

# STK4900 Oblig1

Daniel Heinesen, daniehei

March 13, 2019

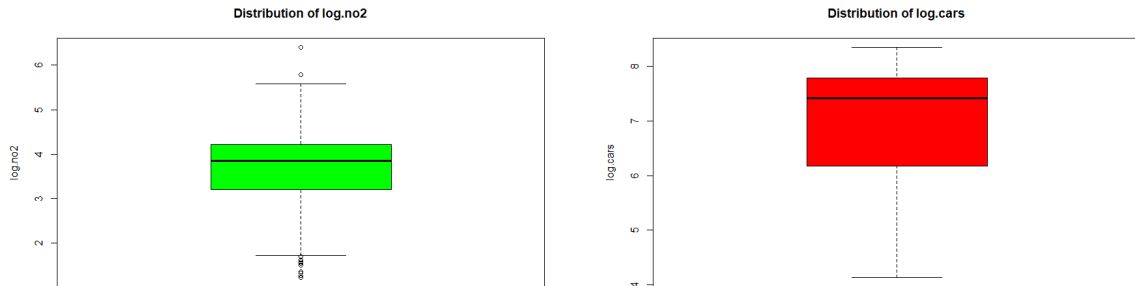
## Exercise 1

a)

log.no2	log.cars
Min. :1.224	Min. :4.127
1st Qu.:3.214	1st Qu.:6.176
Median :3.848	Median :7.425
Mean :3.698	Mean :6.973
3rd Qu.:4.217	3rd Qu.:7.793
Max. :6.395	Max. :8.349

Table 1: Summary of the main features of *log.no2* and *log.cars*

We are looking at the variables *log.no2* and *log.cars*, which represents the logarithm of the measured concentration of NO<sub>2</sub> and the logarithmic number of cars per hour. Table 1 shows a numerical summary of *log.no2* and *log.cars*. Looking at *log.no2* we see that the mean is close to the median, and that they both are more or less equidistant from *min* and *max*, and from the first and third quartiles. This indicates that the distribution is almost symmetric, with a slight skew. Looking at *log.car* we see that this is no longer the case. The mean and median is more different and lay closer to *max* and the third quartile. This means that this distribution is very skewed.



(a) Box plot of the distribution of *log.no2*

(b) Box plot of the distribution of *log.cars*

Figure 1: Box plots showing the distributions for *log.no2* and *log.cars*.

Our ideas about the distribution we got from the numerical summaries, tab. 1, is confirmed by the box plots of the distribution found in fig. 1. Here we can see that *log.no2* is more or less symmetric, while *log.cars* is heavily skewed towards higher values.

We can look at the relationship between *log.no2* and *log.cars*. If we look at fig. 2 we see that there looks to be a linear dependence. We are going to try to fit a linear model to the data, but looking at the plot we see that there seems to be a higher clustering for higher values, which we expect from fig. 1. This means that there seems to be a larger spread for lower values than for the higher, making some of the assumptions used for a linear regression false. We will look at that later.

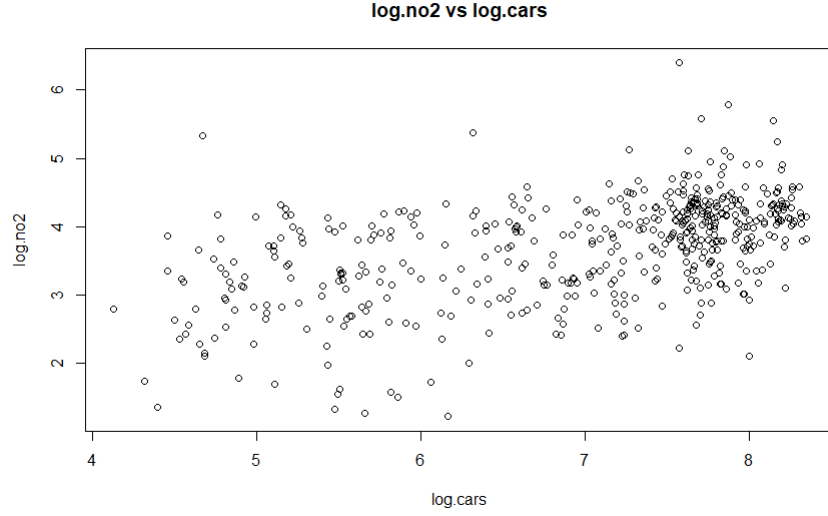


Figure 2: Scatter plot showing the relation between  $\log.no2$  and  $\log.cars$ .

b)

We are not going to fit a linear model to our  $\log.no2$  and  $\log.cars$ .

