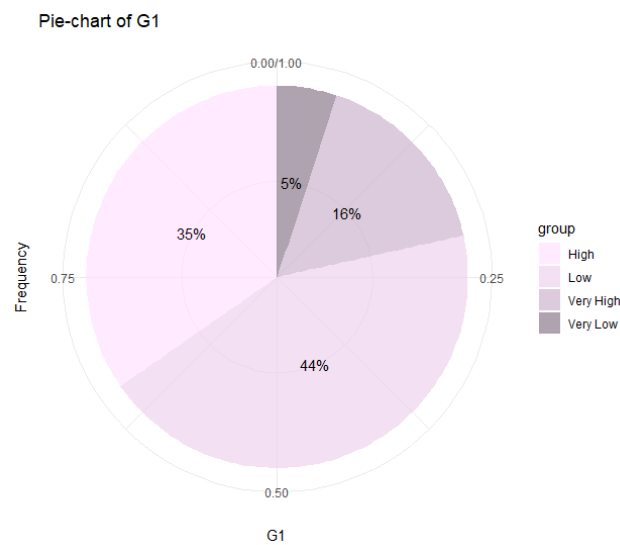
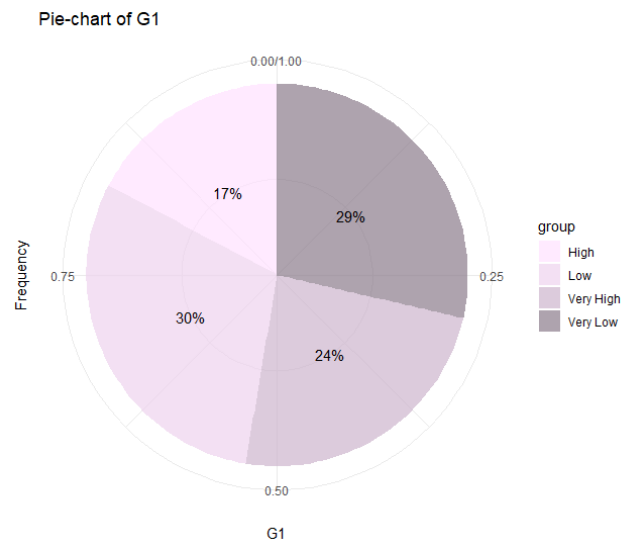


Figure 16: Piechart of G1 - 1<sup>st</sup> method

*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 - 2<sup>nd</sup> 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: 0<sup>th</sup>, 25<sup>th</sup>, 50<sup>th</sup>, 75<sup>th</sup>, and 100<sup>th</sup> 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.

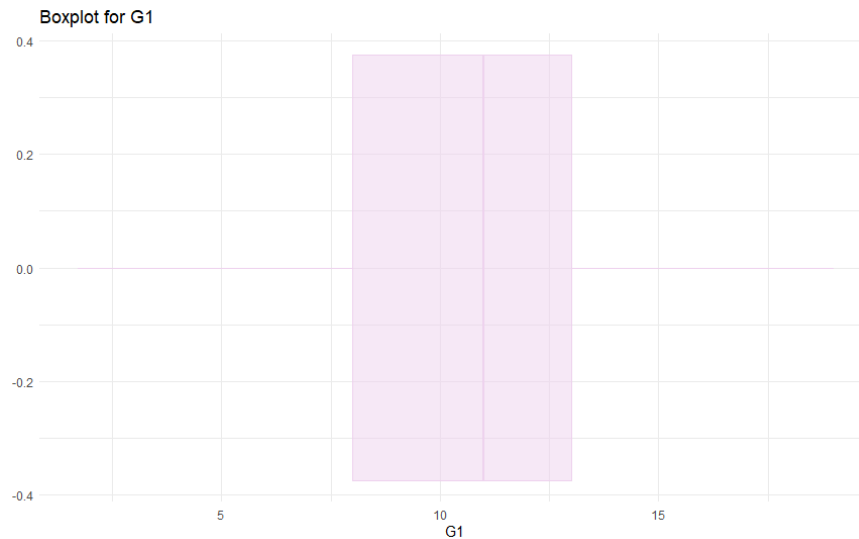


Figure 20: Boxplot of G1

```
> 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.)

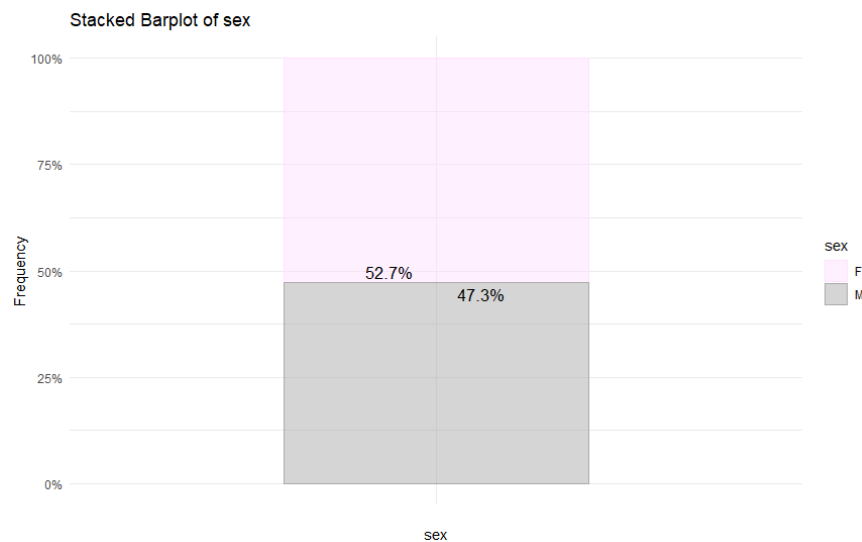


Figure 23: Stacked barplot of sex

c.

Barplots for categorical variables are like histograms for numerical variables.

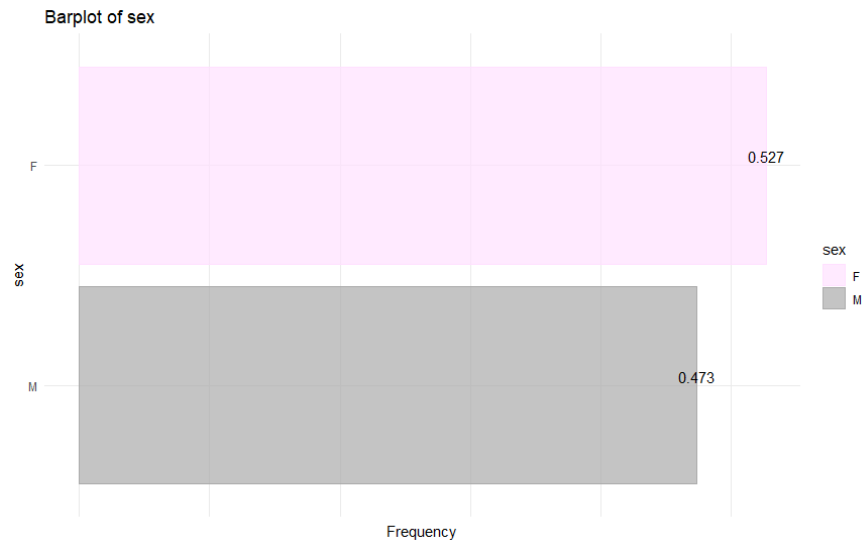


Figure 24: Horizontal barplot of sex

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.

