

# Report Template coursework assignment A - 2018

CS4125 Seminar Research Methodology for Data Science

*Arnar Þór Arnar (4917499), Nathan Buskulić (4947916), Mitchell Deen (4396340)*

*4/3/2019*

## Contents

<b>1</b>	<b>Part 1 - Design and set-up of true experiment</b>	<b>2</b>
1.1	The motivation for the planned research. . . . .	2
1.2	The theory underlying the research. . . . .	2
1.3	Research questions . . . . .	2
1.4	The related conceptual model . . . . .	2
1.5	Experimental Design . . . . .	2
1.6	Experimental procedure . . . . .	2
1.7	Measures . . . . .	3
1.8	Participants . . . . .	3
1.9	Suggested statistical analyses . . . . .	3
<b>2</b>	<b>Part 2 - Generalized linear models</b>	<b>3</b>
2.1	Question 1 Twitter sentiment analysis (Between groups - single factor) . . . . .	3
2.1.1	Collecting tweets, and data preparation . . . . .	3
2.1.2	Conceptual model . . . . .	4
2.1.3	Homogeneity of variance analysis . . . . .	4
2.1.4	Visual inspection . . . . .	5
2.1.5	Mean sentiments . . . . .	7
2.1.6	Linear model . . . . .	7
2.1.7	Post Hoc analysis . . . . .	7
2.1.8	Report section for a scientific publication . . . . .	8
2.2	Question 2 - Website visits (between groups - Two factors) . . . . .	8
2.2.1	Conceptual model . . . . .	8
2.2.2	Visual inspection . . . . .	8
2.2.3	Normality check . . . . .	10
2.2.4	Model analysis . . . . .	10
2.2.5	Simple effect analysis . . . . .	11
2.2.6	Report section for a scientific publication . . . . .	13
2.3	Question 3 - Linear regression analysis . . . . .	13
2.3.1	Conceptual model . . . . .	13
2.3.2	Visual inspection . . . . .	13
2.3.3	Scatter plot . . . . .	16
2.3.4	Linear regression . . . . .	18
2.3.5	Examine assumption . . . . .	19
2.3.6	Impact analysis of individual cases . . . . .	20
2.3.7	Report section for a scientific publication . . . . .	22
2.4	Question 4 - Logistic regression analysis . . . . .	23
2.4.1	Conceptual model . . . . .	23
2.4.2	Logistic regression . . . . .	23
2.4.3	Crosstable predicted and observed responses . . . . .	23
2.4.4	Report section for a scientific publication . . . . .	23
<b>3</b>	<b>Part 3 - Multilevel model</b>	<b>23</b>

3.1	Visual inspection . . . . .	23
3.2	Multilevel analysis . . . . .	25
3.3	Report section for a scientific publication . . . . .	27

# 1 Part 1 - Design and set-up of true experiment

## 1.1 The motivation for the planned research.

(Max 250 words) The coffee is today the most consumed drink in the world and it is told to increase your performance and concentration. We want to challenge this idea and verify scientifically if this is a valid idea. We want to test how coffee consumption (and the level of caffeine inside) affect the result of an IQ test. We are most interesting and seeing what the affect is on TU Delft students like ourselves. So the participants will be recruited from the TU Delft student body. # Add that we are doing that on tudelft student

## 1.2 The theory underlying the research.

(Max 250 words) Preferable based on theories reported in literature There is a large body of literature available on the effects of caffeine on the performance in cognitive tasks. Literature generally supports the idea that coffee improves this performance, see e.g. (Jarvis, 1993; Nehlig, 2010; Rogers et al., 2008). In a brief survey of the relevant literature we did not find any studies specifically addressing students. We would like to investigate this part of the population in more detail.

Jarvis, M. J. (1993). Does caffeine intake enhance absolute levels of cognitive performance?. *Psychopharmacology*, 110(1-2), 45-52. Rogers, P. J., Smith, J. E., Heatherley, S. V., & Pleydell-Pearce, C. W. (2008). Time for tea: mood, blood pressure and cognitive performance effects of caffeine and theanine administered alone and together. *Psychopharmacology*, 195(4), 569. Nehlig, A. (2010). Is caffeine a cognitive enhancer?. *Journal of Alzheimer's Disease*, 20(s1), S85-S94.

## 1.3 Research questions

The research question that will be examined in the experiment (or alternatively the hypothesis that will be tested in the experiment) Does coffee consumption increases IQ test score #Should we add the caffeine level ?

## 1.4 The related conceptual model

Independent variable(s) -> caffeine consumption Dependent variable -> score at IQ test. Mediating variable (at least 1) -> sleepiness feeling Moderating variable (at least 1) -> prior coffee consumption habit/caffeine tolerance

## 1.5 Experimental Design

Note that the study should have a true experimental design The experiment is a two groups, post test only, randomized controlled trail.

## 1.6 Experimental procedure

Describe how the experiment will be executed step by step The participants will be separated into two groups randomly. One group will do the IQ test without any prior coffee consumption while the second group will do the test half an hour after coffee consumption. In the coffee consumption groups, participants will be separated in three subgroups where they will get coffee with different caffeine level. This will allow us to measure the general impact of drinking coffee on an IQ test but it will also allow us to test the difference between each caffeine level.

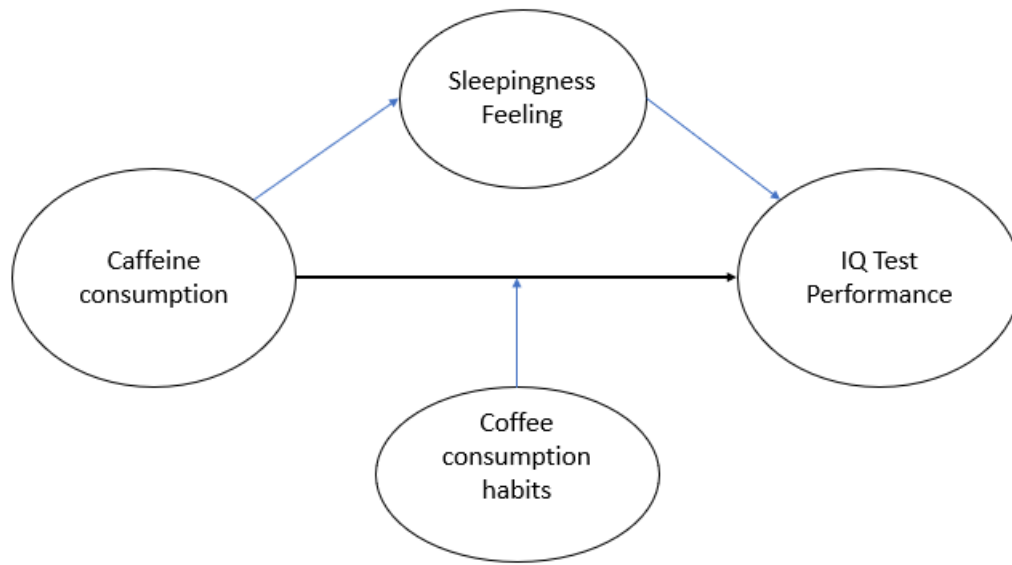


Figure 1: Conceptual model

## 1.7 Measures

Describe the measure that will be used The Coffe consumption will be measured in ml. The perfomance in an IQ test will in a simple integer number on the scale from 0-200 where the mean is around 100. Sleepingness will be given by the participants on the scale from 0-10 where 10 means the highest level of sleepingness. The amount of caffeine will be measured in mg. Prior coffeedrinking habits will be given by participants. They will be asked how much coffee they typically drink on a normal day.

## 1.8 Participants

Describe which participants will recruit in the study and how they will be recruited

Since we are just going to make this experiment on the effects of coffee consumption on students at TU Delft we need to find participants from that group of people. Emails will be sent out to the student body explaining the theory of the experiments and willing volunteers asked to fill in a form. We will try to contact an external company of some sort to get some credit or coupons that we can give to participants as a reward for helping out.

## 1.9 Suggested statistical analyses

Describe the statistical test you suggest to care out on the collected data We will use a one way Analysis of Variance (ANOVA) test between group. Indeed, since the IQ test is follows a gaussian distribution, we just want to compare the mean of each group.

# 2 Part 2 - Generalized linear models

## 2.1 Question 1 Twitter sentiment analysis (Between groups - single factor)

### 2.1.1 Collecting tweets, and data preparation

We collected Tweets for the three celebrities Beyonce, Madonna and " MickaelJackson. The code can be found in the markdown file.

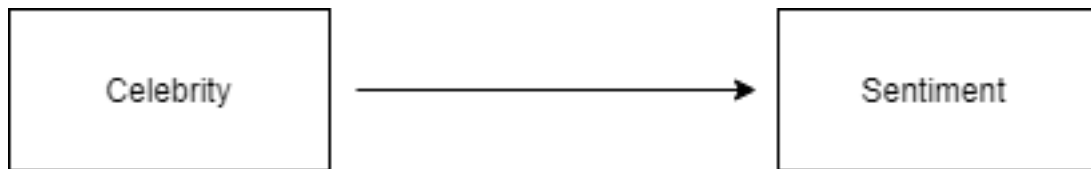


Figure 2: Conceptual model

### 2.1.2 Conceptual model

Make a conceptual model for the following research question: Is there a difference in the sentiment of the tweets related to the different celebrities?

We can see that the sentiment of tweets related to different celebrity is directly connected to the celebrity itself. Therefore the conceptual model is very simple consisting of two variables, “Celebrity” and “Sentiment”.

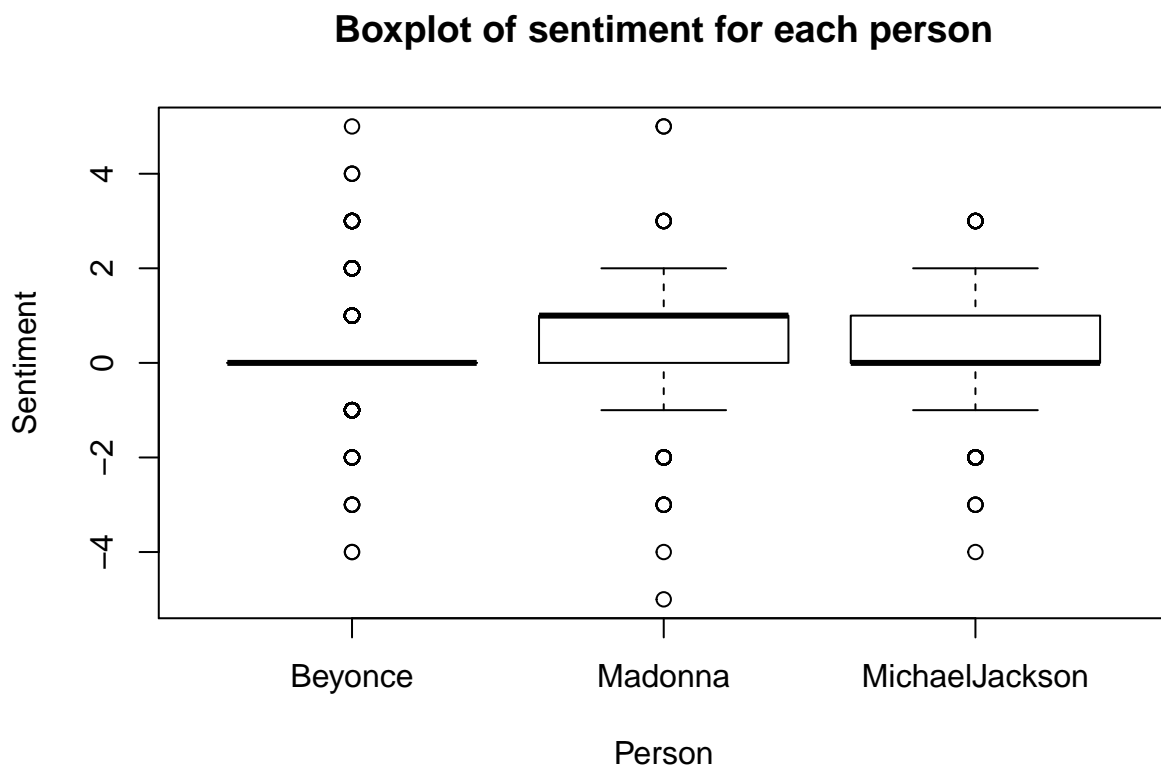
### 2.1.3 Homogeneity of variance analysis

Here we will analyze the homogeneity of variance of sentiments of the tweets of the different celebrities

Let's start by looking at how the boxplot looks for each person and the relevant sentiment score. We can already see from looking at the boxplot that the variance does not seem to be the same for all celebrities. Madonna seems to have the broadest spectrum and the median line hits a bit different places depending on the celebrities.

```

#this was not here in the intermediate report.
#include your code and output in the document
boxplot(score ~ Person, data=semFrame, main="Boxplot of sentiment for each person",
        xlab="Person", ylab="Sentiment")
  
```



```

levene = leveneTest( semFrame$score, semFrame$Person, center = median)
  
```

The Levene test results in a very low p-value  $8.0991132 \times 10^{-7}$ . Therefore the hypothesis of equal variances is rejected and it is concluded that there is a difference between the variances in the population. Therefore the variance is not considered to be homogeneous.

