

Yongkang Long 171022

Question 1

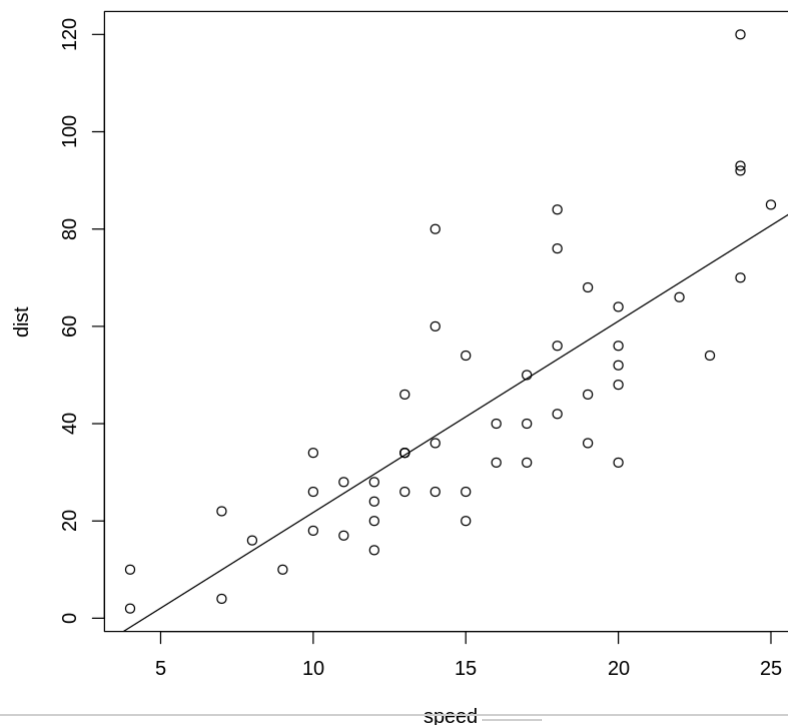
For this question we will use the data set cars.

(i) Plot dist as a function of speed. Fit a simple linear regression model of dist as a function of speed. Add the regression line to the previous plot. Obtain the summary for this regression. Obtain an estimator for the error variance. Observe the value for the intercept and comment.

```
In [1]: 1 library(car)
```

Loading required package: carData

```
In [2]: 1 plot(dist ~ speed, data = cars)
2 abline(lm(dist ~ speed, data = cars))
```



```
In [3]: 1 model1 <- lm(dist ~ speed, data = cars)
        2 summary(model1)
```

Call:
lm(formula = dist ~ speed, data = cars)

Residuals:

	Min	1Q	Median	3Q	Max
	-29.069	-9.525	-2.272	9.215	43.201

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-17.5791	6.7584	-2.601	0.0123 *
speed	3.9324	0.4155	9.464	1.49e-12 ***

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

Residual standard error: 15.38 on 48 degrees of freedom
Multiple R-squared: 0.6511, Adjusted R-squared: 0.6438
F-statistic: 89.57 on 1 and 48 DF, p-value: 1.49e-12

```
In [4]: 1 #the estimated variance is
        2 summary(model1)$sigma^2
```

236.531688564477

1 The incercept is -17.5791, which means that the barking distacne is -17.5791 miles for speed 0. However, it is not true in real. Look at the p-value is 0.0123, marginally significant. If we use alpha=0.01, we accept the null hypothesis that incercept is 0.

(ii)Based on your comments to the previous section, fit a model without an intersect. Draw a scatterplotand add the two regression lines. Obtain a summary for the new regression and comment on thedifferences with the previous model, including the estimated error variance.

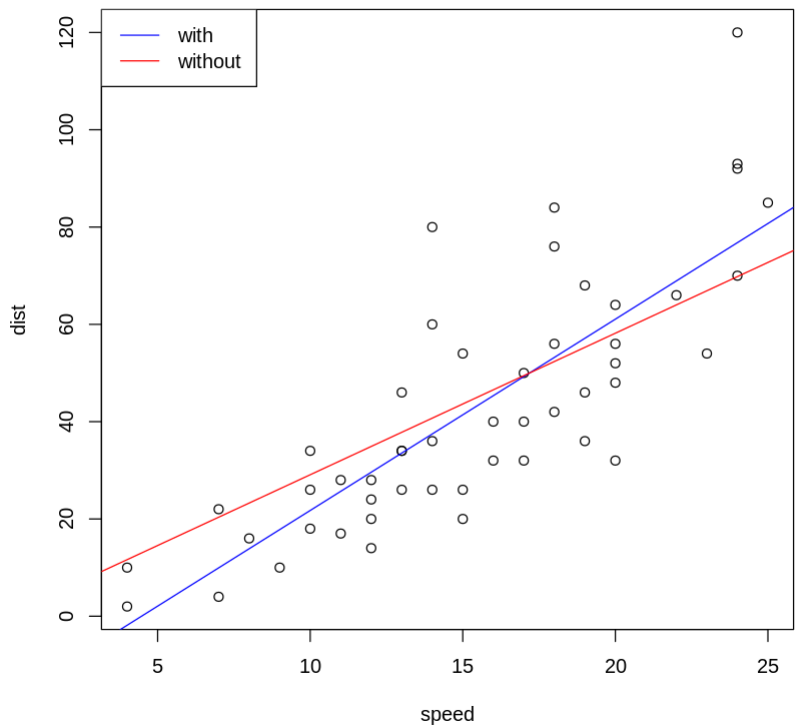
```
In [5]: 1 model2 <- lm(dist ~ -1 + speed, data = cars)
        2 summary(model2)
```