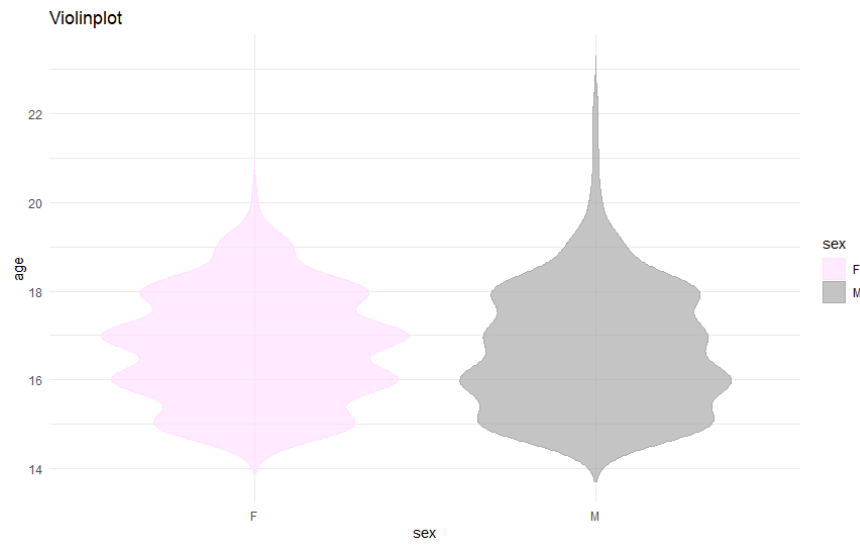


Figure 25: Violin plot of sex

### 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.

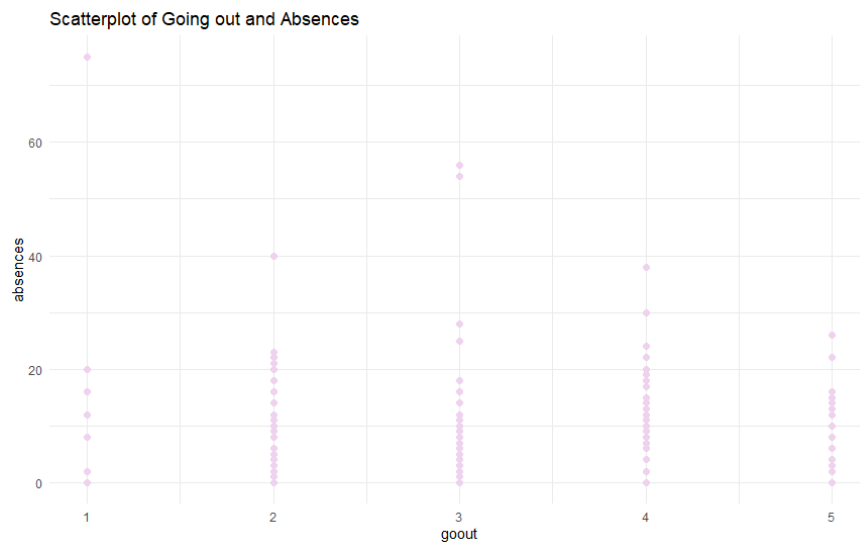


Figure 26: Scatterplot of goout and absences

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.

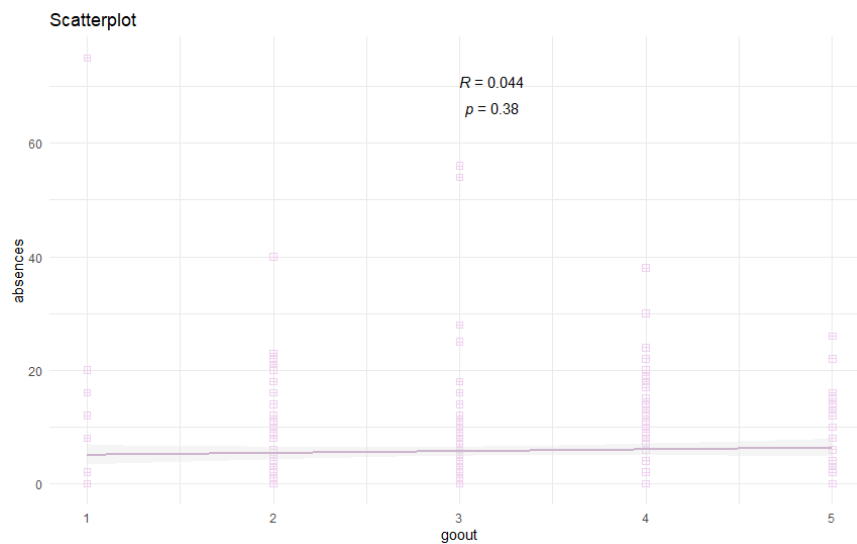


Figure 28: Scatterplot of goout and absences

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  $\rho$  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  $\rho$ .

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,  $\rho$  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*

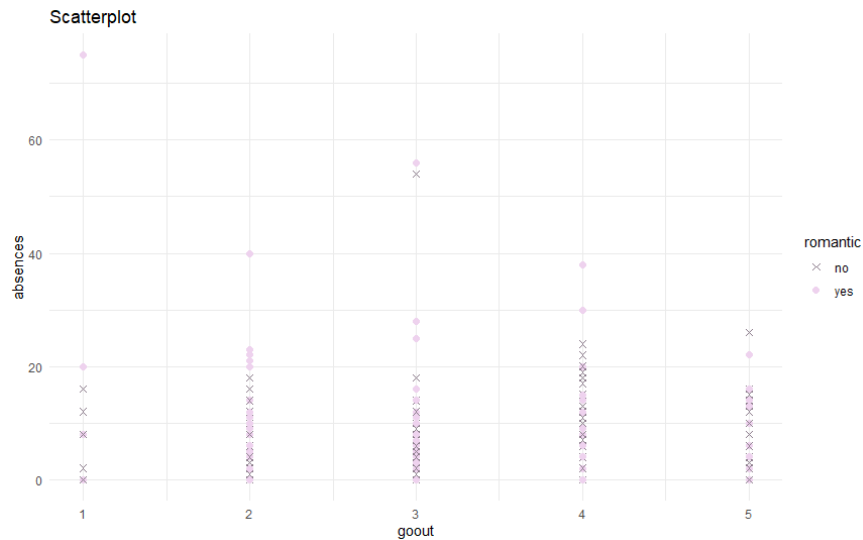


Figure 29: Scatterplot of goout and absences categorized by romantic

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.