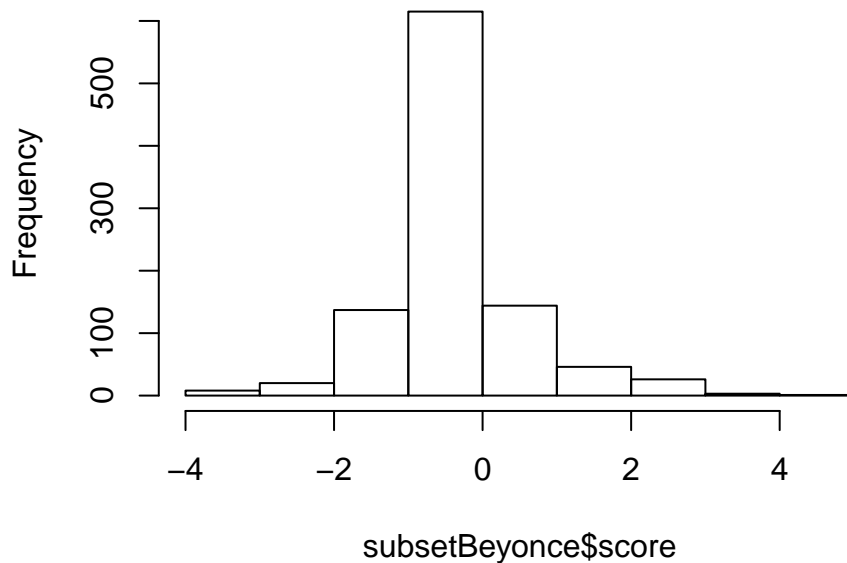
#### 2.1.4 Visual inspection

Looking at the figures here below we see that the sentiment scores for all the celebrities follow a very similar distribution that looks a lot like a normal distribution. But by inspecting the histograms we can see that the distribution is not entirely the same. Therefore we will do a further inspection to see how the distributions differ from each other.

*#include your code and output in the document*

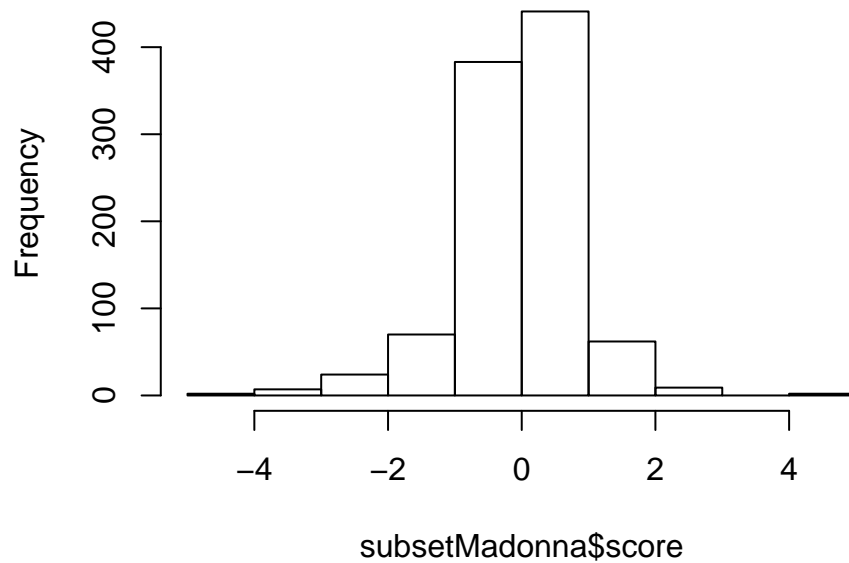
```
hist(subsetBeyonce$score)
```

**Histogram of subsetBeyonce\$score**



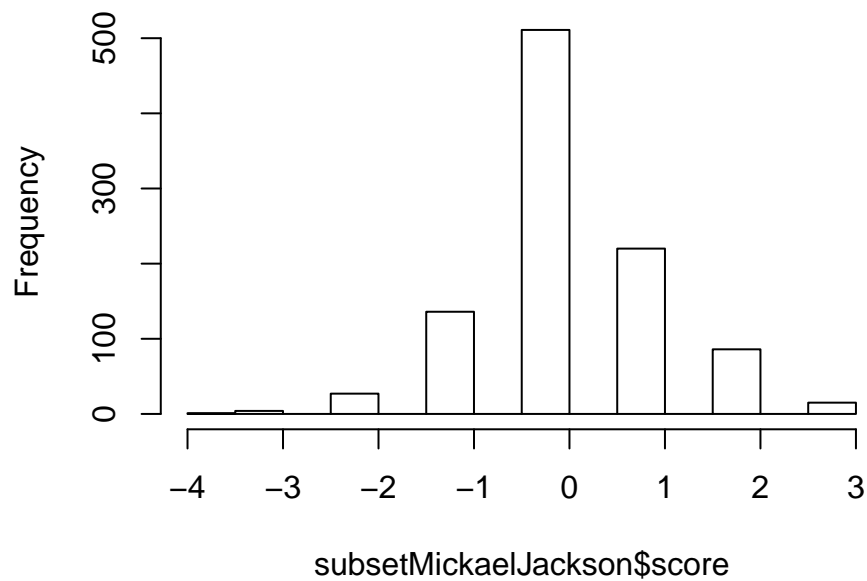
```
hist(subsetMadonna$score)
```

**Histogram of subsetMadonna\$score**



```
hist(subsetMickaelJackson$score)
```

**Histogram of subsetMickaelJackson\$score**

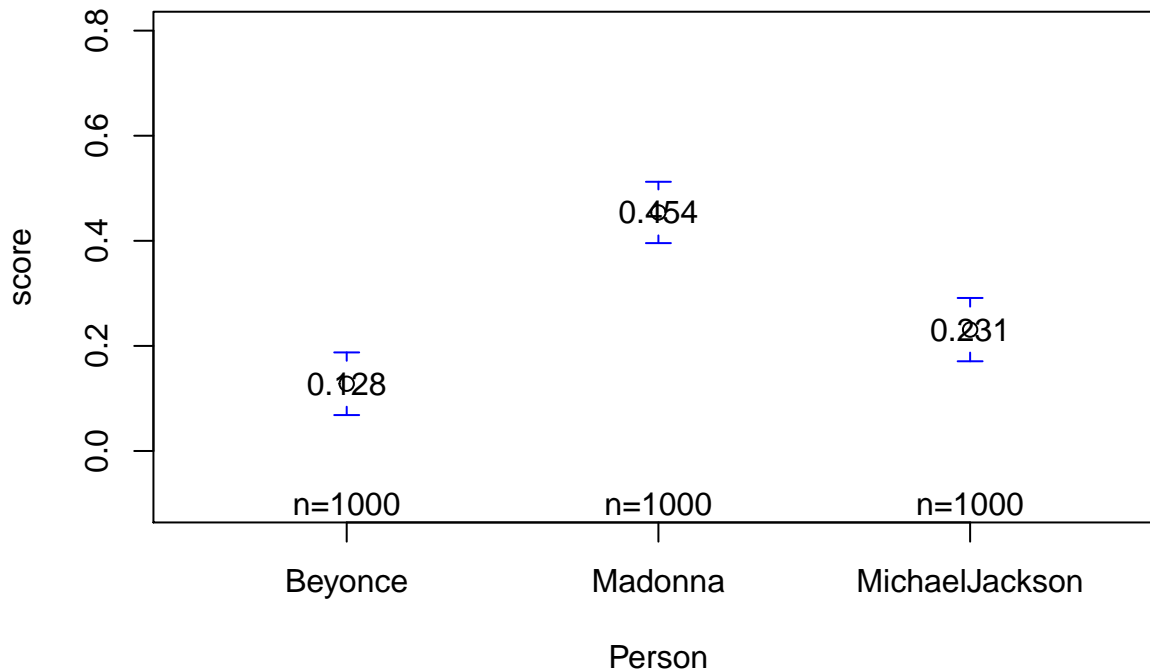


```
meanBeyonce<-mean(subsetBeyonce$score)
meanMadonna<-mean(subsetMadonna$score)
meanMickaelJackson<-mean(subsetMickaelJackson$score)
stdBeyonce<-sd(subsetBeyonce$score)
stdMadonna<-sd(subsetMadonna$score)
stdMickaelJackson<-sd(subsetMickaelJackson$score)
```

### 2.1.5 Mean sentiments

Here below we plot the means of each class using `plotmeans` from the package `gplots`. We can see that the mean for Beyonce is 0.128, for Madonna is 0.454 and for Michael Jackson it is 0.231. Where a lower value means that the sentiment analysis is more negative. Just by looking at the means we see that they are considerably far from each other. As well we have the standard deviations 0.958, 0.897, 0.944 in the same respective order as before.

```
plotmeans(score ~ Person, data = semFrame, mean.labels = TRUE, connect = FALSE, ylim = c(-0.1, 0.8))
```



### 2.1.6 Linear model

```
#include your code and output in the document
model0 <- lm(score ~ 1, data = semFrame) #model without predictor
model1 <- lm(score ~ Person, data = semFrame) #model with predictor
AnovaResults <- anova(model0, model1)
```

Here we have fitted two different linear models to fit our outcome. Model0 is a simple model where the sentiment score stands alone, but model1 has the celebrity in there as well. The calculated f-value,  $F(2, 2997)$  is 31.640 and the p-value is very small and well below .0001. Since the p-value is so small it indicates that the sentiment of tweets is significantly different depending on what celebrity is mentioned in the tweet.