Call:

```
In [6]: 1 #the estimated variance is
        2 summary(model2)$sigma^2

264.36279259209

In [7]: 1 plot(dist ~ speed, data = cars)
        2 abline(model1,col = 'blue')
        3 abline(model2,col = 'red')
        4 legend('topleft',c('with', 'without'),col=c('blue', 'red'),lty=c(1,1))
```



1 The R^2 without intercept is higher than R^2 with intercept and the estimated variance is no difference. It means we prefer the model without intercept. Actually, there is no intercept for dist and speed in real life

(iii) We want to compare the predictive power of these two models. Using the same procedure as in exercise1 of the list for week 9, compare the predictive power of both models and comment on your results

```
In [8]: 1 attach(cars)

In [9]: 1 n = length(dist)
```

```
In [11]: 1 for (i in 1:n) {
          2   xx <- speed[-i]
          3   yy <- dist[-i]
          4   model1 <- lm(yy ~ xx)
          5   model2 <- lm(yy ~ -1 + xx)
          6   pred.values[i,1] = predict(model1,data.frame(xx = speed[i]))
          7   pred.values[i,2] = predict(model2,data.frame(xx = speed[i]))
          8 }
          9 pred.values <- cbind(pred.values,dist)
         10 colnames(pred.values) <- c('P1', 'P2', 'O')
```

```
In [12]: 1 head(pred.values)
```

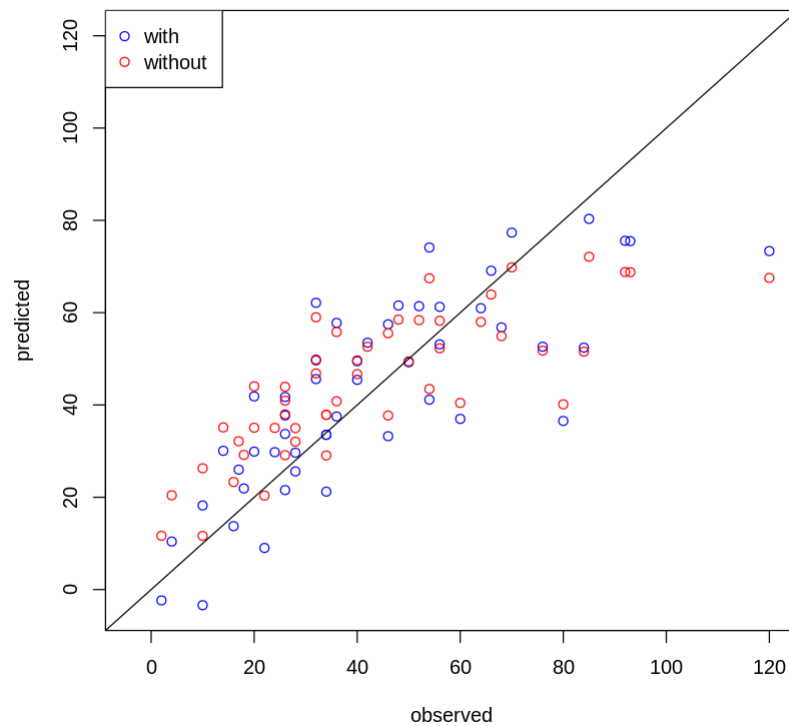
A matrix: 6 × 3 of type dbl

	P1	P2	O
	-2.348991	11.64820	2
	-3.387122	11.63851	10
	10.405805	20.42477	4
	9.019622	20.35784	22
	13.744937	23.30842	16
	18.222888	26.28189	10

```

In [13]: 1 plot(pred.values[,3],pred.values[,1],
2             col = 'blue', ylab = 'predicted',
3             xlab = 'observed',
4             xlim = (range(pred.values)+c(-.5,+.5)),
5             ylim = (range(pred.values)+c(-.5,+.5)))
6 points(pred.values[,3],pred.values[,2], col = 'red')
7 legend('topleft',c('with', 'without'),col=c('blue', 'red'),pch=c(1,1))
8 abline(0,1)

```

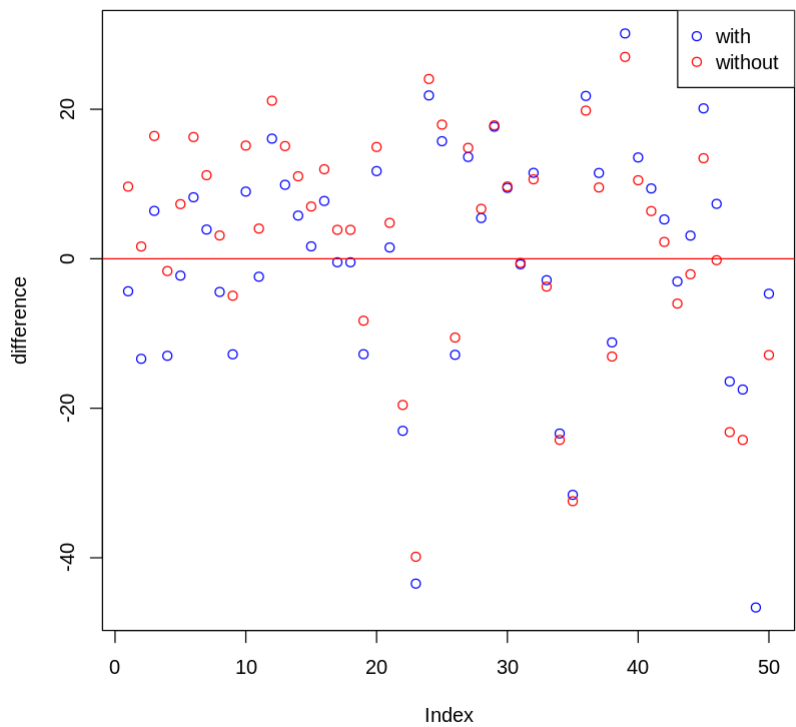


```

1 Only Slightly different between predicted values with and without
  intercept