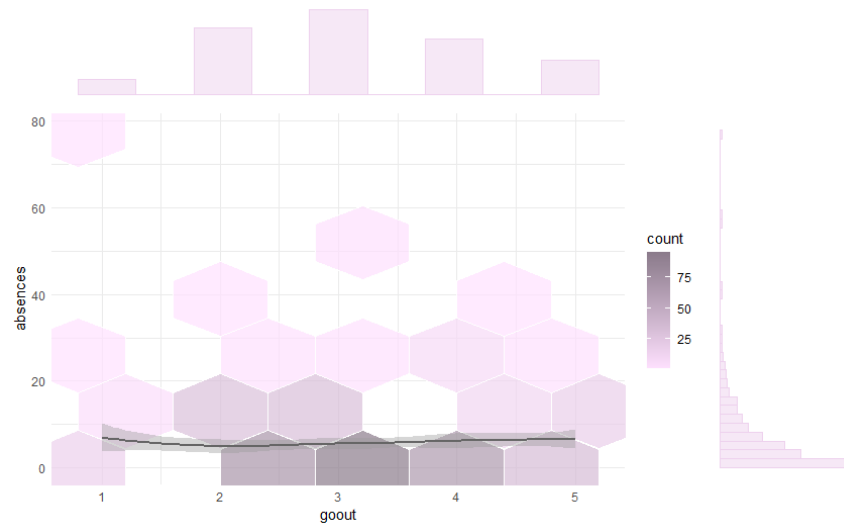


Figure 30: Hexbin plot, binsize = 5

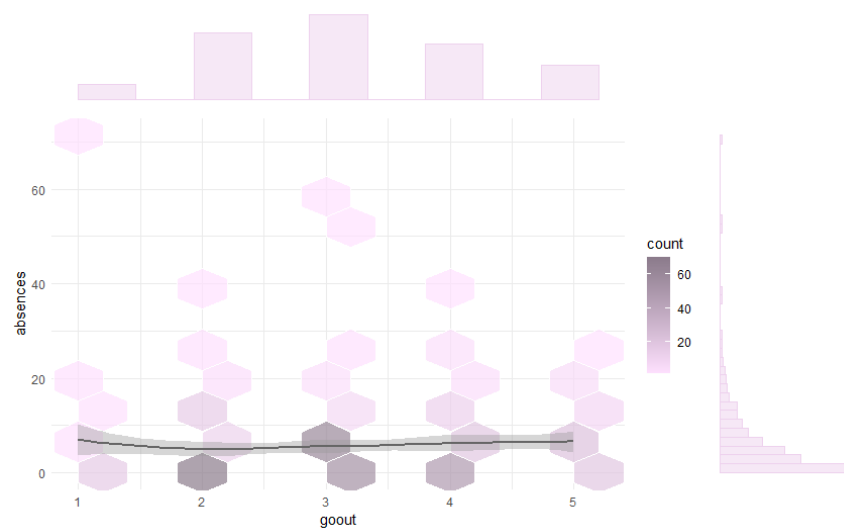


Figure 31: Hexbin plot, binsize = 10

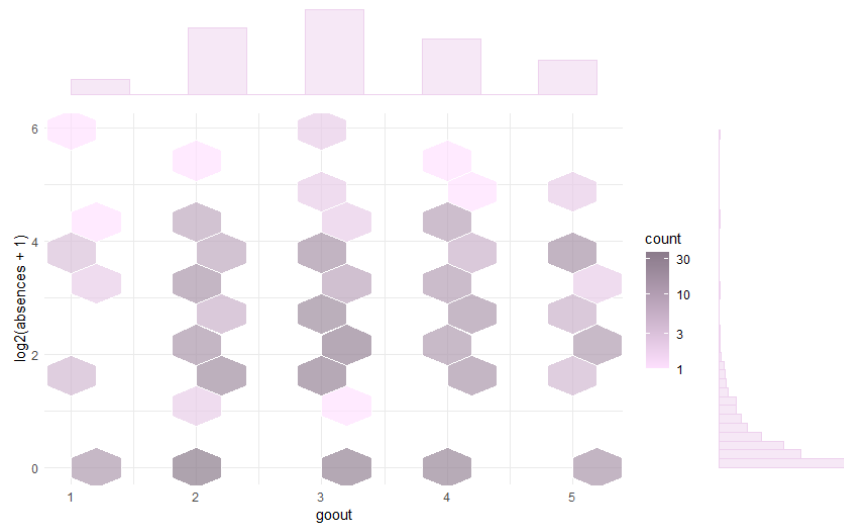


Figure 32: Hexbin plot, binsize = 10 (logarithmic)

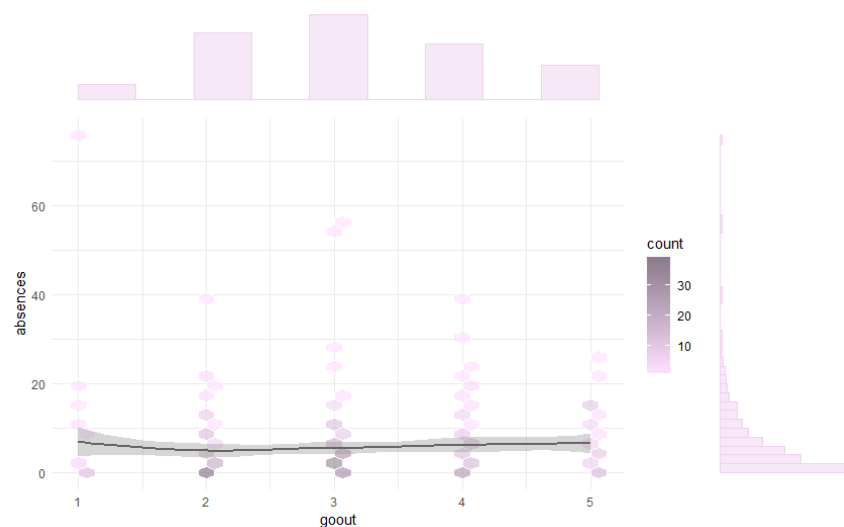


Figure 33: Hexbin plot, binsize = 30

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.