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	1.2331	0.1875	6.57	0.0000
$\log.cars$	0.3535	0.0266	13.30	0.0000
Residual standard error: 0.6454 on 498 degrees of freedom				
Multiple R-squared: 0.2622, Adjusted R-squared: 0.2607				
F-statistic: 177 on 1 and 498 DF, p-value: < 2.2e-16				

Table 2: Summary of the linear model for  $\log.no2$  and  $\log.cars$ .

We start by looking at fig. 3. We can see that the line representing the linear model seems to be a good fit, but the residuals are quite large. If we look at the numerical summary in tab. 2 we see that we get the model

$$\log.no2 = 1.2331 + 0.3535 \cdot \log.cars + \epsilon. \quad (1)$$

This means that given no cars on the road we expect that  $\log.no2 = 1.2331$ , and with each additional (logarithmic) car per hour,  $\log.no2$  increase by 0.3535. Looking at the p-value of  $p \approx 0$  we can conclude that the linear fit is significant. But if we look at the  $R^2 = 0.2622$  we see that this value is rather low. This can indicate that the data doesn't have a linear relationship – which our p-value indicates that our data have – or that  $\log.no2$  isn't described by  $\log.cars$  alone.

c)

As mentioned above, there seems to be a bigger spread in the data in fig. 3 for lower values than for higher ones. This is one of the things we need to look at now.

For a linear model

$$y_i = \alpha + \beta_i x_i + \epsilon_i, \quad (2)$$

where  $\epsilon_i$  is the residuals, we need to have some assumptions.

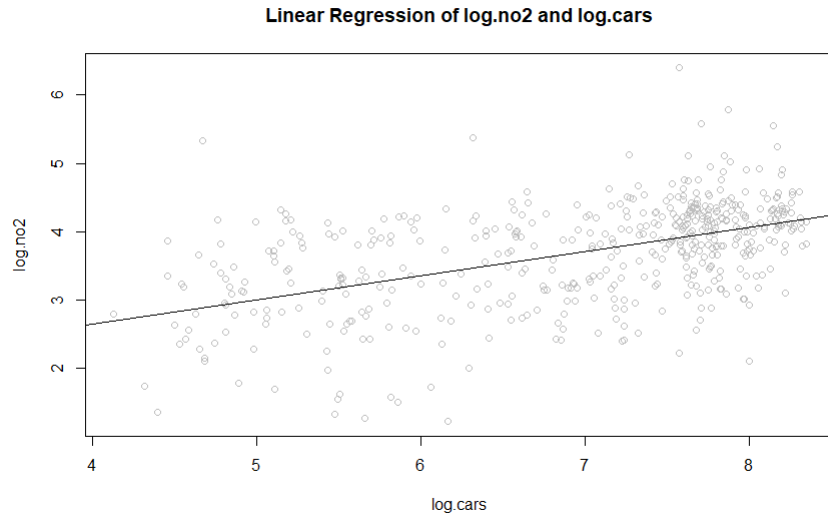


Figure 3: The linear model superimposed on the scatter plot from fig. 2

### Linearity

For a linear regression one of the most important assumptions is that the model is linear in the predictor. We are going to check this with a component-plus-residual plot, CPR plot. This is where we plot the partial residuals  $\beta_i x_i + r_i$  against the predictor  $x_i$ . Since we have a model with only one predictor, this will look like our scatter plot. But R also does a smoothing of the data so that we can look at how it compares to a straight line.