```

```
In [14]: 1 plot(pred.values[,1]-pred.values[,3], col = 'blue', ylab = 'difference')
2 abline(h=0,col='red')
3 points(pred.values[,2]-pred.values[,3], col = 'red')
4 legend('topright',c('with','without'),col=c('blue','red'),pch=c(1,1))
```



1 Only Slightly different between errors with and without intercept

```
In [15]: 1 #mean absolute errors (MAE) of model with intercept
2 mean(abs(pred.values[,1]-pred.values[,3]))
```

12.0591786486375

```
In [16]: 1 #mean absolute errors (MAS) of model with outintercept
2 mean(abs(pred.values[,2]-pred.values[,3]))
```

12.9782367734261

1 No difference between MAE with two models. Removing intercept is powerful

Question 2

For this question use the data set data1

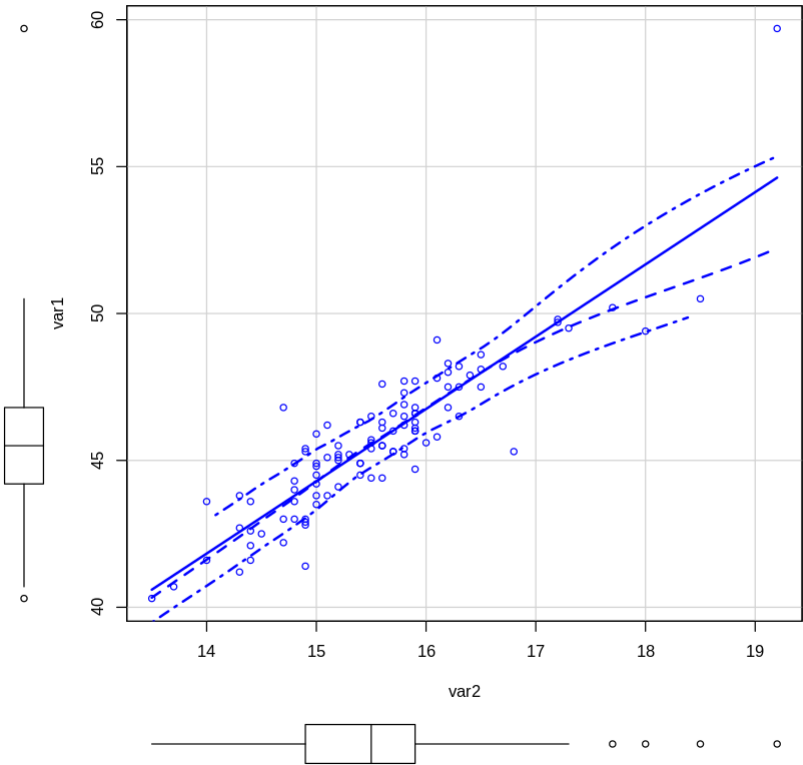
function of var2. Fit a simple linear regression and add the line to the plot. Comment. Obtain a summary of the regression.

```
In [17]: 1 data1 = read.table('data1')
         2 head(data1)
         3 attach(data1)
```

A data.frame: 6 × 2

	var1	var2
	<dbl>	<dbl>
1	46.8	15.9
2	45.2	15.2
3	46.6	15.9
4	44.9	15.0
5	46.1	15.6
6	45.1	15.2

```
In [18]: 1 scatterplot(var1 ~ var2, data = data1)
```



1	The local smooth regression line fit the regression line in low values but is different in high values		
---	--	--	--

No documentation for 'cex' in specified packages and libraries: you could try '??cex'

```
In [19]: 1 model1 <- lm(var1 ~ var2, data = data1)
2 plot(var1 ~ var2, data = data1)
3 abline(model1,col = 'blue')
4 summary(model1)
```

Call:
lm(formula = var1 ~ var2, data = data1)

Residuals:

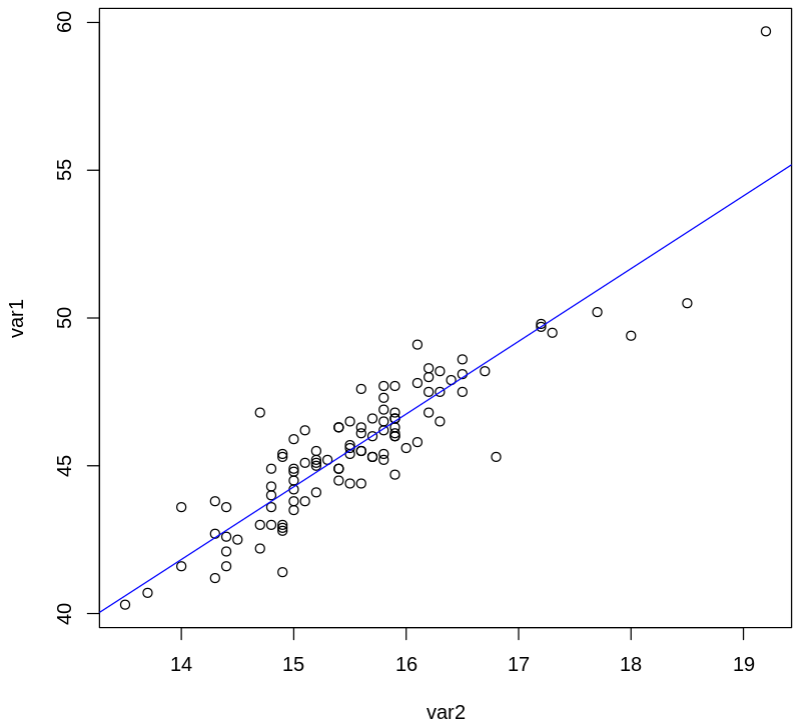
Min	1Q	Median	3Q	Max
-3.4183	-0.7043	-0.0072	0.6049	5.0765