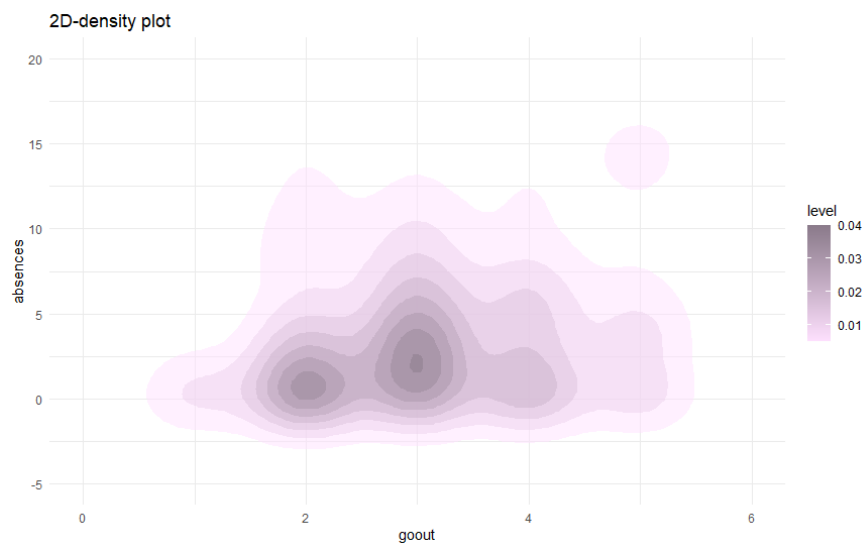


Figure 34: 2D density plot of goout and absences

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

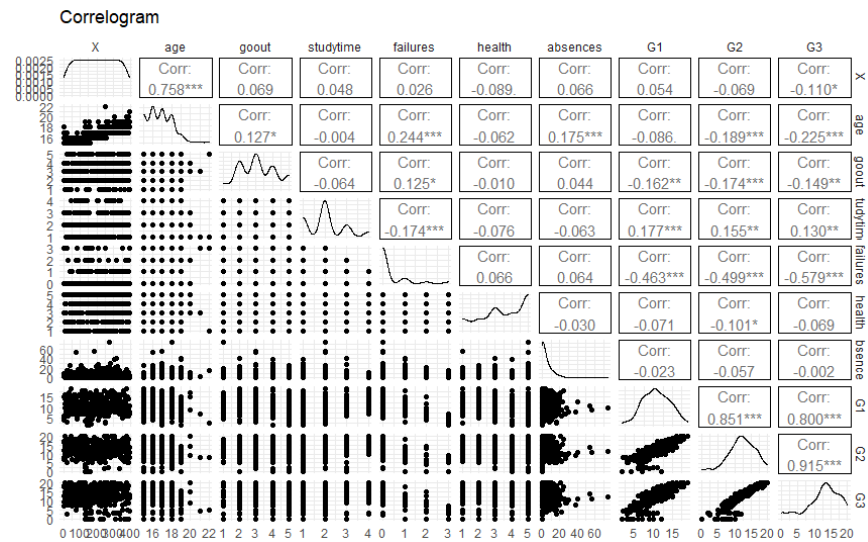


Figure 35: Bivariate Correlogram with Pearson correlation

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

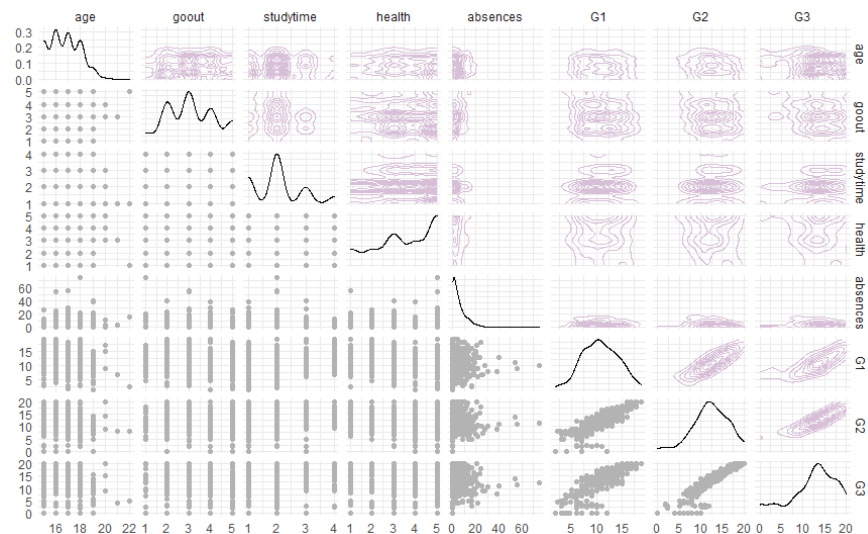


Figure 36: Bivariate Correlogram with density - scatterplot

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 :

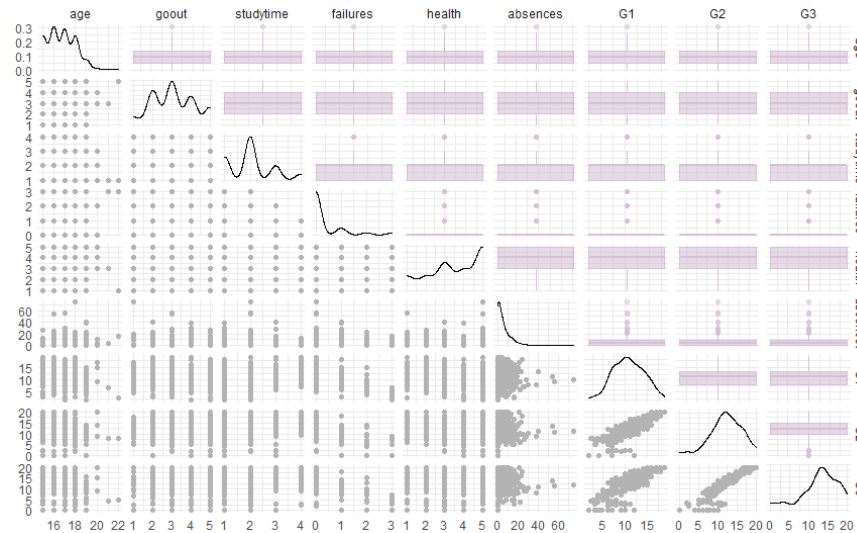


Figure 37: Bivariate Correlogram with barplot - scatterplot

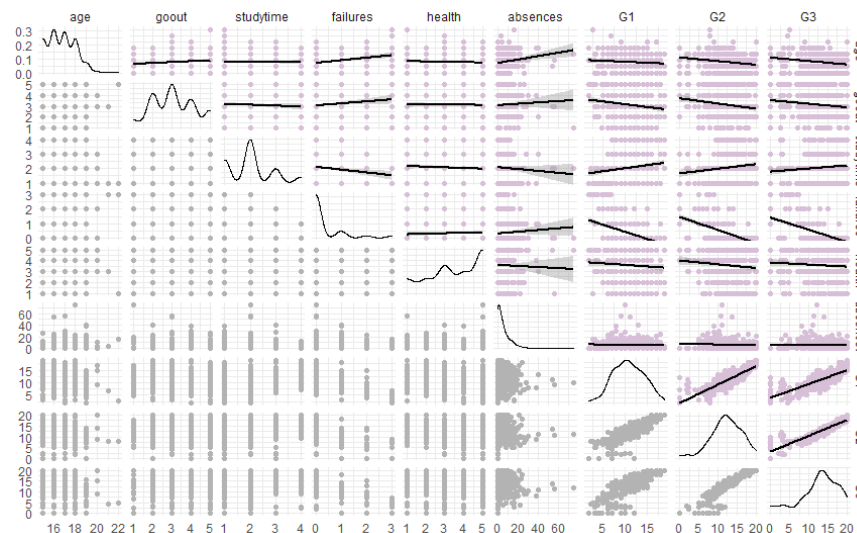


Figure 38: Bivariate Correlogram with linear association - scatterplot

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  $G3$ .