### 2.1.7 Post Hoc analysis

Now a post-hoc analysis is performed to examine which of the tweets differ from other celebrity tweets

```
#include your code and output in the document
BonferroniResults <- pairwise.t.test(semFrame$score, semFrame$Person, paired = FALSE, p.adjust.method =
BonferroniP <- BonferroniResults$p.value
BonferroniP
```

```
##              Beyonce      Madonna
## Madonna      1.062510e-13      NA
## MichaelJackson 4.861709e-02 6.105306e-07
```

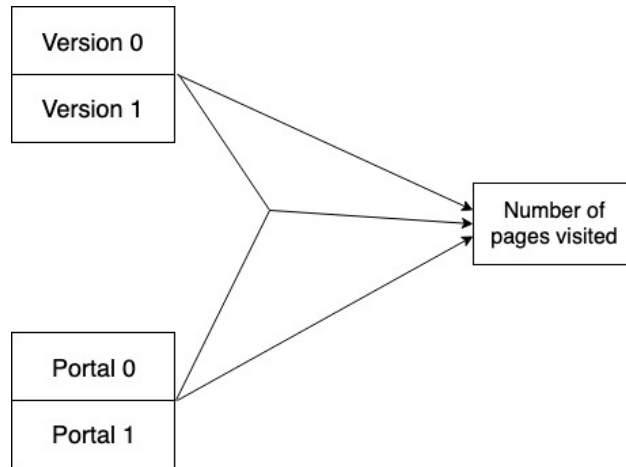


Figure 3: Model for the webvisit question

We chose to use the Bonferroni correction to conduct this post-hoc analysis. There the p-values are multiplied by the number of comparisons.

According to results here above all of the celebrity pair comparisons have a low p-value that indicates again that the sentiment is dependent on what celebrity it is in relation to.

### 2.1.8 Report section for a scientific publication

We analysed sentiment scores of tweets for three different celebrities. We had scores for Madonna (M=0.459, SD=0.897), Beyonce (M=0.132, SD=0.958) and Michael Jackson (M=0.245, SD=0.944). The means and standard deviations are different between celebrities so they were inspected further.

A linear model was fitted on the number of the sentiment score, comparing the difference when taking the relative celebrity in account and not. We first conducted an Anova test and obtained the results ( $F(2,2997) = 31.65$ ,  $p < .0001$ ) which states a significant difference in the sentiment, depending on which celebrity it is for. Then a Post Hoc analysis by the means of Bonferroni was conducted. There we again got p-values ( $< .001$ ,  $.02$ ,  $< .001$ ) that show us that the difference of scores is significant depending on celebrity.

## 2.2 Question 2 - Website visits (between groups - Two factors)

### 2.2.1 Conceptual model

The model can be found in the figure below.

### 2.2.2 Visual inspection

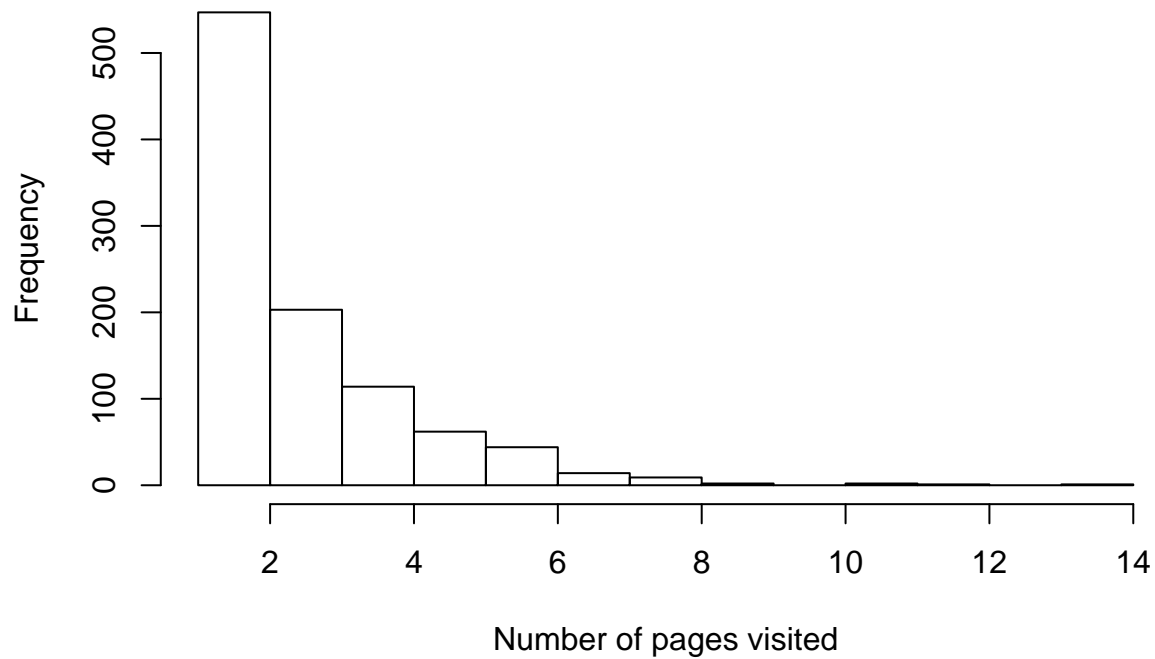
We will first examine the data and look at the distributions.

```
myData <- read.csv("webvisita.csv",header=TRUE)
# We transform into factors what need to be.
myData$user <- factor(myData$user)
myData$version <- factor(myData$version, levels=c(0:1), labels=c("old","new"))
myData$portal <- factor(myData$portal, levels=c(0:1),labels=c("consumer","company"))

hist(myData$pages, xlab="Number of pages visited", main = "Histogram of the number of pages visited")
```

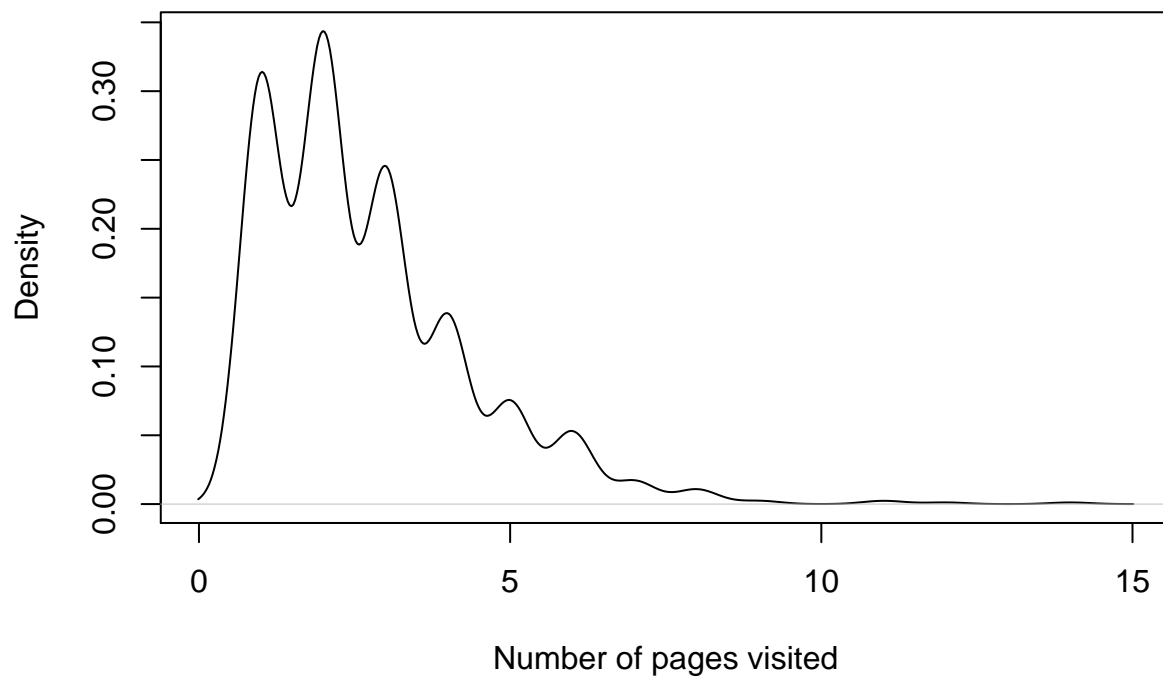


**Histogram of the number of pages visited**



```
plot(density(myData$pages), xlab="Number of pages visited", main = "density of the number of pages visited")
```

**density of the number of pages visited**



It appears that the data do not look like normally distributed, We will provide more analysis to understand these distributions.

### 2.2.3 Normality check

We can see that the data does not seem to come from a normal distribution, thus we will do a normality test.

```
shapiro.test(myData$pages)
```

```
##  
## Shapiro-Wilk normality test  
##  
## data: myData$pages  
## W = 0.8436, p-value < 2.2e-16
```

The really small p-value indicates here that there is a high probability that this data do not come from a normal distribution. From what we can visually see it seems that the data comes from a Poisson distribution, which is an important information for the Model analysis.

### 2.2.4 Model analysis

Since the data are not normally distributed, we cannot use a simple linear model which assume the normality of the data. Thus we will fit generalized linear model with a poisson distribution assumption.

```
# We create all the different models  
model0 <- glm(pages ~ 1, data=myData, family="poisson")  
model1 <- glm(pages ~ version, data=myData, family="poisson")  
model2 <- glm(pages ~ portal, data=myData, family="poisson")  
model3 <- glm(pages ~ version + portal, data=myData, family="poisson")  
model4 <- glm(pages ~ version + portal + version:portal, data = myData, family="poisson")
```

Since we are using generalized models, we cannot use an F-test and we decided to use a Chi-Square test instead in our analysis.

```
pander(anova(model0,model1,test="Chisq"),caption = "Version as main effect on the number of pages visited")
```

Table 1: Version as main effect on the number of pages visited

Resid. Df	Resid. Dev	Df	Deviance	Pr(>Chi)
998	938.6	NA	NA	NA
997	931.2	1	7.36	0.006671

```
pander(anova(model0,model2,test="Chisq"),caption = "Portal as main effect on the number of pages visited")
```

Table 2: Portal as main effect on the number of pages visited

Resid. Df	Resid. Dev	Df	Deviance	Pr(>Chi)
998	938.6	NA	NA	NA
997	878.2	1	60.39	7.793e-15

```
pander(anova(model3,model4,test="Chisq"),caption = "Interaction effect on top of the two main effect")
```

Table 3: Interaction effect on top of the two main effect

Resid. Df	Resid. Dev	Df	Deviance	Pr(>Chi)
996	872.4	NA	NA	NA
995	843.5	1	28.98	7.318e-08

```
pander(anova(model4, test="Chisq"),caption = "Effect of version, portal and interaction effect on the n
```

Table 4: Effect of version, portal and interaction effect on the number of pages visited