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	7.3826	1.9110	3.863	0.000199 ***
var2	2.4605	0.1227	20.060	< 2e-16 ***

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

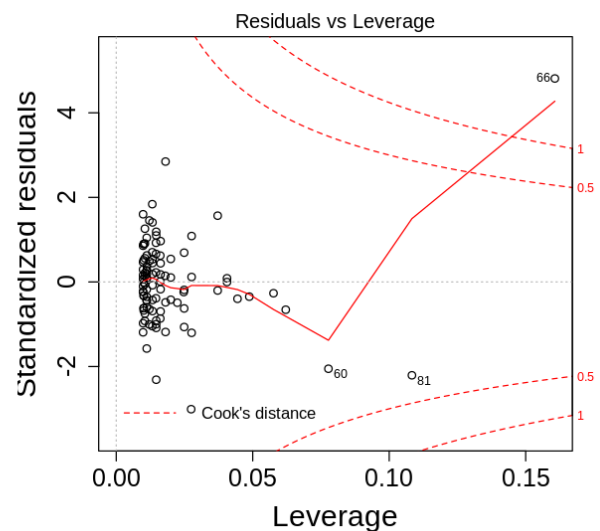
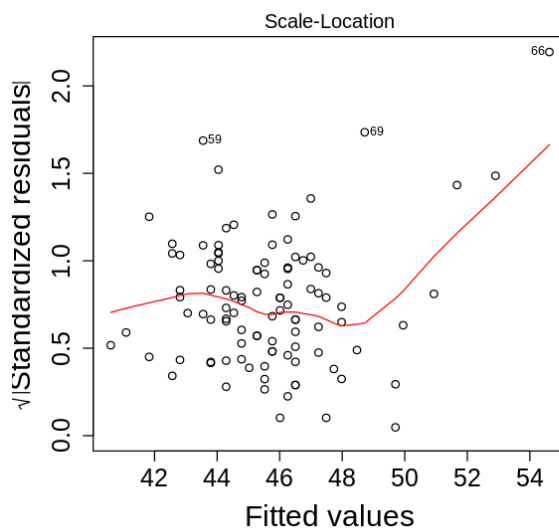
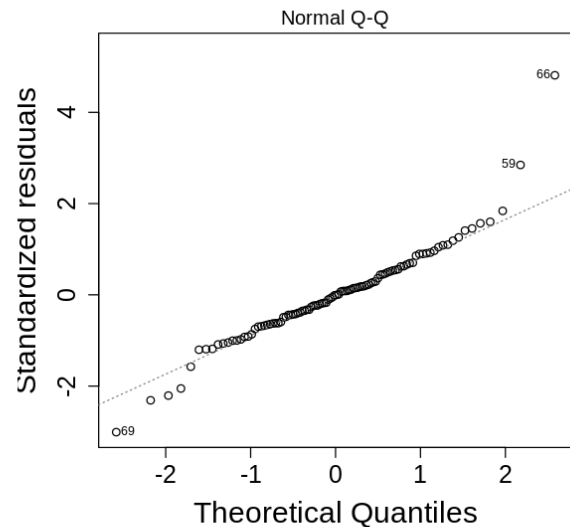
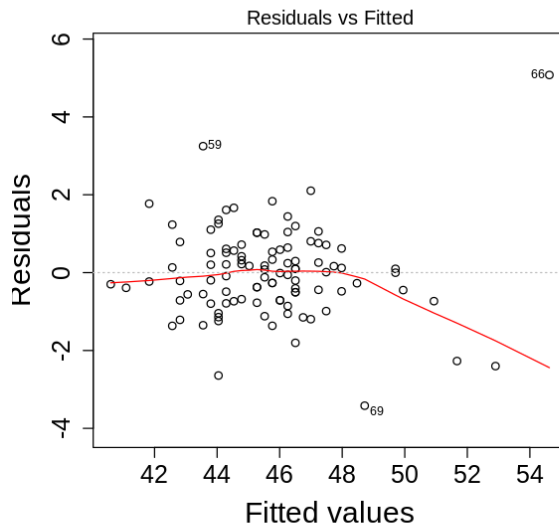
Residual standard error: 1.152 on 100 degrees of freedom
Multiple R-squared: 0.801, Adjusted R-squared: 0.799
F-statistic: 402.4 on 1 and 100 DF, p-value: < 2.2e-16



```
1 The p value is significant and R^2 is 0.8, it seems good. However,
we can observe an outlier obviously which means that there is room
for improvement.
```