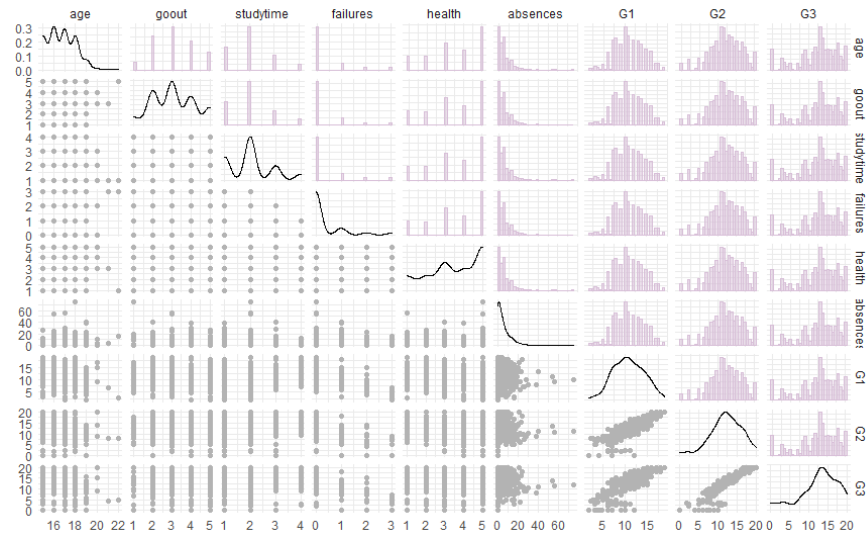


Figure 39: Bivariate Correlogram with barplot - scatterplot

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 )

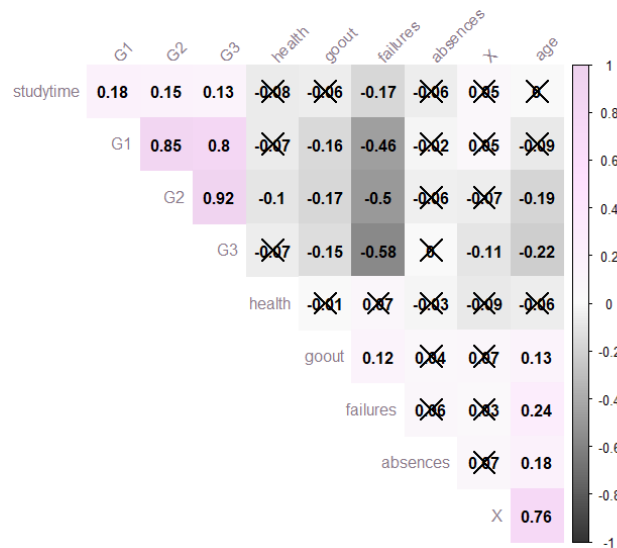


Figure 40: Heatmap correlogram of numerical values

c.



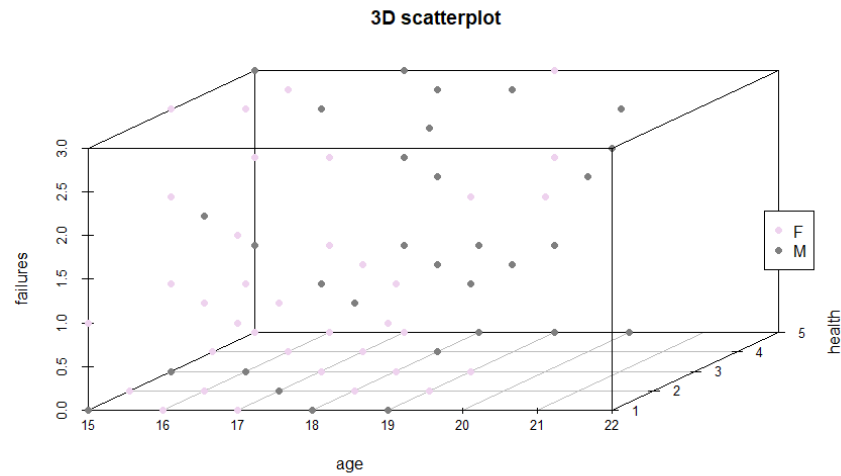


Figure 41: 3D scatterplot of age, failures and health colored by sex

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.

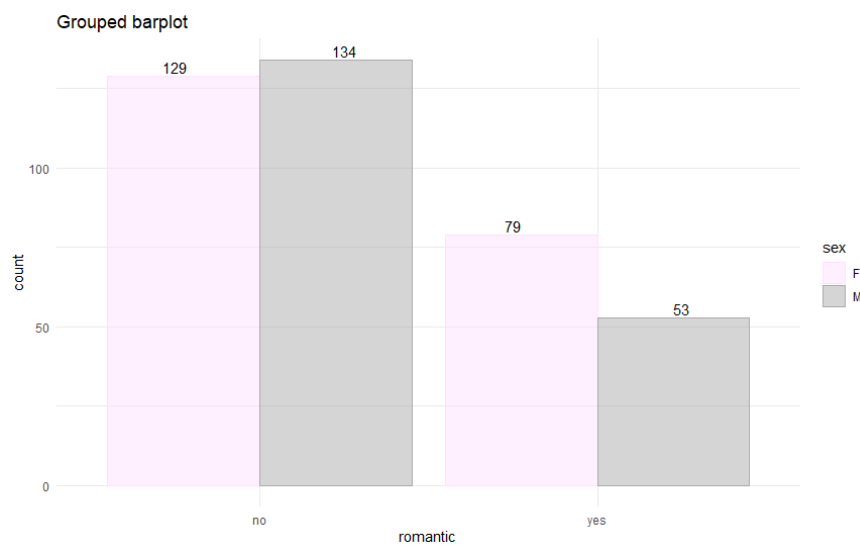


Figure 43: Grouped barplot of sex and romantic

c.