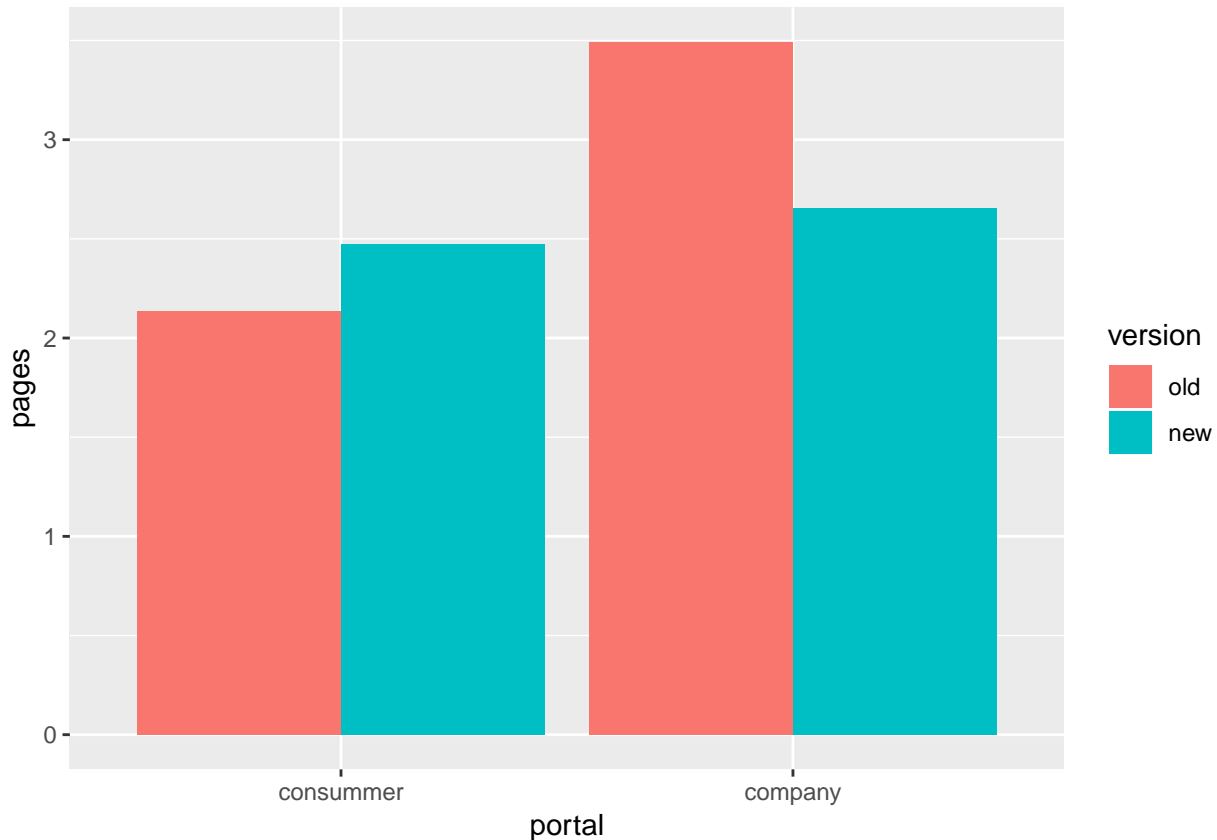
	Df	Deviance	Resid. Df	Resid. Dev	Pr(>Chi)
<b>NULL</b>	NA	NA	998	938.6	NA
<b>version</b>	1	7.36	997	931.2	0.006671
<b>portal</b>	1	58.78	996	872.4	1.764e-14
<b>version:portal</b>	1	28.98	995	843.5	7.318e-08

We can observe a significant main effect of the version ( $p < 0.01$ ) and the portal ( $p < 0.01$ ). We also see a significant two-way interaction effect ( $p < 0.01$ ), we will thus perform a simple effect analysis to better understand this interaction effect.

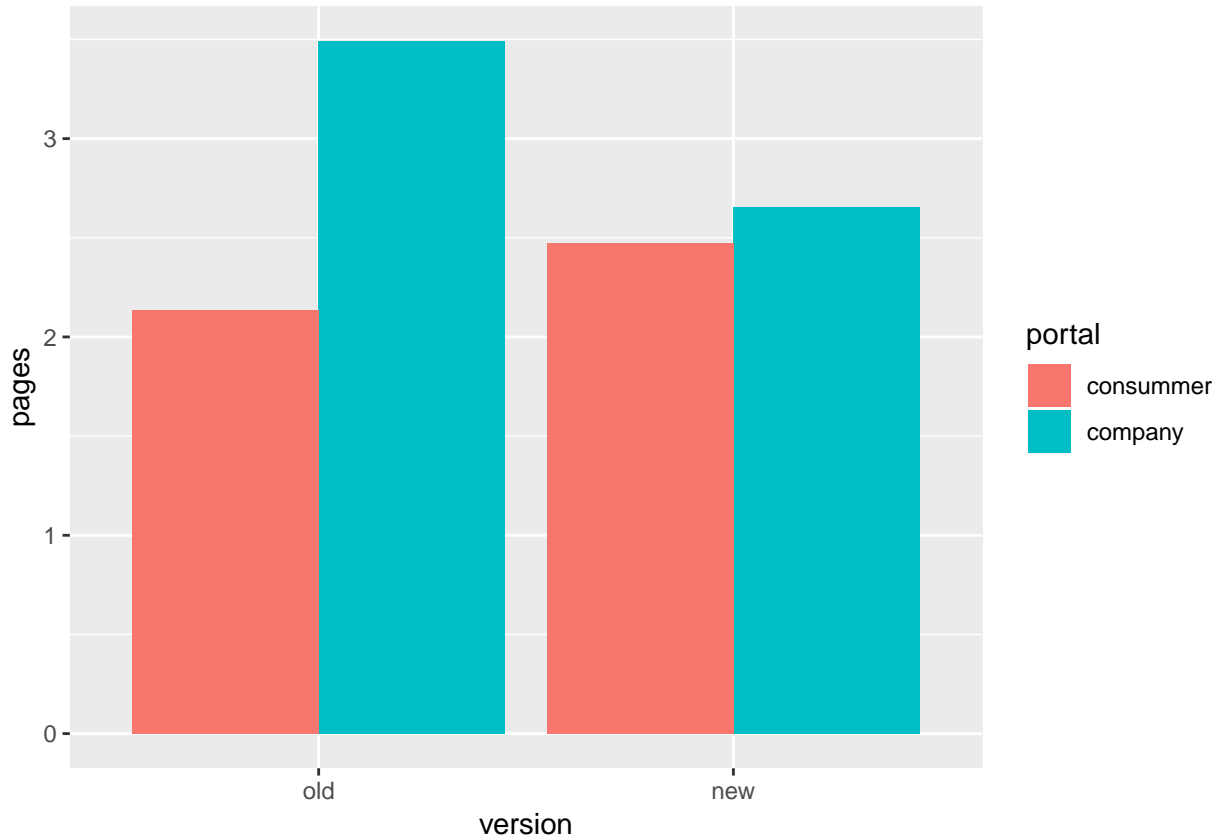
### 2.2.5 Simple effect analysis

Since we found a significant two-way interaction effect we will conduct a simple effect analysis to explore this interaction effect. We will first plot the effect of the means according to the version and the portal in two different figures to know which effect we should explore in more details.

```
bar <- ggplot(myData, aes(portal, pages, fill = version))
bar + stat_summary(fun.y = mean, geom = "bar", position="dodge")
```



```
bar <- ggplot(myData, aes(version, pages, fill = portal))
bar + stat_summary(fun.y = mean, geom = "bar", position="dodge")
```



From the two figures, we decided that the most interesting effect to explore was the one that compare the two versions and try to see if the change of version change the fact that there is a significant difference in the number of pages visited according to the portal.

```
myData$simple <- interaction(myData$version, myData$portal) #merge two factors

contrastOldVersion <- c(1,0,-1,0) #Only the old version data
contrastNewVersion <- c(0,1,0,-1) #Only the new version data

SimpleEff <- cbind(contrastOldVersion,contrastNewVersion)
contrasts(myData$simple) <- SimpleEff #now we link the two contrasts with the factor simple
pander(simpleEffectModel <- glm(pages ~ simple , data = myData, na.action = na.exclude, family = "poisson"))
```

Table 5: Simple effect analysis

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	0.9729	0.01969	49.4	0
simplecontrastOldVersion	-0.246	0.02698	-9.119	7.551e-20
simplecontrastNewVersion	-0.0347	0.0287	-1.209	0.2267
simple	-0.06349	0.03939	-1.612	0.107

The results of the simple effect analysis shows that while there were a significant difference ( $P < 0.001$ ) between the two portals in the old version, we cannot find that significant difference anymore in the new version.

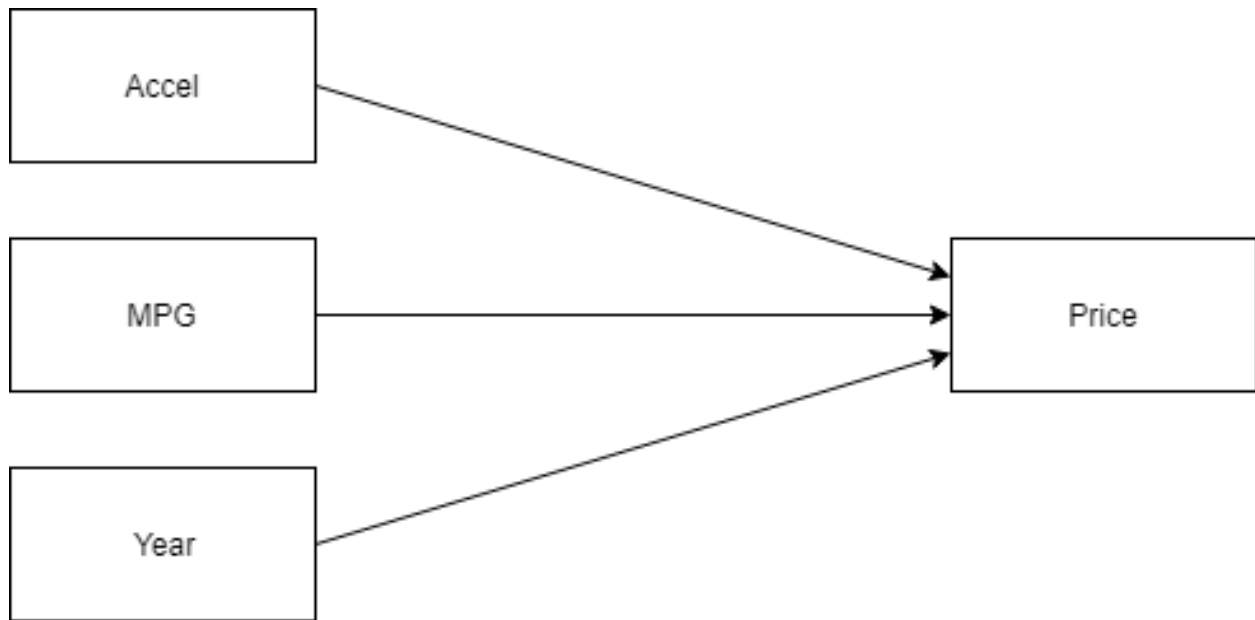


Figure 4: Conceptual model of the four considered variables

### 2.2.6 Report section for a scientific publication

A generalized linear model was fitted on the number of pages visited of a website, taking the version of the website (an old version and a new one) and the portal that was used to access the website (a portal for consumers and one for companies) as independent variables, and including a two-way interaction between these variables. The analysis found a significant main effect for the version ( $\chi^2(1,997) = 7.36$ ,  $p. < 0.01$ ) and for the portal ( $\chi^2(1,997) = 60.39$ ,  $p. < 0.001$ ). It also found a significant two-way interaction effect ( $\chi^2(1,995) = 28.98$ ,  $p. < 0.001$ ) between these two variables. The two-way interaction was further examined by a Simple Effect analysis on the effect of the version on the differences in portals. It found a significant difference for the portal given the old version ( $z = -9.1$ ,  $p. < 0.001$ ) but no difference could be found when using the new version ( $z = -1.2$ ,  $p. = 0.23$ ).