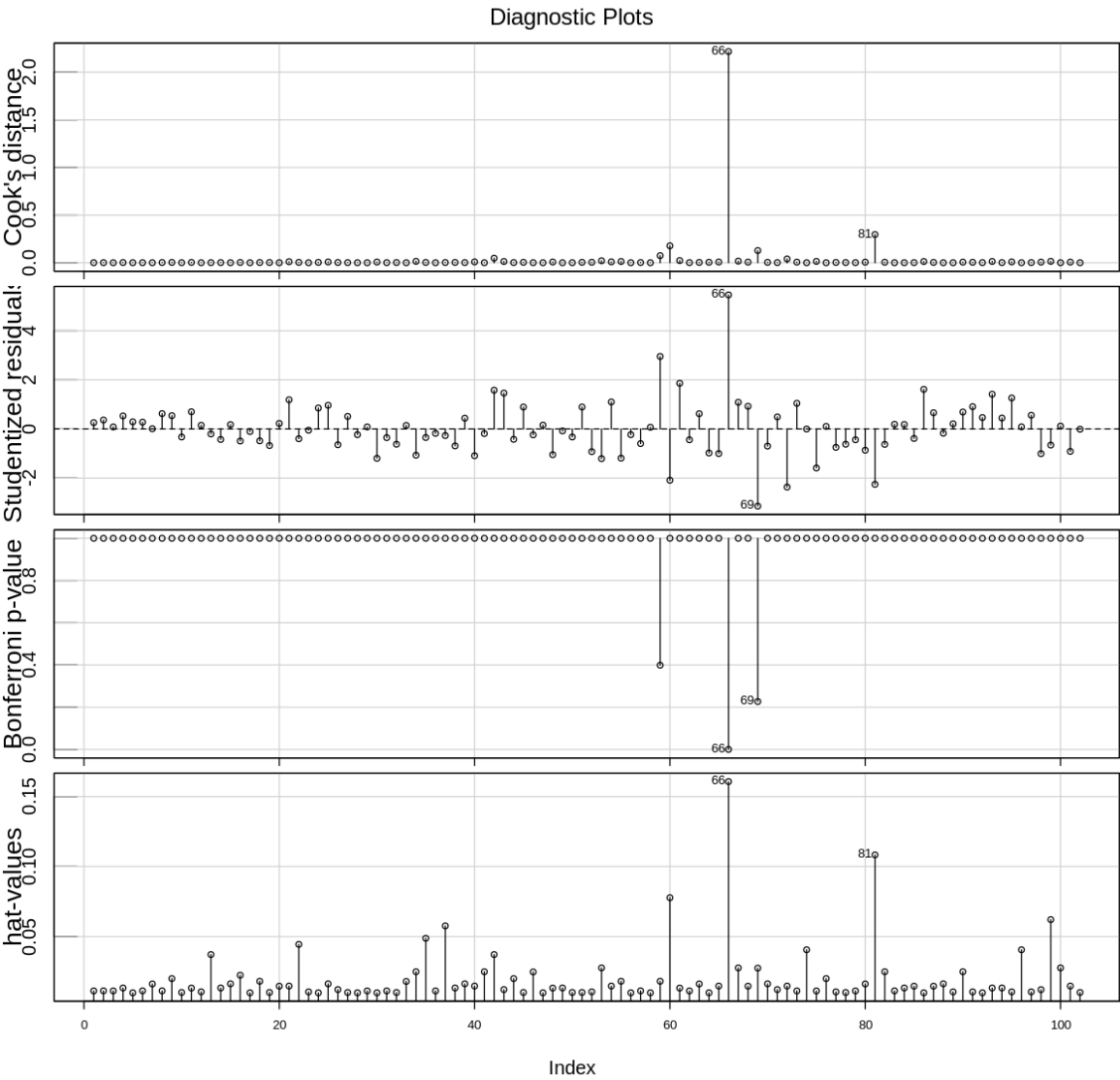

(ii) Draw the diagnostic plots. Do you identify any point as an outlier? If you do, which point is this? Can you identify this point in the initial scatterplot?

```
In [20]: 1 options(repr.plot.width=10, repr.plot.height=10)
2 par(mfrow = c(2,2))
3 plot(model1,cex.axis=1.5,cex.lab=1.8,ps=10)
4 par(mfrow = c(1,1))
```



- 1 The residuals plot and QQ plot shows the variance are uniform and normality is valid. Outlier 66 has very large cook distance. 60 and 81 maybe outlier but with cook's distance < 0.5