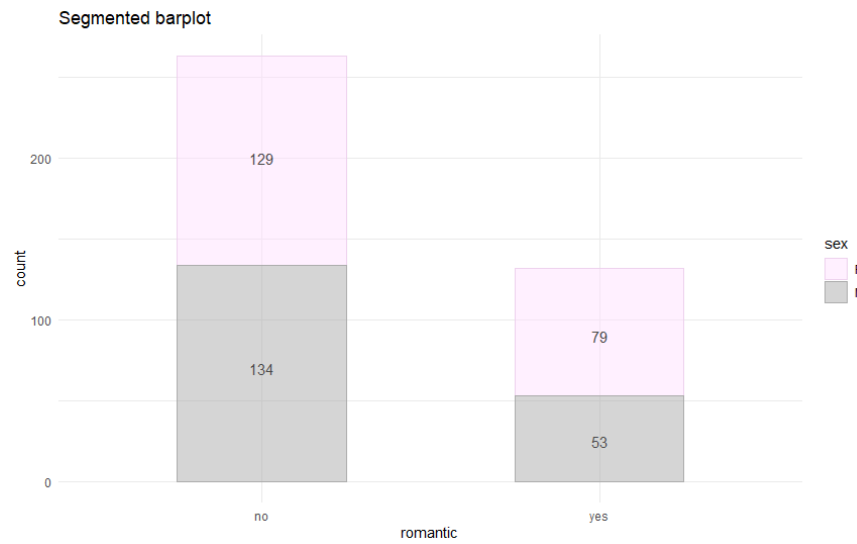


Figure 44: Segmented barplot of sex and romantic

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!

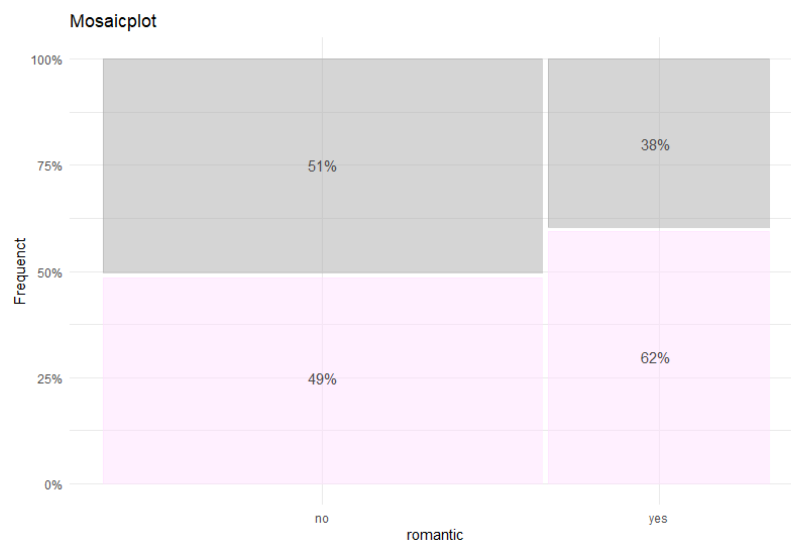


Figure 45: Mosaicplot of sex and romantic

## 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 :
  - $n < 30 \rightarrow t\text{-test}$  ,  $n > 30 \rightarrow z\text{-test}$
  - 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.

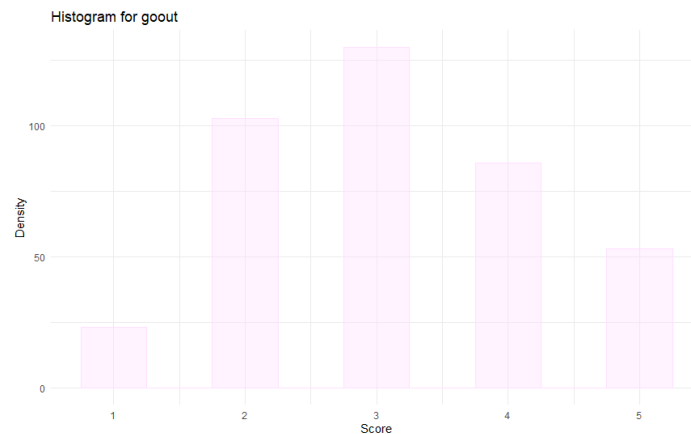


Figure 46: Histogram of goout

**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.**

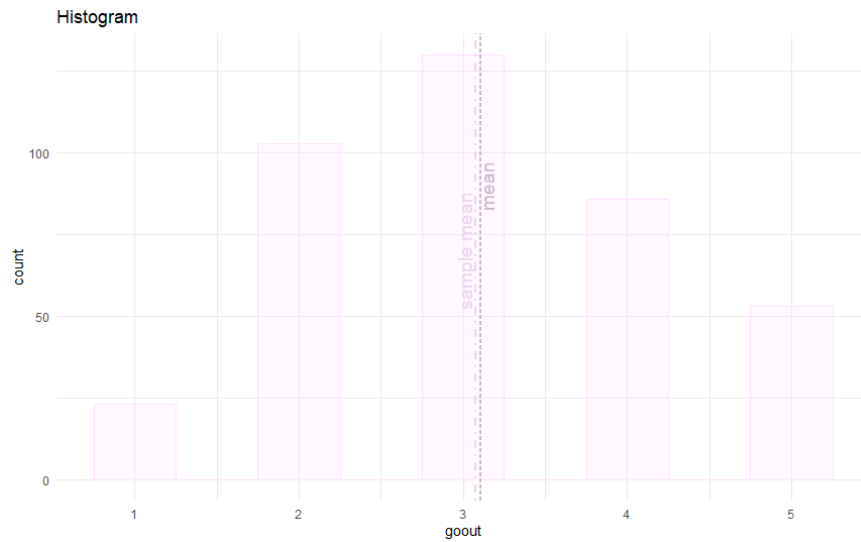