## 2.3 Question 3 - Linear regression analysis

### 2.3.1 Conceptual model

For this assignment we retrieved a data set from <http://www.stat.ufl.edu/~winner/datasets.html>. The dataset contains facts about 153 hybrid cars, including their price, year built, acceleration data and fuel consumption; those are the four quantitative variables that will be the subject of the linear model in this question. We would like to predict the price of the car (response variable), using data on acceleration rate of the car, the fuel consumption and the year that it was built. The conceptual model for this research looks like this:

### 2.3.2 Visual inspection

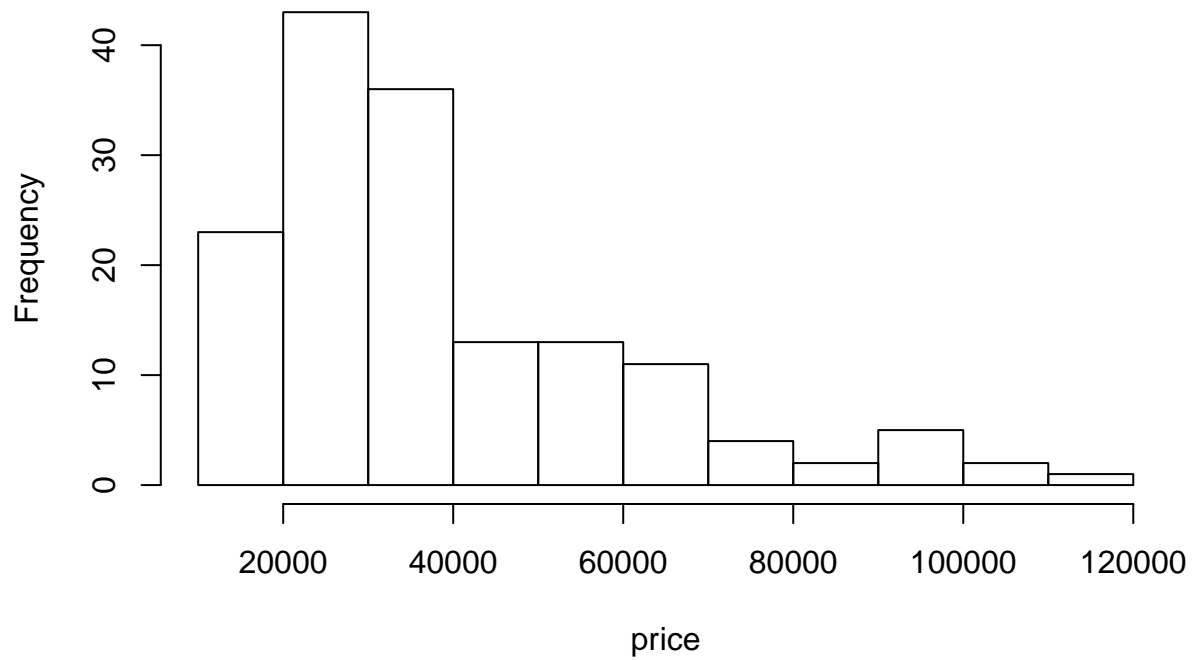
The distribution of the independent variable is displayed

```

# Reading in the necessary packages
library(readr)
d <- read_csv("hybrid_reg.csv")
mpg <- d$mpg
year <- d$year
accel <- d$accelrate
price <- d$msrp
  
```

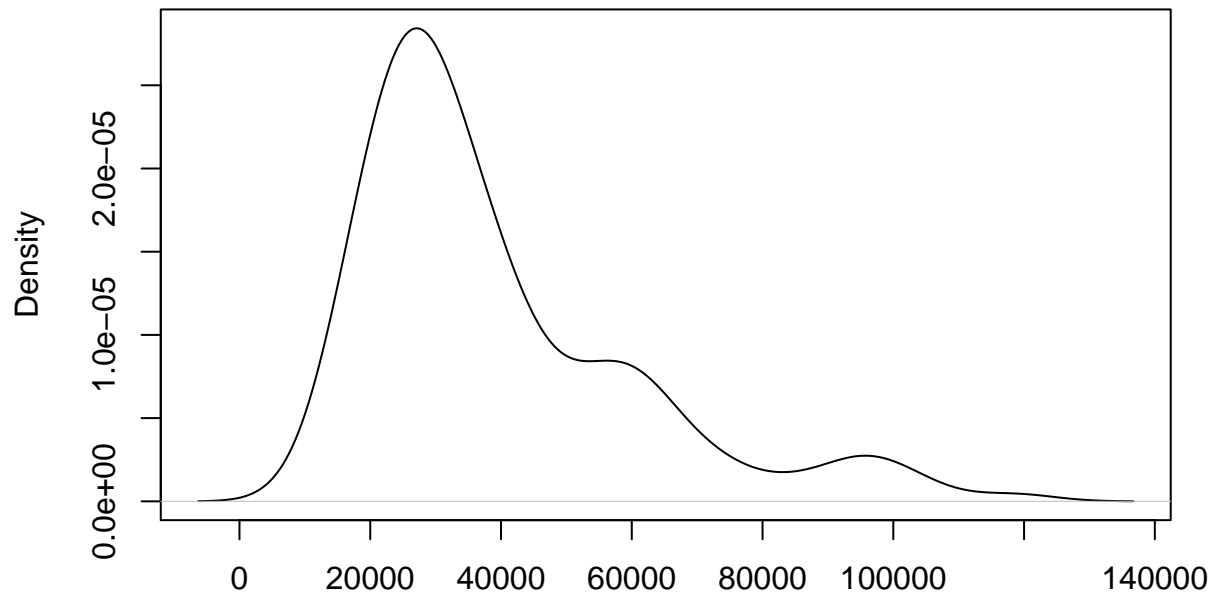
```
# Histogram of the distribution of the price variable  
hist(price)
```

**Histogram of price**



```
# Density plot of the price variable  
plot(density(price),main="Density plot of price")
```

## Density plot of price

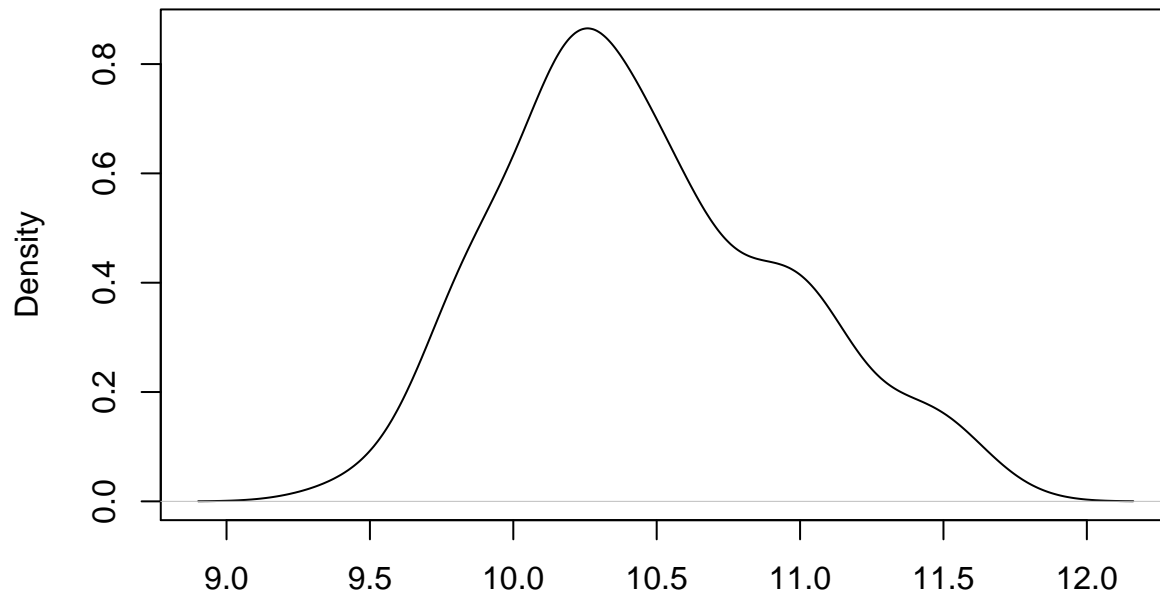


N = 153 Bandwidth = 6055

Visual inspection of the plots reveals that the distribution of price deviates from a normal distribution. Especially the right tail of the density distribution has more mass than it should have. Since the distribution is right skewed, a logarithmic transformation is effective in increasing the normality; the result can be seen in the figure below.

```
# Density plot of the transformed price variable
plot(density(log(price)),main="Density plot of log(price)")
```

### Density plot of log(price)



N = 153 Bandwidth = 0.1596

```
shapiro.test(price)
```

```
##  
## Shapiro-Wilk normality test  
##  
## data: price  
## W = 0.85261, p-value = 4.345e-11
```

```
shapiro.test(log(price))
```

```
##  
## Shapiro-Wilk normality test  
##  
## data: log(price)  
## W = 0.97322, p-value = 0.004424
```

As we will see later on, the logarithmic transformation is also necessary to justify the choice of linear regression, as without it the assumptions for linear regression do not hold.

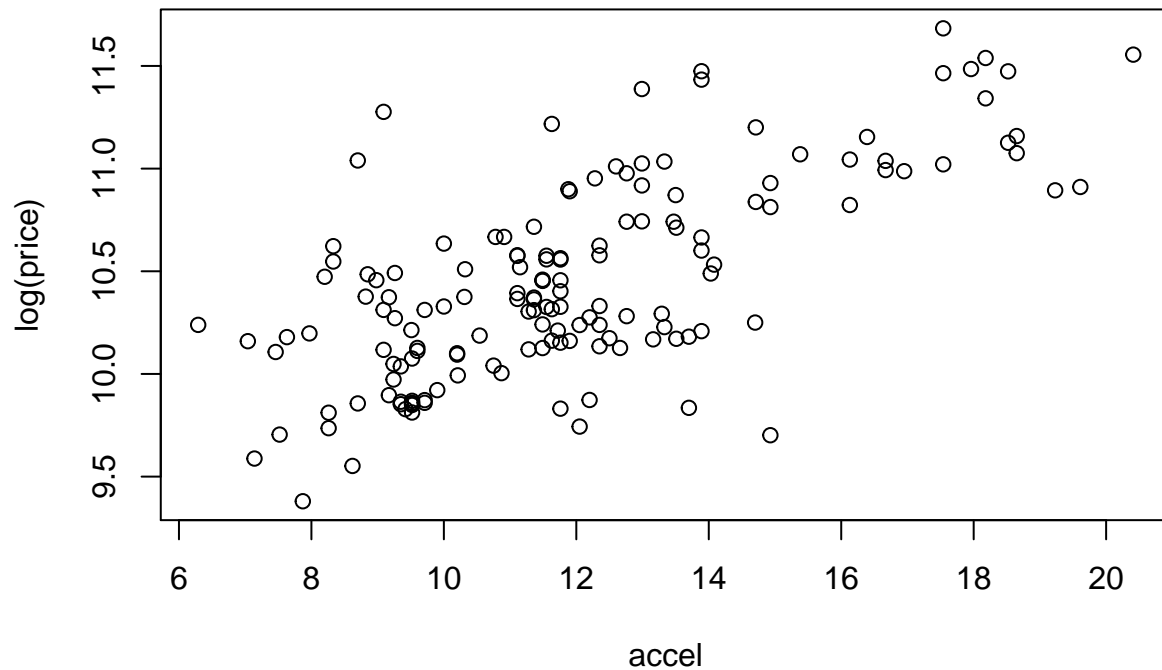
#### 2.3.3 Scatter plot

Using scatter plots we can visually examine the relationship between two variables. The following figures show the scatter plots of the response variable price paired with each of the predictors.

```
plot(log(price) ~ accel, main="log(price) vs accel")
```