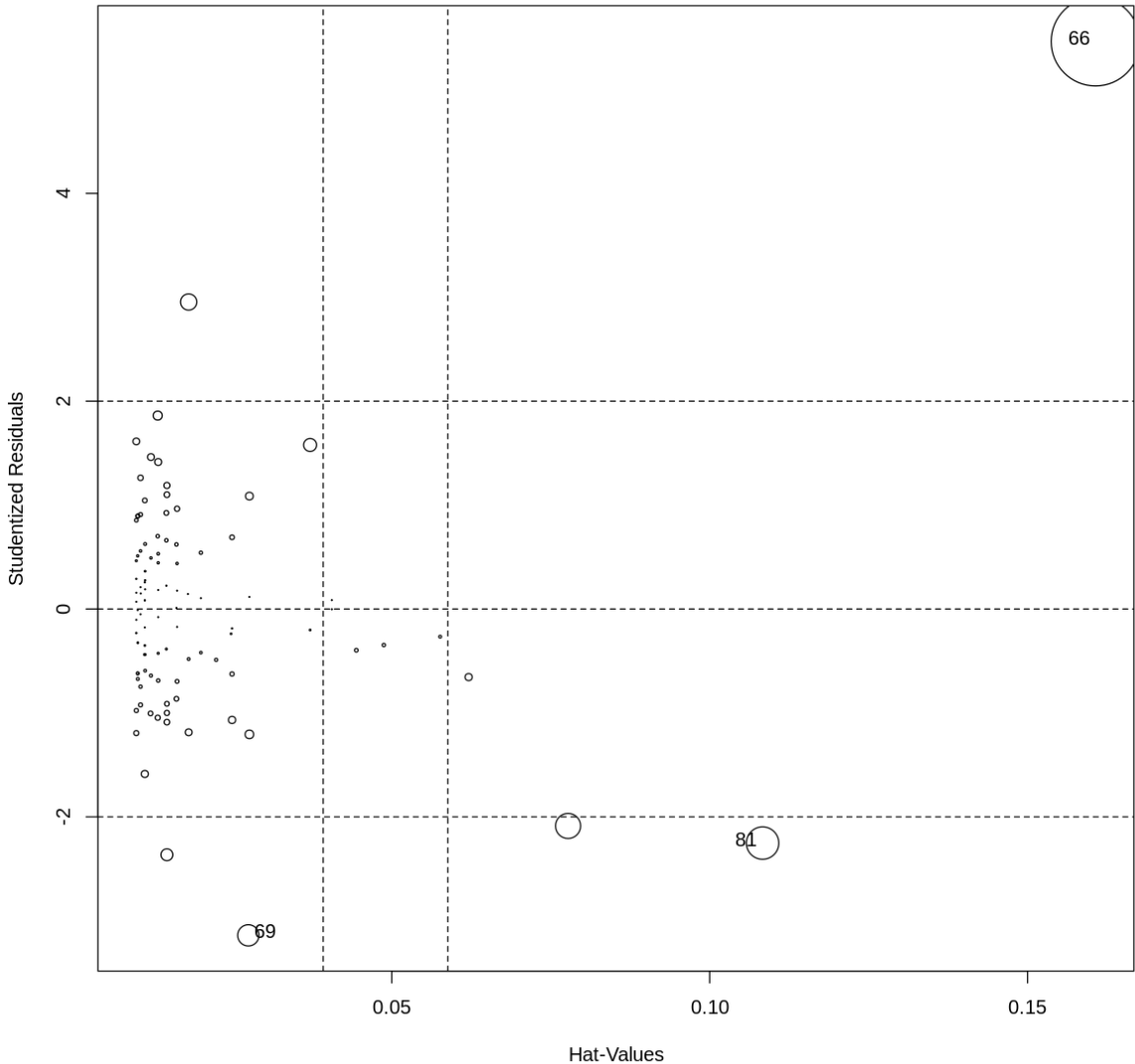
```
In [21]: 1 influenceIndexPlot(model1,cex.lab=2,cex.axis=1.5)
```



```
In [22]: 1 influencePlot(model1)
```

A data.frame: 3 × 3

	StudRes	Hat	CookD
	<dbl>	<dbl>	<dbl>
66	5.461408	0.16068783	2.2163170
69	-3.140454	0.02744525	0.1278291
81	-2.252573	0.10832633	0.2961511

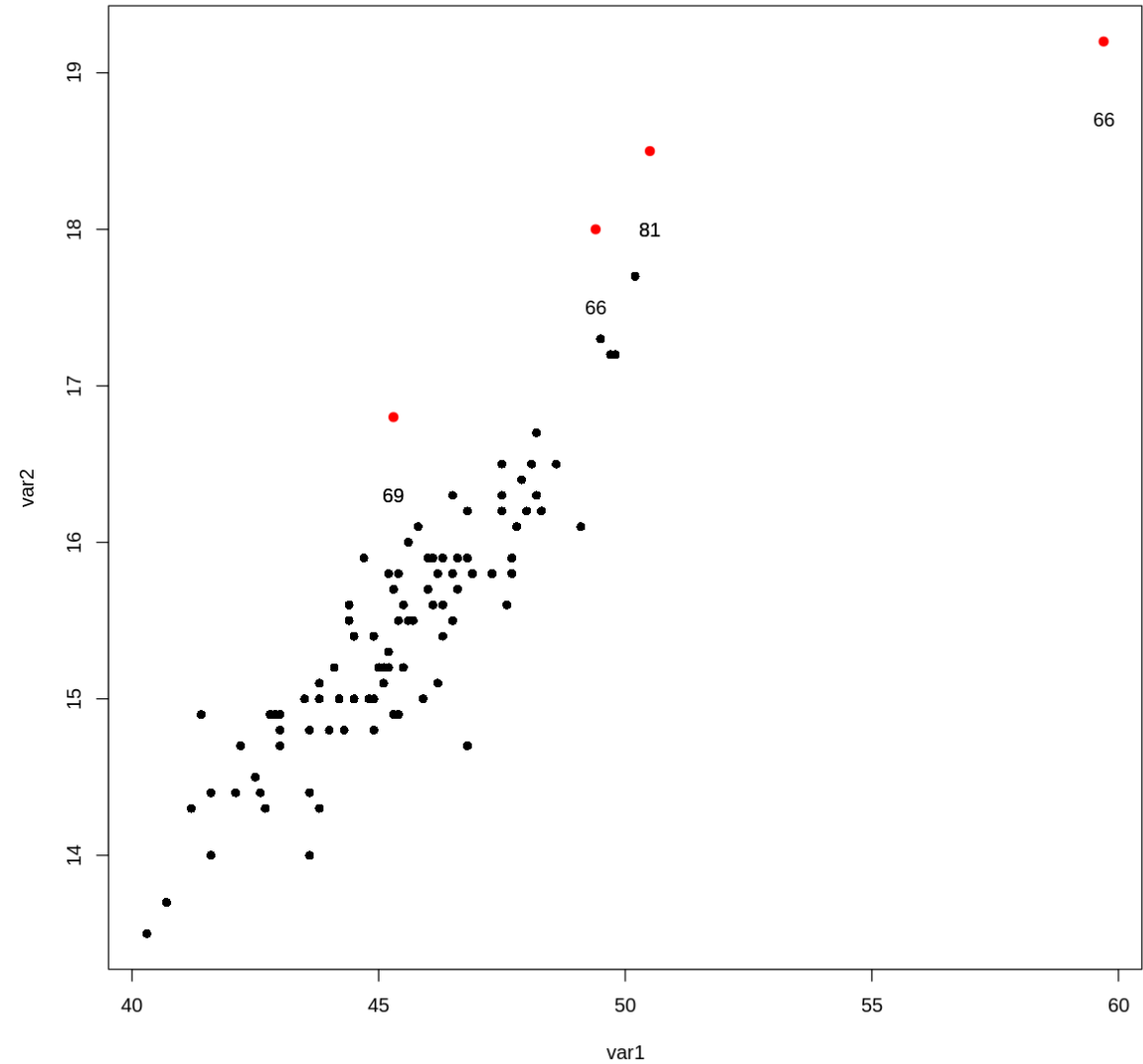


```
1 We now get three outliers 60,69,66,and 81
```

```
In [23]: 1 plot(var1 ,var2,pch=16)
2 #summary(model2)
3 points(var1[c(60,66,69,81)], var2[c(60,66,69,81)],pch=19, col='red')
4 text(var1[60],var2[60]-.5, '66');
```