Figure 49: Histogram of goout marked with actual mean and sample mean

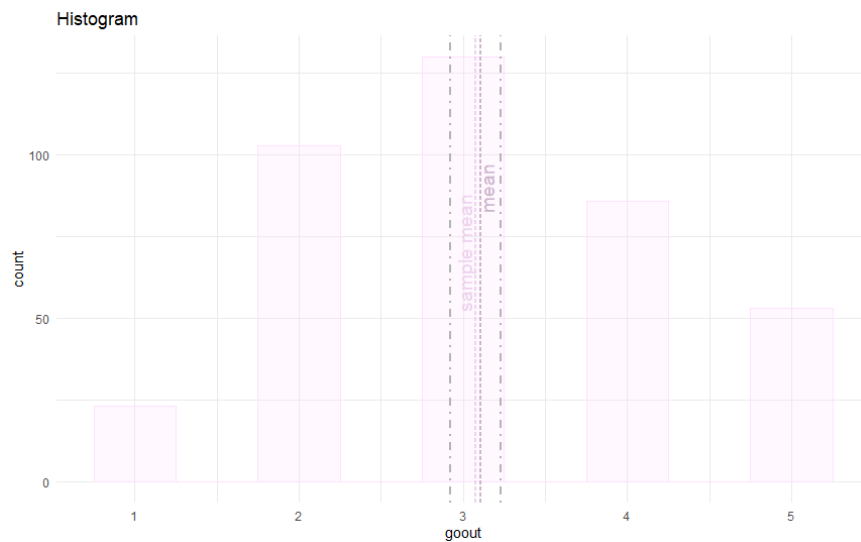


Figure 50: Histogram of goout marked with CI, actual mean and sample mean

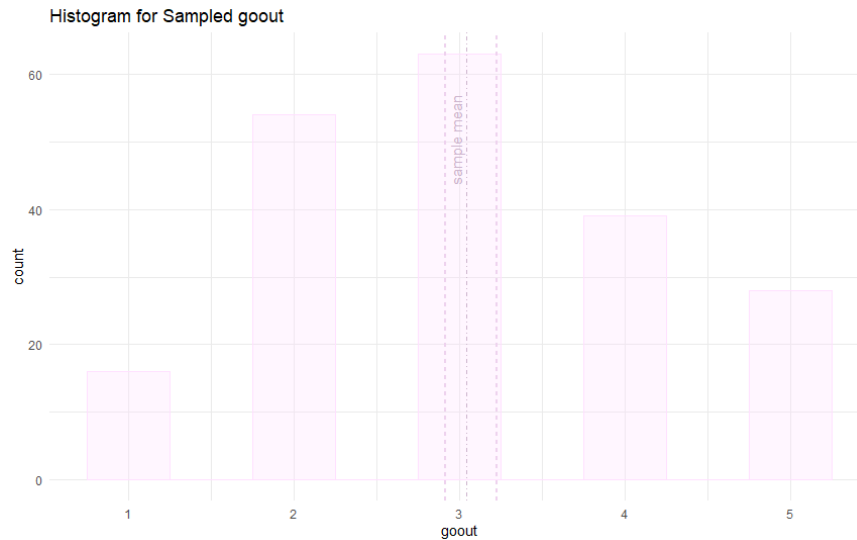


Figure 51: Histogram of sampled goout marked with CI and sample mean

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.

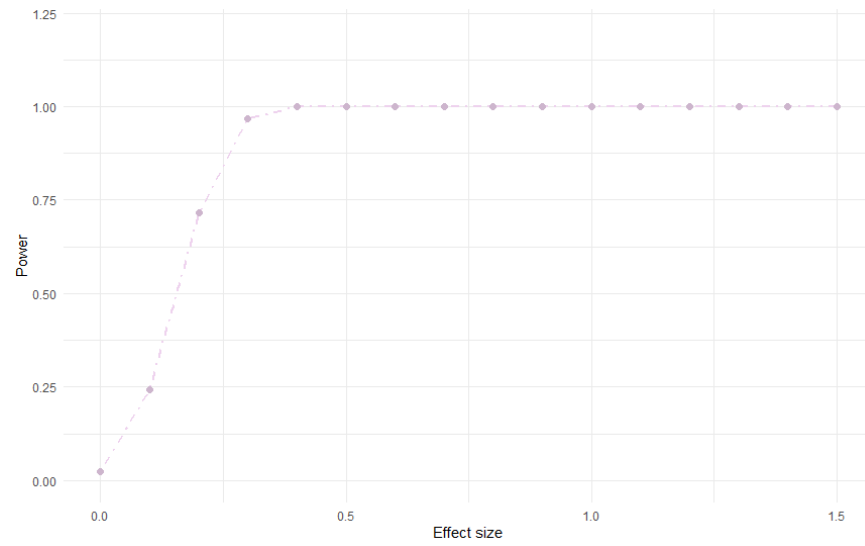


Figure 56: Relationship between effect size and power



## 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*

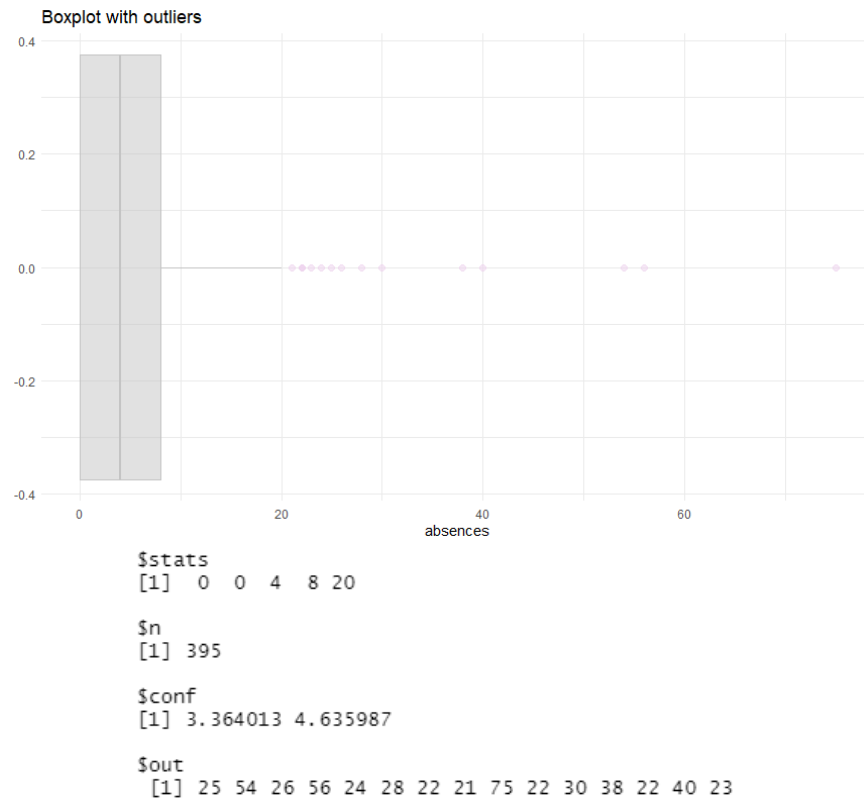


Figure 61: boxplot of absences with stat