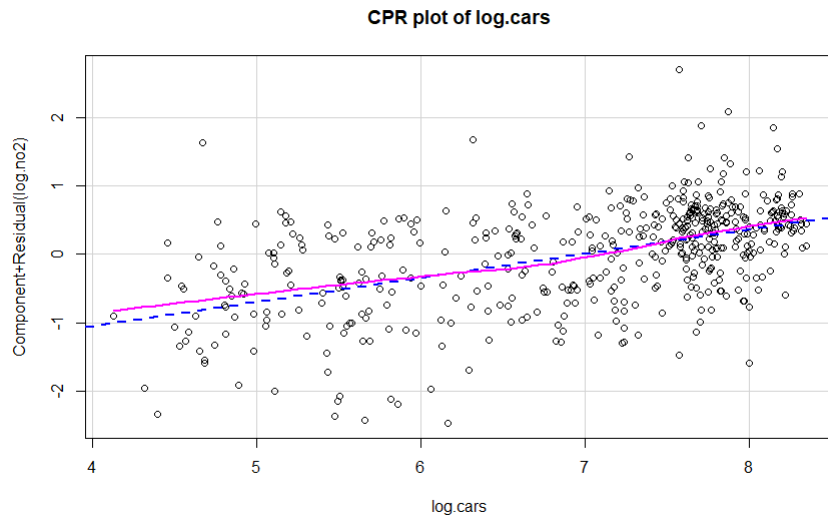


Figure 4: Component-plus-residual plot of *log.cars*.

From fig. 4 we see that the data is somewhat linear. There is some deviation from linearity for the lower values, but the deviation is not so great that we say that the model is non-linear.

## Homoscedasticity

For our model we want our residuals to be normally distributed with the same distribution independent of the value of *log.cars*. To check this we are going to plot the standardized residuals – from a R function – against the fitted values. We expect that the residuals there are random around a mean for all the fitted values, meaning that we can fit a straight, horizontal line to the data.

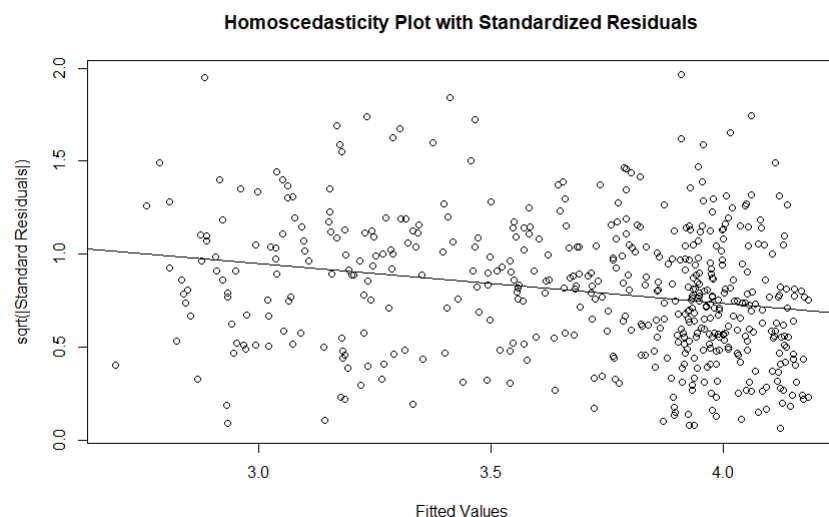


Figure 5: Figure showing  $\sqrt{|\text{Standardized Residuals}|}$  against the fitted values.

If we look at fig. 5 we see that the fitted line is not horizontal, but instead decreasing. This indicates that the residuals are smaller for larger values of *log.cars*. This seems to verify what we expected from the start: That there was a larger spread for the smaller values of *log.cars* than for the larger.

## Normality

The last assumption we have is that the residuals are normal.

Looking at fig. 6 we see that the distribution looks more or less normal, but the median is somewhat off center, which indicates skewness. Fig. 7 shows that this is the case. The distribution is far from normal, and seems to be skew towards the negative values. From both these figures we can see that neither the median nor the mean is 0, which is something we would like it to be.

The last nail in the coffin is fig. 8. If the residuals were normally distributed, we would expect the points to follow the straight line. The S-shape indicates a light tail, and the bend we see in the middle indicates a left skew – which corresponds well with our other plots –. This means that our assumptions that the residuals are normally distributed is false.

So we have seen that while our model is close to linear, the residuals are neither homoscedastic nor normal. So to use a linear model for our data is difficult to justify.