No documentation for 'cex' in specified packages and libraries; you could try "?cex"

```
5 text(var1[69],var2[69]-.5,'69');
6 text(var1[81],var2[81]-.5,'81');
7 text(var1[66],var2[66]-.5,'66');
8 text(var1[69],var2[69]-.5,'69');
9 text(var1[81],var2[81]-.5,'81');
```



1 We could not identify 60, 69 and 81 point in the initial scatterplot

(iii)Fit a new regression model excluding the outlier(s) that you identified in the previous section. Draw a scatterplot with both regression lines. Compare the summary tables. Draw the diagnostic plots and comment.

In [24]:

No documentation for 'cex' in specified packages and libraries: you could try '??cex'

1 options(repr.plot.width=10, repr.plot.height=10)

2 data_60 = data[-c(60,66,69,81),]

```
3 model2 <- lm(var1 ~ var2, data = data1_eo)
4 plot(var1 ~ var2,data=data1,pch=16)
5 abline(model1,col = 'red',lty=3, lwd=2)
6 abline(model2,col = 'blue',lty=4, lwd=2)
7 legend('topleft',c('model1','model2'),col=c('red','blue'),lty=c(1,1)
8 summary(model2)
9
```

Call:
lm(formula = var1 ~ var2, data = data1_eo)

Residuals:

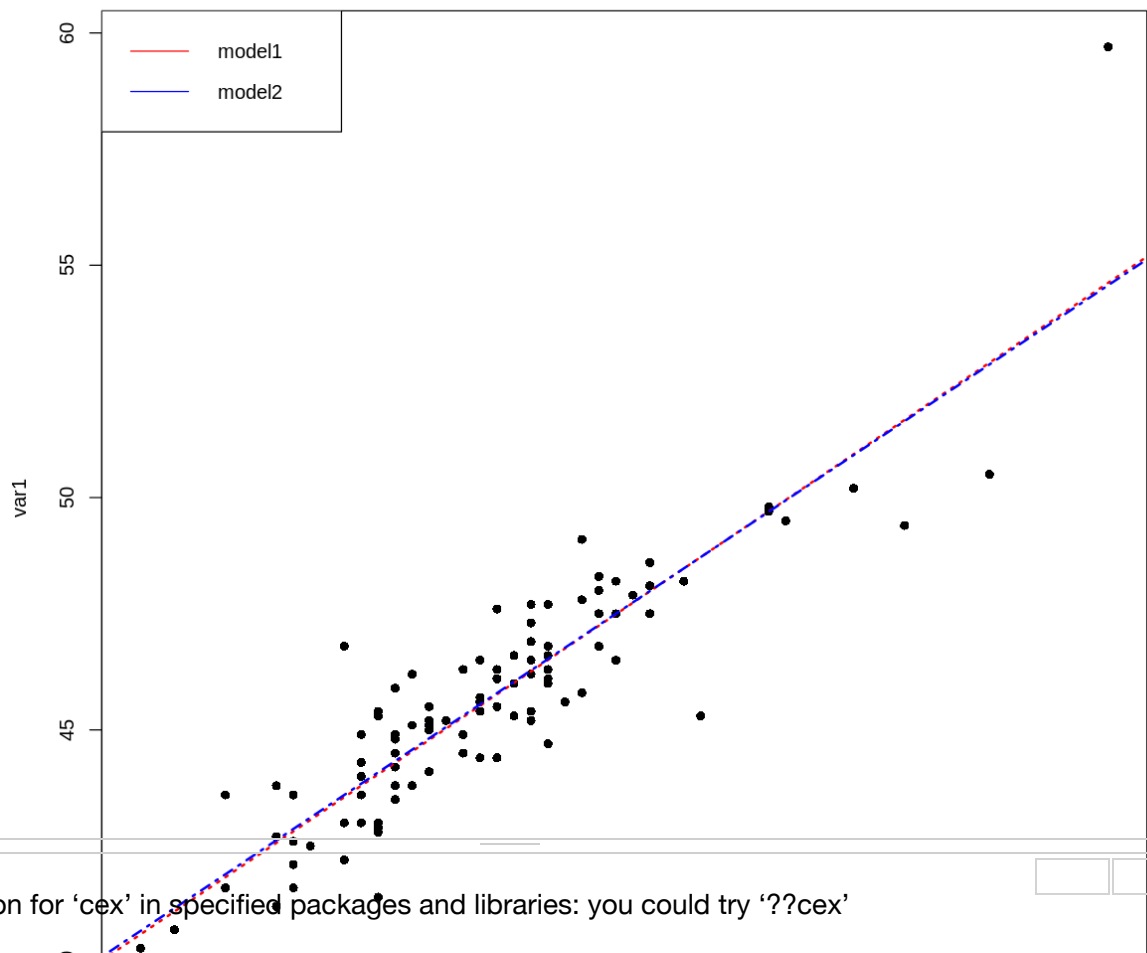
Min	1Q	Median	3Q	Max
-2.6866	-0.6084	0.0034	0.5680	3.2010