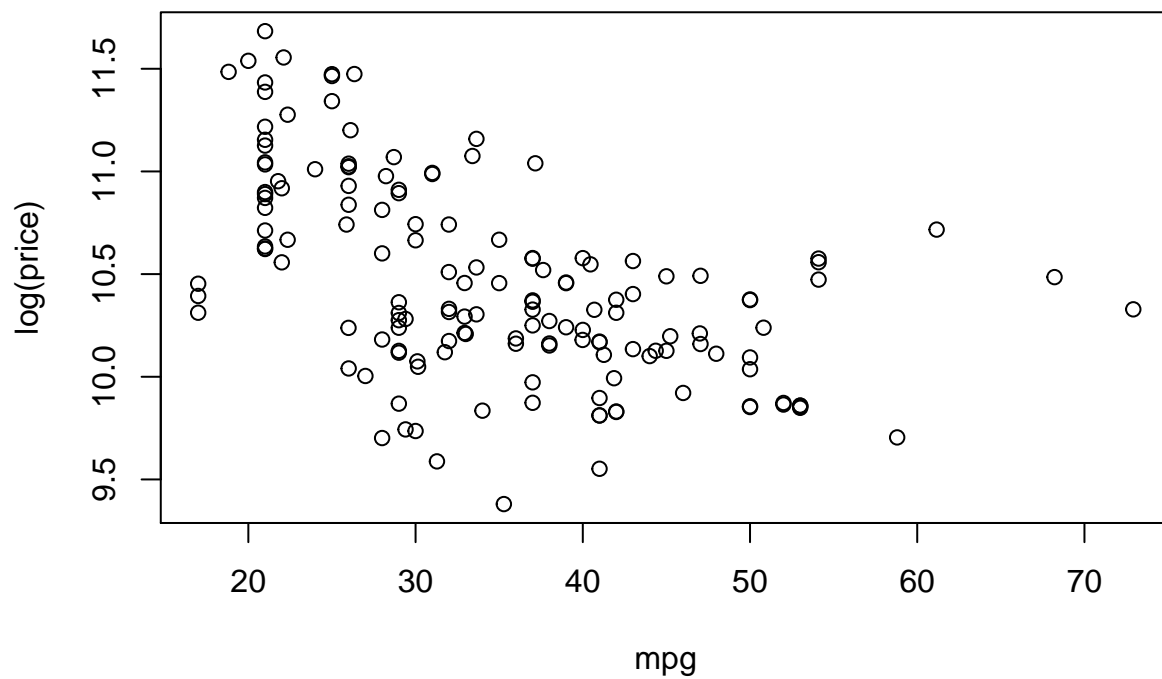


**log(price) vs accel**

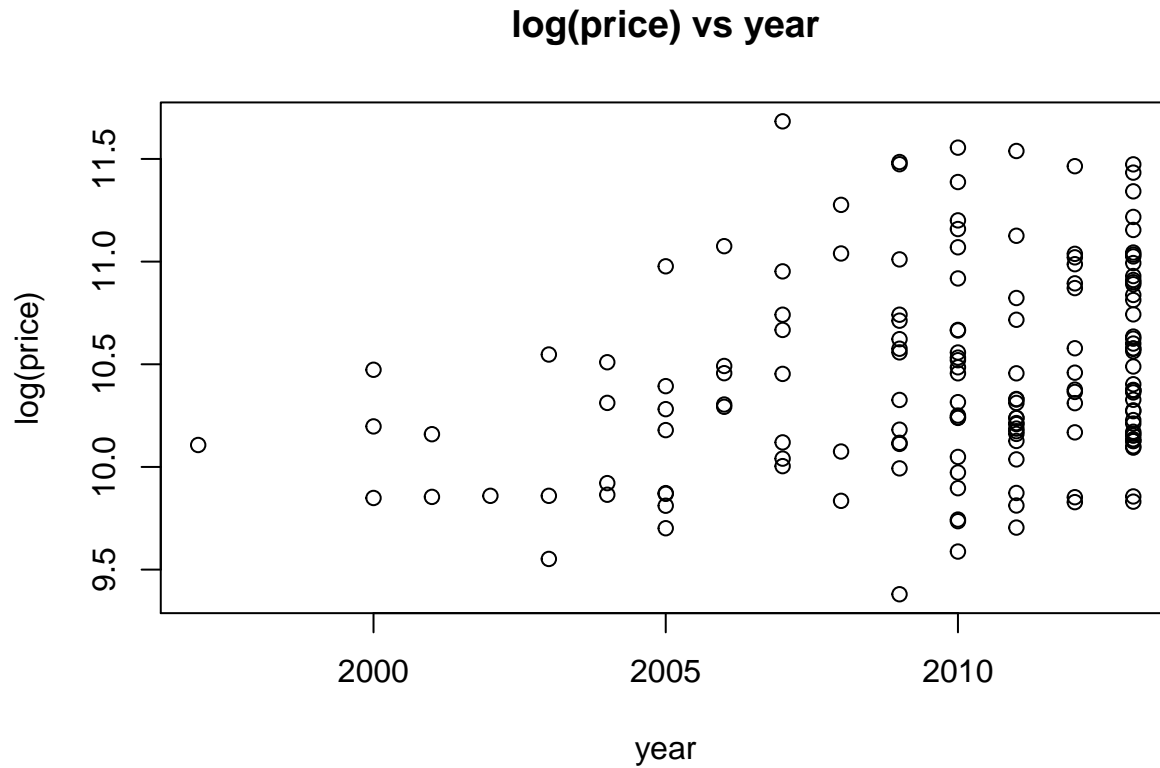


```
plot(log(price) ~ mpg, main="log(price) vs mpg")
```

**log(price) vs mpg**



```
plot(log(price) ~ year, main="log(price) vs year")
```



#### 2.3.4 Linear regression

From the anova table we can see that the accel and mpg variables are able to significantly improve the model. The year variable is not able to explain any additional significant variance in the price variable. Therefore we exclude the price variable from further analysis.

```
library(pander)
model0 <- lm(log(price) ~ 1)
model1 <- lm(log(price) ~ accel)
model2 <- lm(log(price) ~ accel + mpg)
model3 <- lm(log(price) ~ accel + mpg + year)

pander(anova(model0,model1,model2,model3),
       caption = "Model comparison to predict the price of a car")
```

Table 6: Model comparison to predict the price of a car

Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
152	35.77	NA	NA	NA	NA
151	18.83	1	16.94	146.9	5.776e-24
150	17.19	1	1.639	14.22	0.0002343
149	17.18	1	0.01013	0.08785	0.7673

```
library(QuantPsyc)
pander(confint(model2),
       caption = "#95% confidence interval of the estimates")
```

Table 7: #95% confidence interval of the estimates

	2.5 %	97.5 %
(Intercept)	9.328	10.13
accel	0.07142	0.1142
mpg	-0.0167	-0.00524

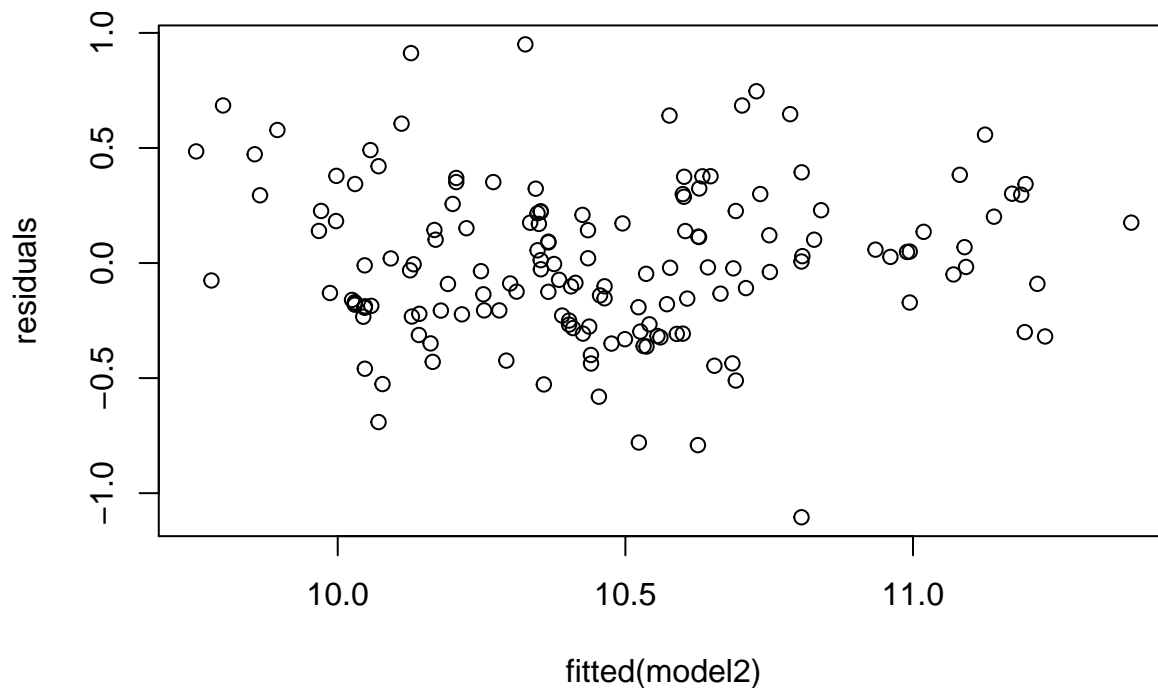
```
pander(lm.beta(model2),
       caption = "standardised regression coefficients") # standardised regression coefficients
```

accel	mpg
0.5626	-0.2482

### 2.3.5 Examine assumption

The residuals vs fitted plot is a useful tool in examining the linearity and equal variances assumptions.

```
residuals = resid(model2)
plot(residuals ~ fitted(model2))
```

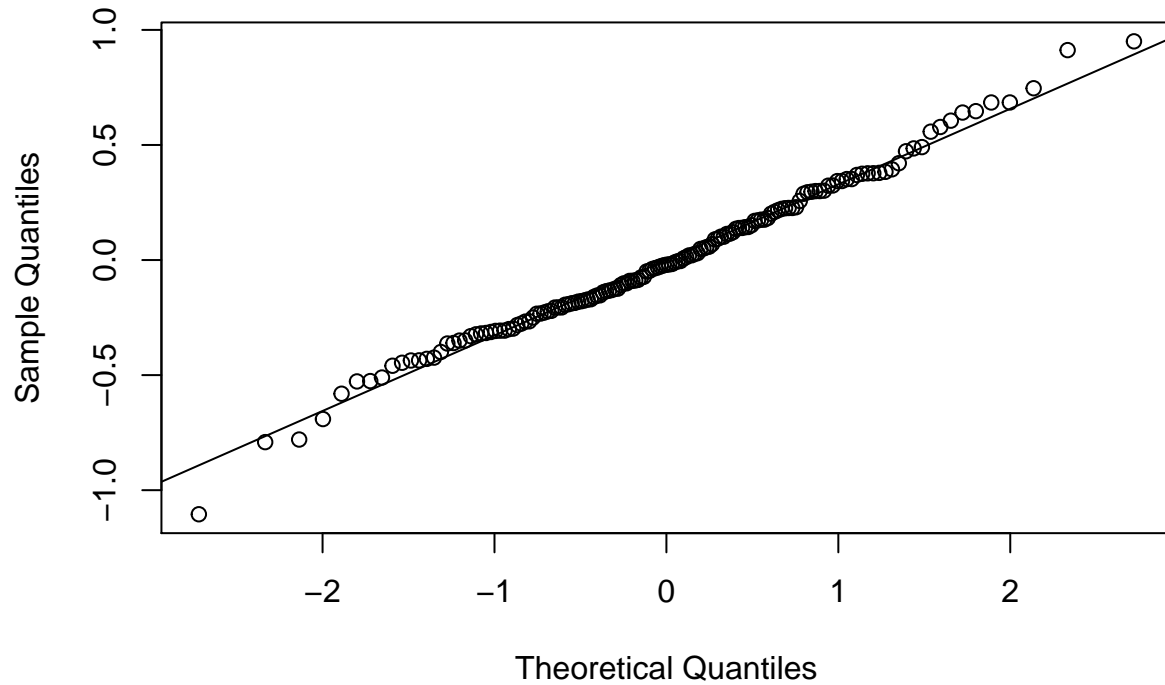


normality we can check the qq-plot:

```
library(car)
qqnorm(residuals)
qqline(residuals)
```

For

## Normal Q-Q Plot



Multi-

collinearity:

```
vif(model12)
```

```
##      accel      mpg  
## 1.34428 1.34428
```

```
1/vif(model12) # Tolerance
```

```
##      accel      mpg  
## 0.7438928 0.7438928
```