## d)

Our data consists of more possible predictor variables than just *log.cars*. We have three more variables we can use in a multiple regression. To see which if these variables we are going to use, we are going to do a *forward selection* of the variables and their logarithm: We start by making a linear model with each predictor variable (and their logarithm), using **Cross Validated  $R^2$**  we will find the best model. We then go through the remaining variables, adding that to the



Figure 6: Figure a box plot of the residuals.

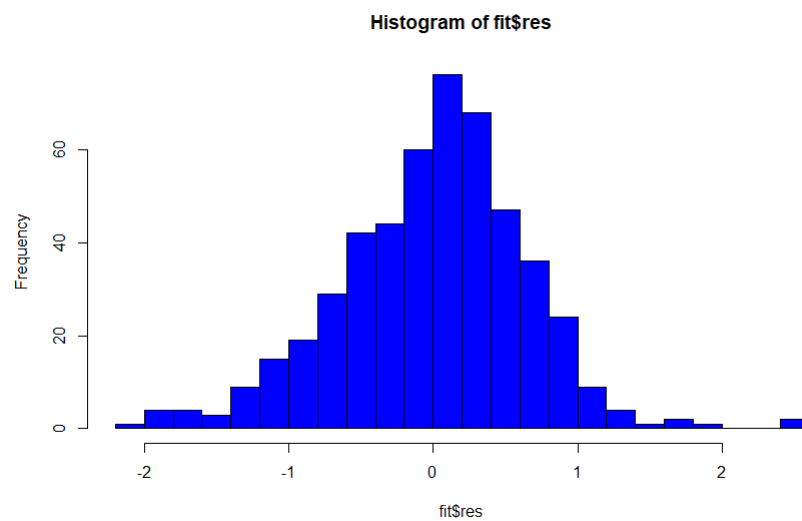


Figure 7: Figure the histogram of the residuals.

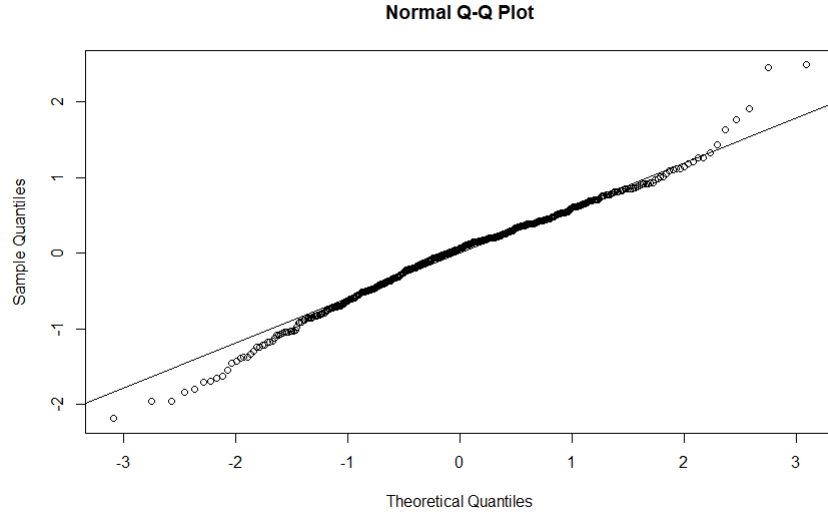


Figure 8: Figure a QQ-plot of the residuals.

best predictor and see which variable improves the fit. This is done recursively until there is no improvement in the cross validated  $R^2$ . We then have our best model.

Due to my lack in familiarity with R, this was done in the stupid way of check each combination one by one.