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)

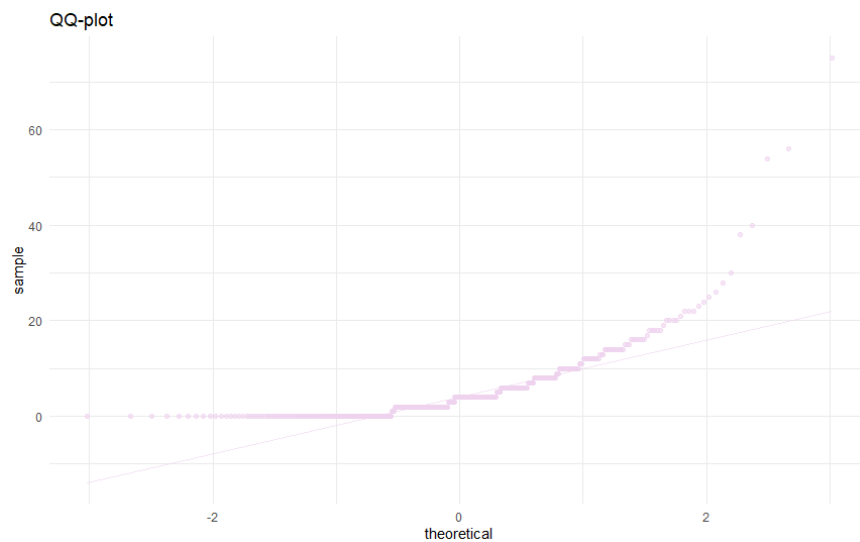
**c.**

Figure 65: QQ-plot of absences

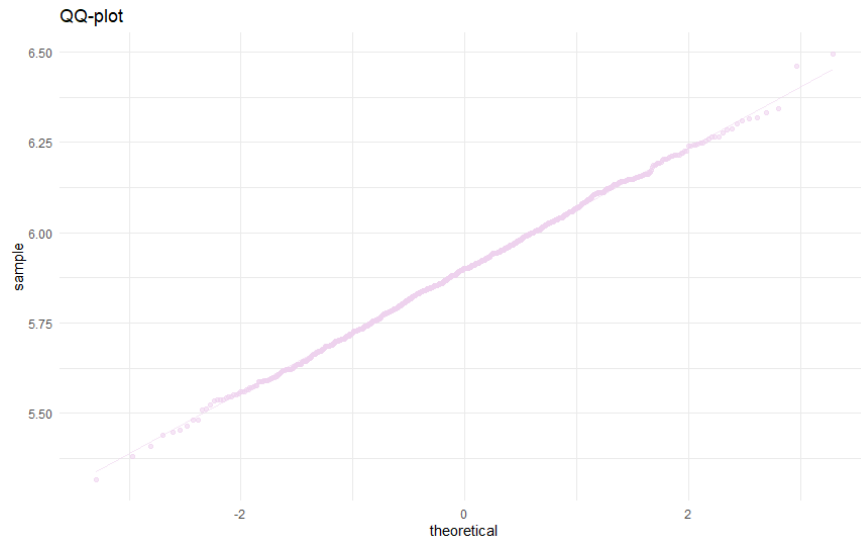


Figure 66: QQ-plot of mean of bootstrapped samples

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.

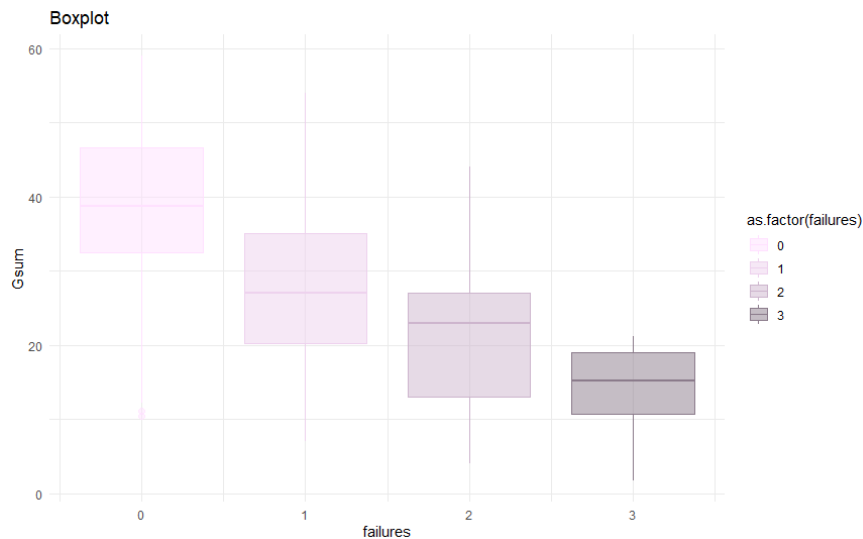


Figure 68: Boxplot grouped by number of failures

```
>
> 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  $< 0.05$ .

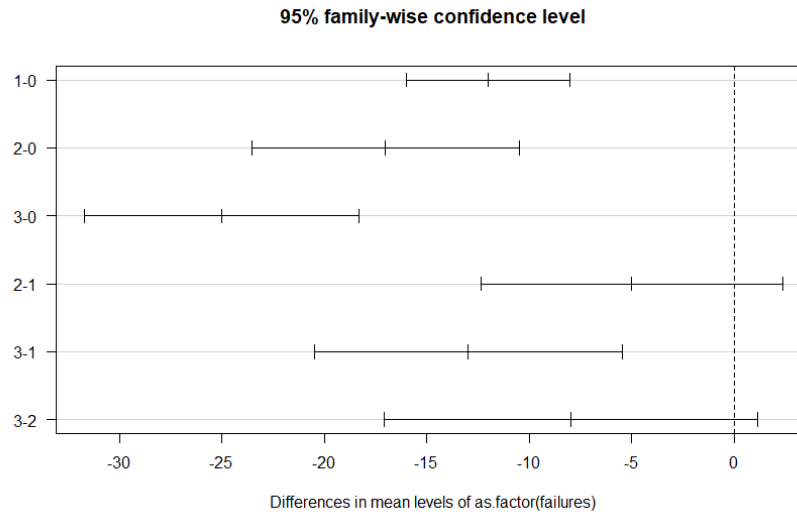


Figure 70: Pairwise confidence level(0.95%)



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