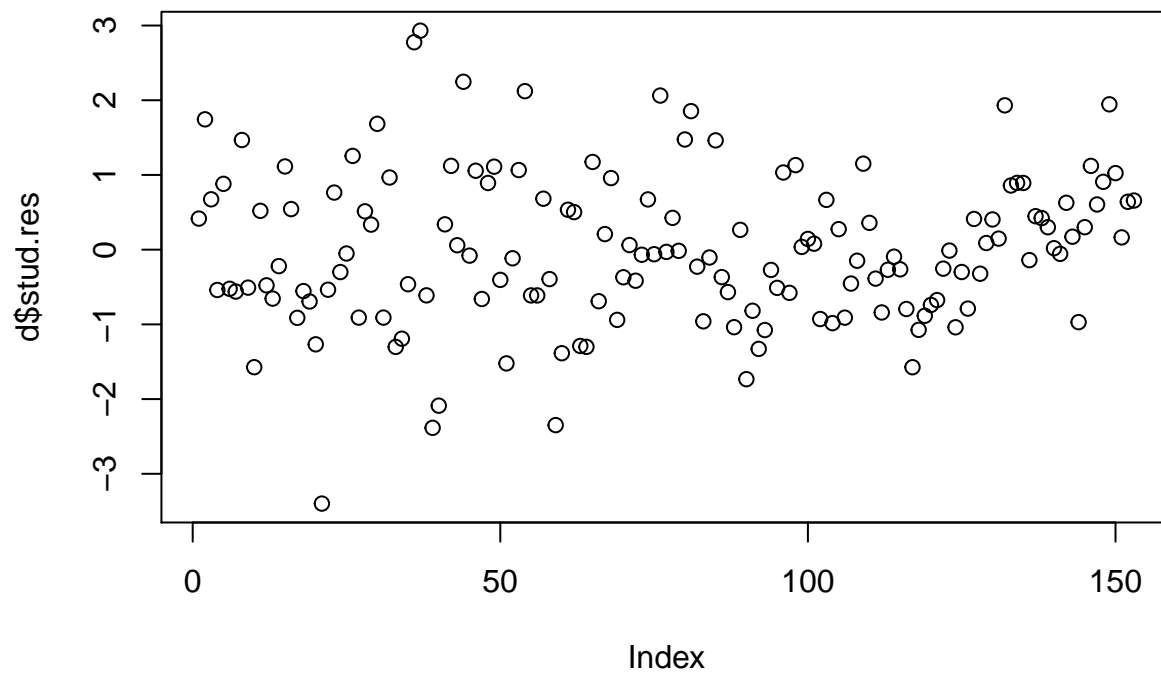
Autocorrelation:

```
durbinWatsonTest(model12)
```

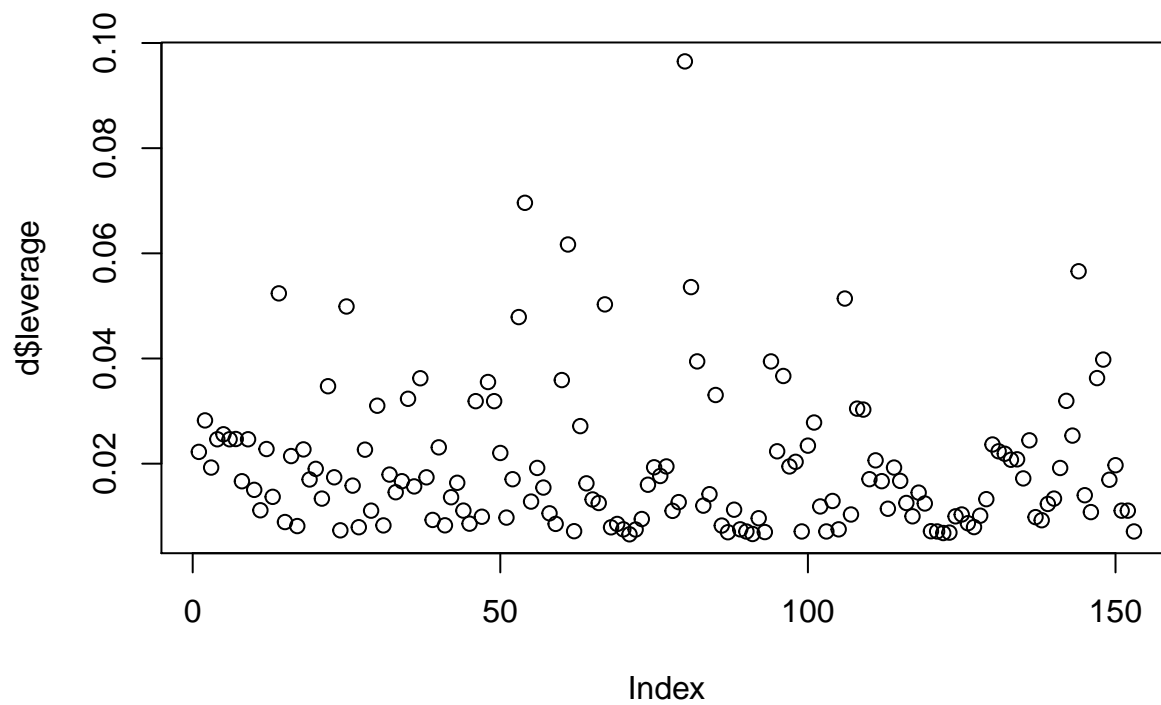
```
## lag Autocorrelation D-W Statistic p-value  
## 1      0.2303107      1.535374    0.004  
## Alternative hypothesis: rho != 0
```

### 2.3.6 Impact analysis of individual cases

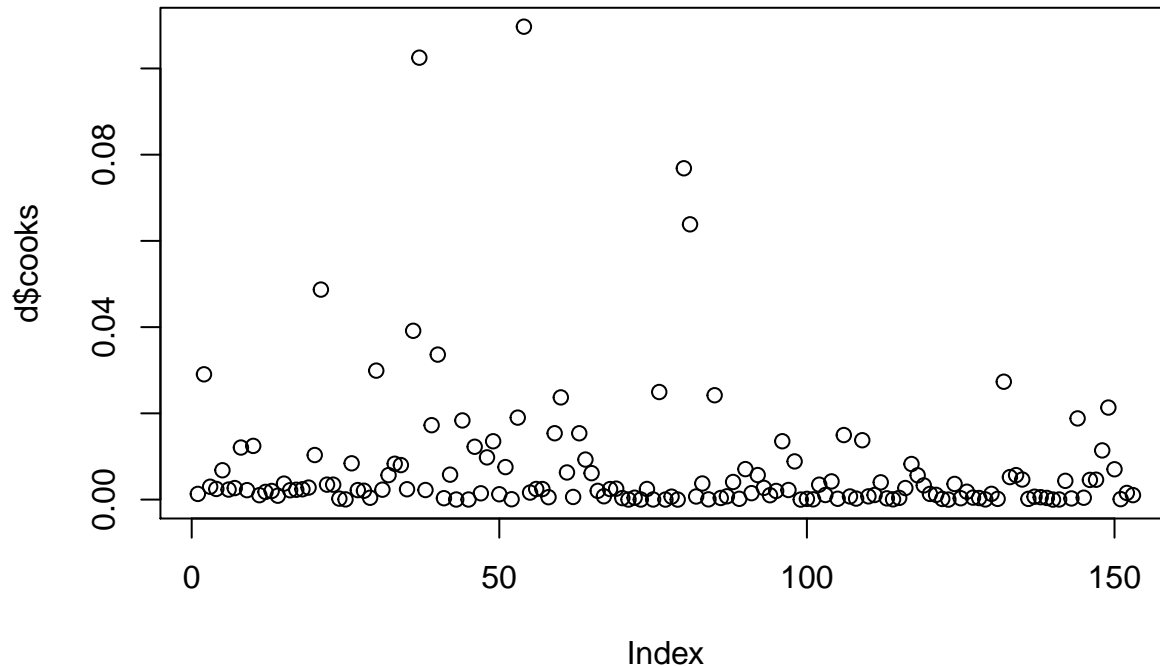
```
d$stud.res<-rstudent(model12)  
plot(d$stud.res)
```



```
d$leverage<-hatvalues(model2)  
plot(d$leverage)
```



```
d$cooks<-cooks.distance(model2)  
plot(d$cooks)
```



### 2.3.7 Report section for a scientific publication

In this section we briefly present the results of fitting a linear regression model in order to predict the price of cars using data on their acceleration rate, their fuel efficiency in miles per gallon, and their build year.

First of all we can conclude that a logarithmic transformation was necessary to increase normality of the distribution of the price variable. While the distribution still significantly differs from a normal one ( $W = 0.973$ ,  $p = 0.004$ ), it is an improvement over the original distribution ( $W = 0.853$ ,  $p = 4.3e-11$ ). Furthermore, the transformation was necessary to justify the choice of performing linear regression, as without it the assumptions for linear regression do not hold, especially the linearity assumption.

Inspecting the scatter plots, it becomes clear that a linear relationship between the natural logarithm of price and the year the car was built is absent. The scatter plots of the other two variables show some indication that a relationship might exist.

Fitting the model revealed that the year variable is indeed not able to explain any additional variance in the price variable on top of the accel and mpg variables ( $F = 0.088$ ,  $p = 0.77$ ). Based on this result we decided to exclude the independent variable year from the model. Hence we end up with the following model:

$$\log(\text{price}) = 9.73 + 0.093 \times \text{accel} - 0.011 \times \text{mpg}$$

Checking the assumptions of the linear regression, we found that the distribution of the residuals is normal with expected value 0 and (roughly) constant variance. Testing for independence showed a violation of the assumption ( $D-W = 1.54$ ,  $p = 0.008$ ). Violating the independence assumption is quite problematic, but for the sake of the exercise we will continue the analysis. Additionally, no multicollinearity could be found in our model.

Analysis of influential and leverage points revealed no severe outlying cases that undermine the linear regression model.

The interpretation of the coefficients is slightly tricky, since we are dealing with a transformed dependent variable. Instead of additive, the model becomes multiplicative, and each coefficients has to be interpreted as an exponent (i.e. the intercept becomes  $e^{9.73} = 16,815$ ). The standardized coefficients tell us that the effect of one higher standard deviation in the acceleration rate has about twice the effect on the price of one higher standard deviation in the fuel efficiency.

To conclude, we were able to formulate a linear model that is to some extent able to predict the price of a car based on its acceleration rate and its fuel efficiency. Since not all assumptions of linear regression were met, interpretation of the results requires caution.

## 2.4 Question 4 - Logistic regression analysis

```
#include your code and output in the document
Data <- read.csv("port_taiwan.csv",header=TRUE)

Data$year <- factor(Data$year, levels=c(2003:2006), labels=c("year2003","year2004","year2005","year2006"))

#remove one port so out data can pass as dichotomous
PortData <- subset(Data,(port != "3"))
PortData$port <- factor(PortData$port, levels=c(1:2),labels=c("1","2"))
```

### 2.4.1 Conceptual model

Make a conceptual model underlying this research question

### 2.4.2 Logistic regression

Conduct a logistic regression, examine whether adding individual indicators in the model improves the model compared to Null model. Make a final model with only significant predictor(s). For this model, calculate the pseudo R-square. Calculate the odd ratio for the predictors and their confidence interval

```
#include your code and output in the document
```

### 2.4.3 Crosstable predicted and observed responses

Make a crosstable of the predicted and observed response

```
#include your code and output in the document
```

### 2.4.4 Report section for a scientific publication

Write a small section for a scientific publication, in which you report the results of the analyses, and explain the conclusions that can be drawn.