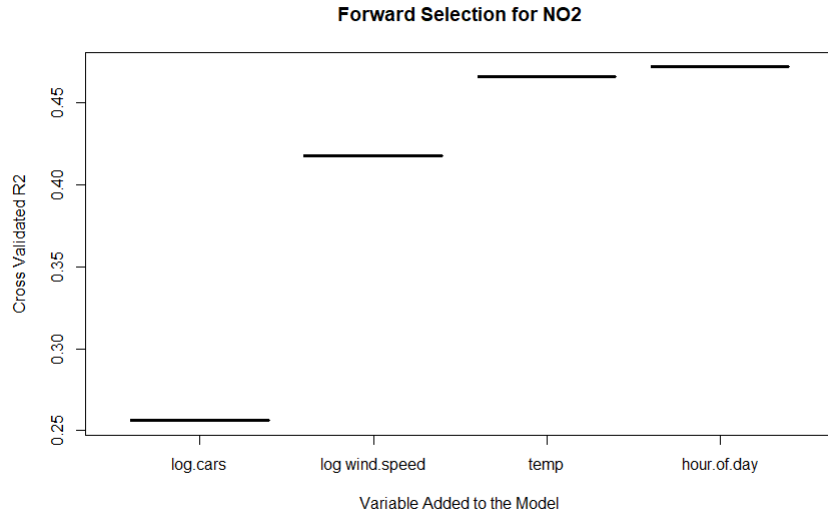


Figure 9: Figure of the increase in the cross validated  $R^2$  as variable are added. The x-axis shows which predictor variable is added to the model.

We see from fig. 9 that as we add more predictor variables the model becomes better and better. We also notice that to get the best model we also had to take the logarithm of *wind.speed*. The model we are left with then is found in tab. 3, and can be written as

$$\log.no2 = 1.07 + 0.46 \cdot \log.cars - 0.42 \cdot \log(wind.speed) - 0.03 \cdot temp - 0.01 \cdot hour.of.day + \epsilon \quad (3)$$

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	1.0709	0.1710	6.26	0.0000
log.cars	0.4572	0.0279	16.37	0.0000
log(wind.speed)	-0.4194	0.0364	-11.53	0.0000
temp	-0.0267	0.0038	-6.96	0.0000
hour.of.day	-0.0123	0.0044	-2.81	0.0051

Residual standard error: 0.542 on 495 degrees of freedom  
Multiple R-squared: 0.483, Adjusted R-squared: 0.479  
F-statistic: 116 on 4 and 495 DF, p-value: <2e-16

Table 3: Table showing the final model for our multiple regression. This is the model which gives the best cross validation  $R^2$ .

e)

### Checking our Assumptions

We will now check if our model assumptions are correct. We will do the exact same analysis as in exercise 1c. We will do it a bit quicker now since we have done it once before.

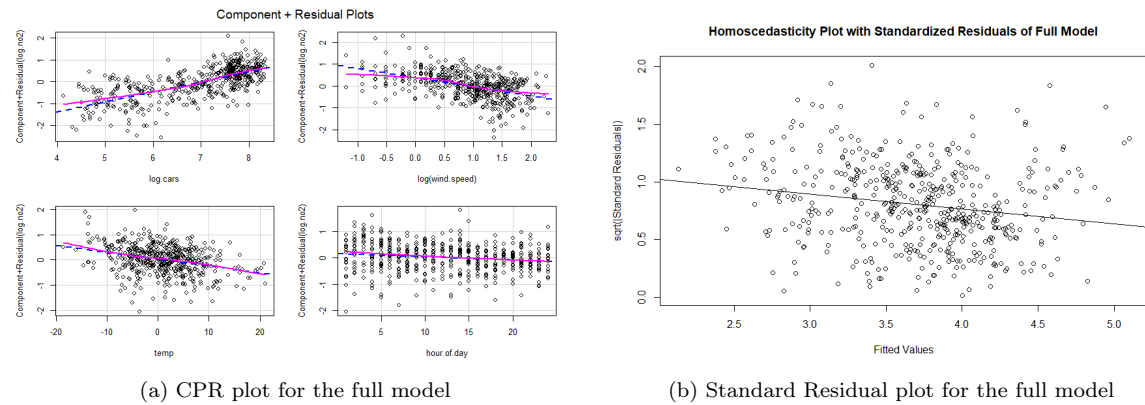


Figure 10: Plots judging the linearity and homoscedasticity of the residuals of the full model.

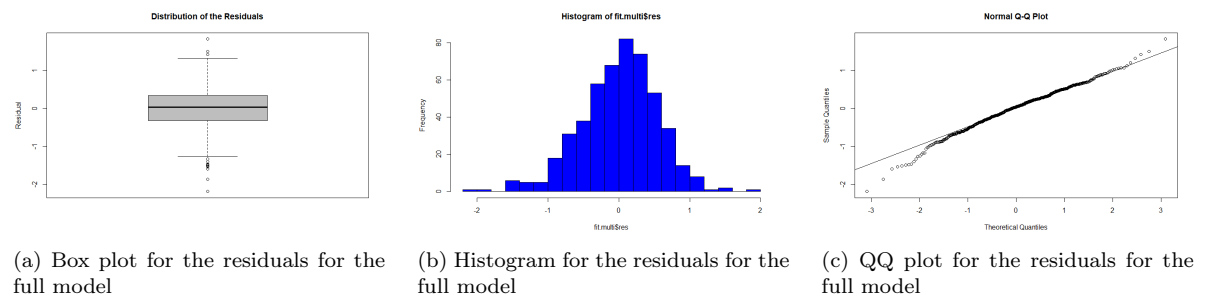


Figure 11: Plots judging the normality of the residuals of the full model.

Fig. 10a shows that *temp* and *hour.of.day* seems to be more or less linear, while *logwind.speed* has some small deviations from linearity. It is close enough that we can say that our model is linear.

Fig. 10b shows that we still have the problem that the residuals are decreasing, which is problematic for our assumption of homoscedasticity.

From the plots in fig. 11 we see that the residuals are a bit non-normal still. If we only look at the values close to the mean we see that the QQ-plot indicates these to be normal. It is the outliers – as judged by the box plot – that seems to be breaking normality. To excluding these, we can say that the residuals are normally distributed for the full model.

### Interpreting our Model

If we look at (3) and tab. 3 we can see that the base value of  $\log.no2$  is 1.07, and for each (logarithmic) car per hour this value is increased with 0.46, and with each extra (logarithm)  $m/s$  of wind speed it is decreased with 0.42. With  $1^\circ C$  increase in temperature the  $\log.no2$  is decreased with 0.03. And lastly with one later in the day, it is decreased with 0.01.

We see that we get a  $R^2$  value of 0.483, which is significantly higher than with our model in tab. 1. It is still not close to 1, so there is some other factors with helps determine  $\log.no2$ .

So all in all, we have a good model, which fullfill most of our assumptions of a linear model.

## Exercise 2

a)

The data we are looking at is the systolic blood pressure for 36 people divided into three groups determined by their age. These groups are 30 – 45 years, 46 – 59 years and 60 – 75 years.

Age Group	N	Mean	St. Dev.	Min	Pctl(25)	Pctl(75)	Max
All	36	138.806	25.749	104	117.5	156.2	214
30 – 45 years	12	122.167	15.338	104	112	129	160
46 – 59 years	12	139.083	22.625	108	121.5	157.8	174
60 – 75 years	12	155.167	27.719	110	138	164	214

Table 4: A table showing different measures for the distribution of the systolic blood pressure for the three age groups, and the entire group as a whole. The increase in all the measures for the higher age groups seems to indicate a correlation between age and blood pressure.

If we look at table 4 we see the numerical summary of the blood pressure of the three age groups, and the group as a whole. From the mean column, we see that mean blood pressure seems to increase with increasing age group. This is also the case for all the other measures see in the table. This seems to indicate that there is some relationship between systolic blood pressure and age.

From fig. 12 we see that there is a clear increase in blood pressure with age. Both the median, max and min increases with age. This is exactly what we saw in tab. 4 – a difference is that in the table the mean was given, while in the box plot the median is given –. What we can see clearer from the plot is that the distributions also become wider as age increase. The interquartile width of the middle age group is quite wide, and overlaps with the other groups. This can make this group a bit more difficult to get a significant result from.

All in all, both the numerical summary, tab. 4, and the box plot, fig. 12, seems to show the same: A increase in systolic blood pressure as age increases.

b)

We now want to look close at whether blood pressure varies between the different age groups. Above we used qualitative reasoning to to argue that this is the case. Now we want a quantified measure of this. To do this we will use a one-way ANOVA.

This may be just repetition of the caption...



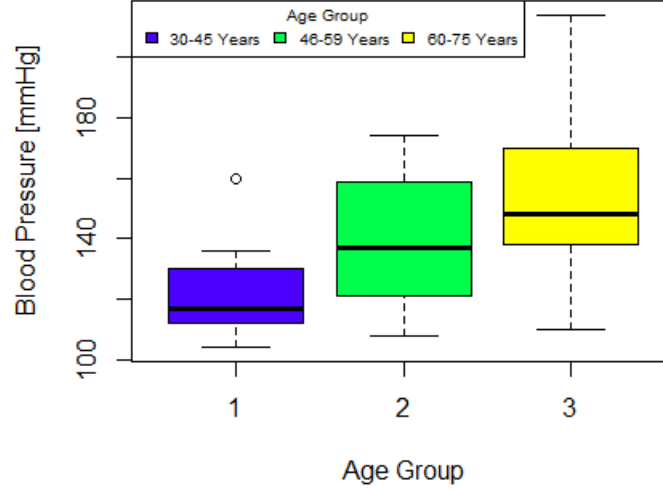


Figure 12: Boxplots showing the distribution of the blood pressure of the three age groups.