## 3 Part 3 - Multilevel model

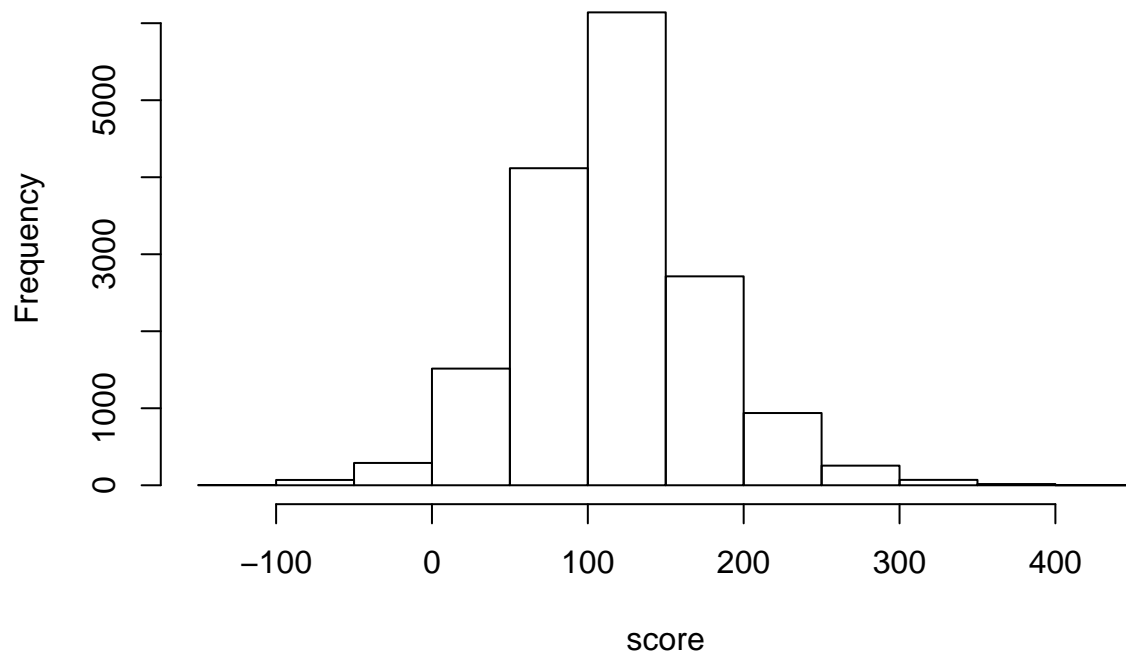
### 3.1 Visual inspection

First, we inspect the distribution of the scores.

```
set1 <- read.csv("set1.csv", header = TRUE)
set1$Subject <- factor(set1$Subject)
# Should we consider the session as factor ?
#set1$session <- factor(set1$session)

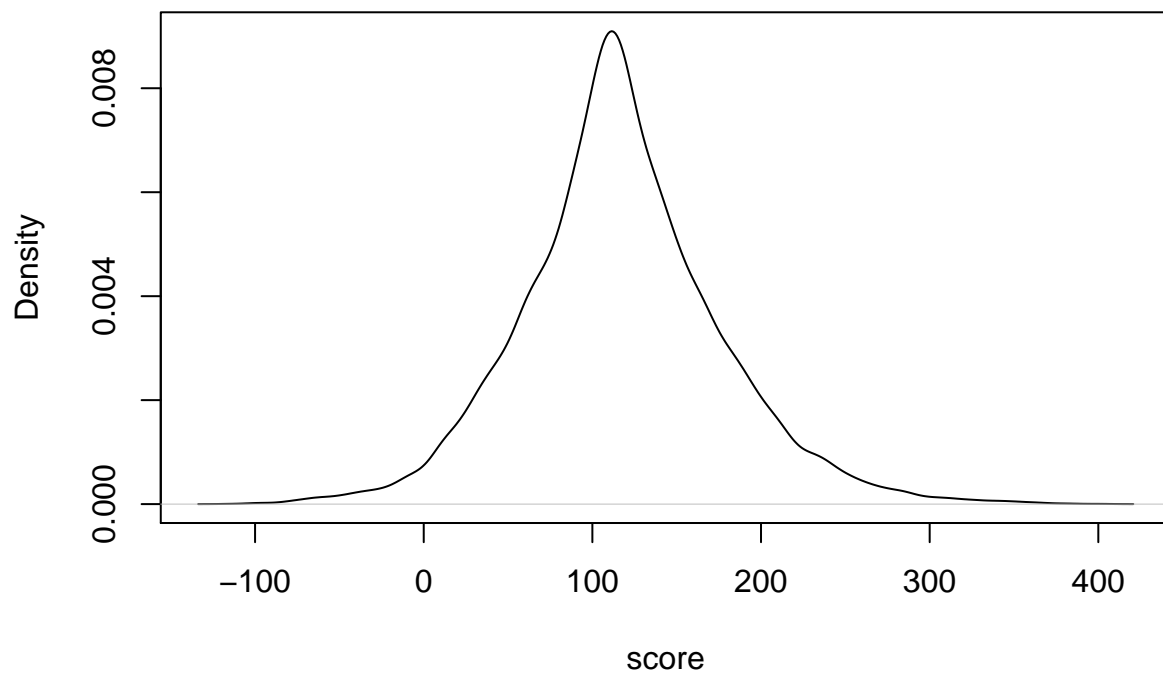
hist(set1$score, xlab="score", main="Histogram of the scores")
```

### Histogram of the scores



```
plot(density(set1$score), xlab="score", main="density of the scores")
```

### density of the scores

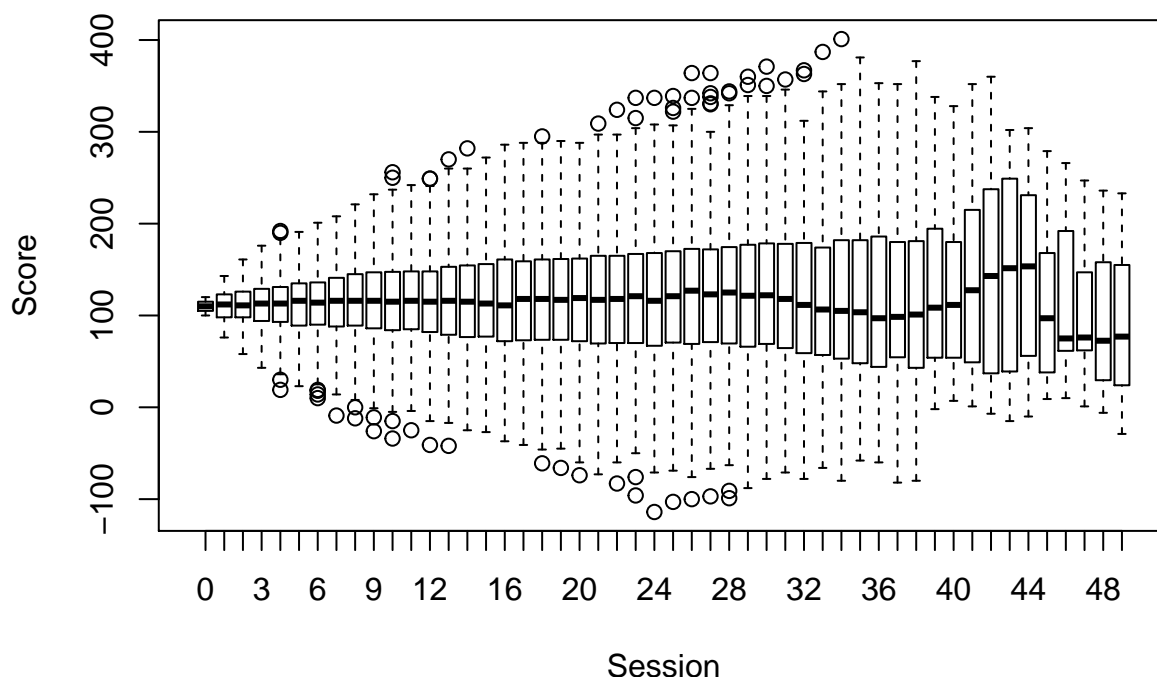


If we just look at the score distribution, they seems to be normally distributed and centered around 150. We will now look at how the session affect the scores with an assumption of iid of the variables (which is false).



```
#Plot of the relationship between session and score
# Assuming iid of the variables which is not true, shouldn't do that.
boxplot(score ~ session, data=set1, xlab="Session", ylab="Score", main="relationship between score and session")
```

## relationship between score and session if we assume iid



It seems that even without looking at each person individually the session seems to have an impact on the scores. We will further analyse this effect by taking into account the evolution of each person now.

### 3.2 Multilevel analysis

Conduct multilevel analysis and calculate 95% confidence intervals, determine:

We will conduct a multilevel analysis of the dataset where each subject get a random intercept. We first want to know if the session has an impact on people scores.

```
randomIntercept <- lme(score ~ 1, data = set1, random = ~1|Subject, method="ML")
addSession <- update(randomIntercept, .~. + session)
pander(anova(randomIntercept,addSession), caption="comparisons of models when session is added as a fixed factor")
```

Table 9: comparisons of models when session is added as a fixed factor (continued below)

	call	Model	df	AIC
<b>randomIntercept</b>	lme.formula(fixed = score ~ 1, data = set1, random = ~1   Subject, method = "ML")	1	3	162711
<b>addSession</b>	lme.formula(fixed = score ~ session, data = set1, random = ~1   Subject, method = "ML")	2	4	162545

	BIC	logLik	Test	L.Ratio	p-value
<b>randomIntercept</b>	162734	-81352		NA	NA
<b>addSession</b>	162576	-81269	1 vs 2	167.7	2.317e-38

The addition of the session significantly improves the model ( $p. < 0.001$ ), we will now verify the 95% confidence bound.

```
intervals(addSession)
```

```
## Approximate 95% confidence intervals
##
## Fixed effects:
##           lower      est.      upper
## (Intercept) 106.8665678 111.0675622 115.2685566
## session      0.3126229  0.3682005  0.4237781
## attr("label")
## [1] "Fixed effects:"
##
## Random Effects:
## Level: Subject
##           lower      est.      upper
## sd((Intercept)) 43.67471 46.5146 49.53914
##
## Within-group standard error:
## lower      est.      upper
## 34.68269 35.06933 35.46028
```

We see that for the fixed effect, the session deviates from 0 in the 95% interval which means that there is a significant impact of the session on people's scores.

We will now look if there is a significant variance between the participants in their score.

```
intervals(randomIntercept)
```

```
## Approximate 95% confidence intervals
##
## Fixed effects:
##           lower      est.      upper
## (Intercept) 112.7019 116.8139 120.9259
## attr("label")
## [1] "Fixed effects:"
##
## Random Effects:
## Level: Subject
##           lower      est.      upper
## sd((Intercept)) 43.68633 46.52747 49.55338
##
## Within-group standard error:
## lower      est.      upper
## 34.86891 35.25763 35.65067
```

We can see that in the Random effect, the standard deviation of the intercept does not include 0 in the 95% interval, thus there is a significant variance between the participants in their score.

### 3.3 Report section for a scientific publication

We fitted a linear Mixed-Effects model with a random intercept for each subject and we then build a new model that used the session as independent variables to predict the scores. It appears that there is a significant main effect of the session over participants score ( $\chi^2(1) = 167.7$ ,  $p. < 0.001$ ) when compared to the baseline at a 95% confidence interval ( $\text{estimate}(\text{session}) = [0.31; 0.42]$ ). We can also show that there is a significant variance between participants in their score ( $\text{sd}(\text{intercept}) = [43.69; 49.55]$ ) at a 95% confidence interval. From these results we can draw the conclusion that the session in which the participants is has an impact on his score which can be interpreted as improvement over each exercise session since the estimates are positive. We can also conclude that each participant is different and that indeed, we cannot consider the observations to be independent.