For the ANOVA we have the following hypotheses:

- $H_0 : \mu_1 = \mu_2 = \mu_3$
- $H_a : \mu_1 \neq \mu_2 \neq \mu_3$

In words: We have a null hypothesis that the mean blood pressure of all the groups are the same, and an alternative hypothesis that this is not that case. For the ANOVA to be valid, we need to assume that:

- The observations are independent of one another
- The observations in one group are a random sample with a normal distribution  $N(\mu_k, \sigma^2)$

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
Age Group	2	6535.39	3267.69	6.47	0.0043
Residuals	33	16670.25	505.16		

Table 5: A table showing the result of an ANOVA test on the blood pressure dataset. We see that we have a high F-value, leading to a small, significant p-value.

From tab. 5 we see that we get a small p-value of  $p = 0.0043$ . This is below the  $p = 0.05$  mark which is standard to use as the limit of significance. This means that we can throw out that null hypothesis, and conclude that blood pressure indeed varies across the age groups <sup>1</sup>

c)

We are now going to try to use a categorical regression on the data set. We are going to use treatment-contrast with group 1, the youngest, as the reference group. This means that our model will look like

<sup>1</sup>It is normally bad practice to interpret the p-value in an article, but since this is an oblig, it is done here.

$$y_i = \mu_1 + (\mu_2 - \mu_1) \cdot x_{1,i} + (\mu_3 - \mu_1) \cdot x_{2,i} + \epsilon_i, \quad (4)$$

where  $\epsilon_i$  is a normally distributed error term,  $\mu_j$  is the mean of the different age groups, and

$$x_{j-1,i} = \begin{cases} 1 & \text{for } i \text{ in group } j \\ 0 & \text{else} \end{cases}. \quad (5)$$

This means that if patient is in group 2, then  $x_{1,i} = 1$  and  $x_{2,i} = 0$ , and vis versa if the patient is in group 3. If the patient is in group 1, both will be 0.

	Estimate	Std. Error	t value	Pr(>  t )
(Intercept)	122.1667	6.4882	18.83	0.0000
Age Group 2	16.9167	9.1757	1.84	0.0742
Age Group 3	33.0000	9.1757	3.60	0.0010
Residual standard error: 22.48 on 33 degrees of freedom				
Multiple R-squared: 0.2816, Adjusted R-squared: 0.2381				
F-statistic: 6.469 on 2 and 33 DF, p-value: 0.004263				

Table 6: The summary of the categorical regression of the blood pressure with respect to the age groups.

In tab. 6 we see the summary of the regression model with the age groups as a categorical predictor variable. The intercept in the table is the mean of the youngest age group  $\mu_1 = 122.167$  – our reference group –. The slopes of the two variables can be interpreted as how much the mean of blood pressure increases if you are in that age group, compared to the reference group. E.g. if you are in age group 2, we will expect your blood pressure to be 33.0 mmHg higher than the mean of the youngest group, i.e. it is expected to be 155.167 mmHg.

We can look at the p-values of tab. 6 we see that the p-value for age group 3 is  $p = 0.001$ . This means that the increase of 33 mmHg is significant. But for group 2 the p-value is just 0.0742, which means that this increase is not significant, so we can not say that we have an 17 mmHg increase compared with the reference. This might have to do with what we see in the box plots 12, where the distribution of the second age group is very wide. The reason might also be due to a non linear dependence on age.

We also see that we have a low  $R^2 = 0.2816$ , which means that the blood pressure is not fully described by the age groups.

But all in all we see that the we have a good regression model for group 1 and 3, with the change in group 2 not being significant. So there is some significant variation between the age groups, as we saw with the ANOVA.

Look  
over and  
rewrte...