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	7.7628	1.9086	4.067	9.75e-05 ***
var2	2.4378	0.1234	19.756	< 2e-16 ***

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

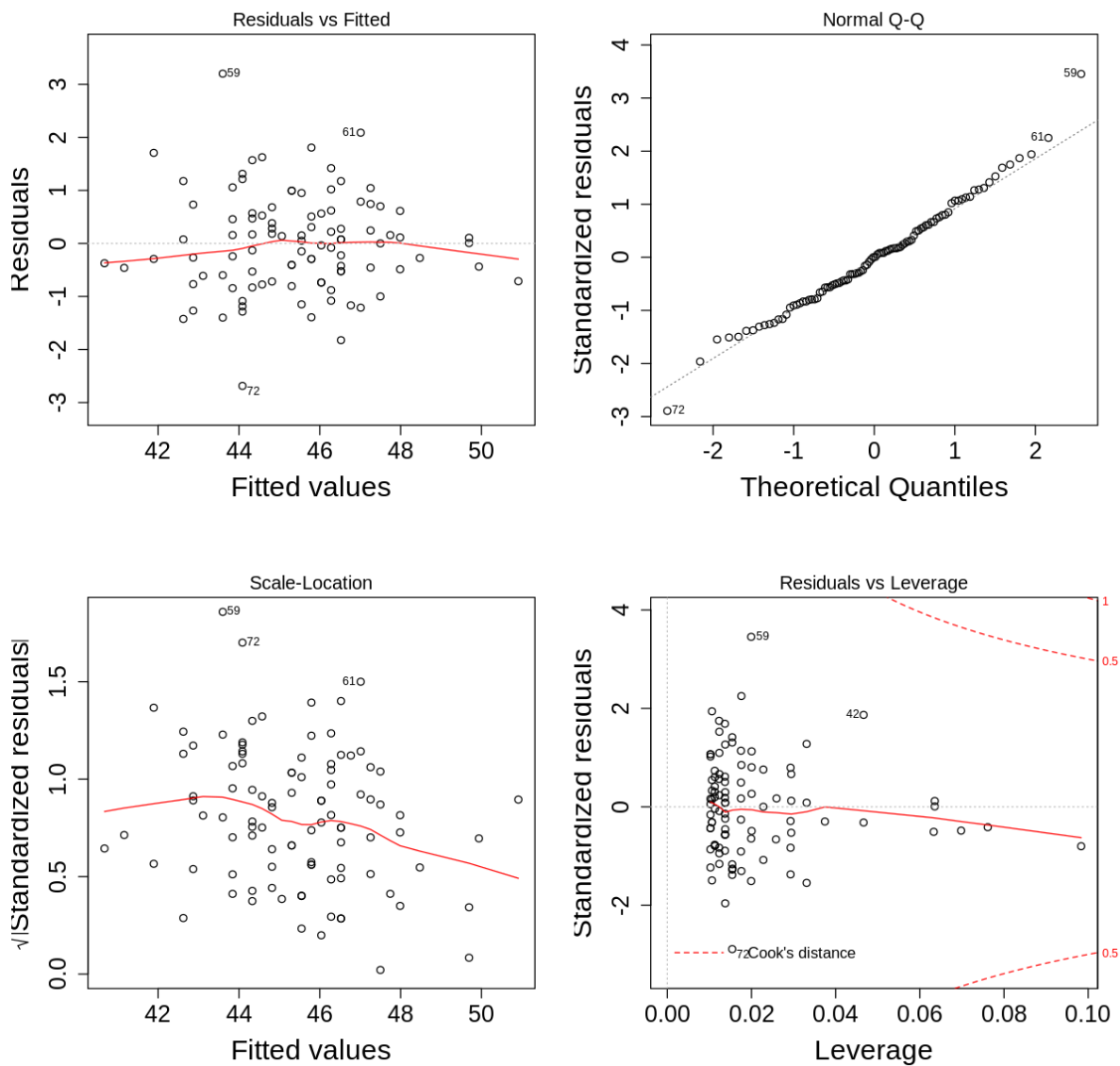
Residual standard error: 0.9361 on 96 degrees of freedom
Multiple R-squared: 0.8026, Adjusted R-squared: 0.8005
F-statistic: 390.3 on 1 and 96 DF, p-value: < 2.2e-16



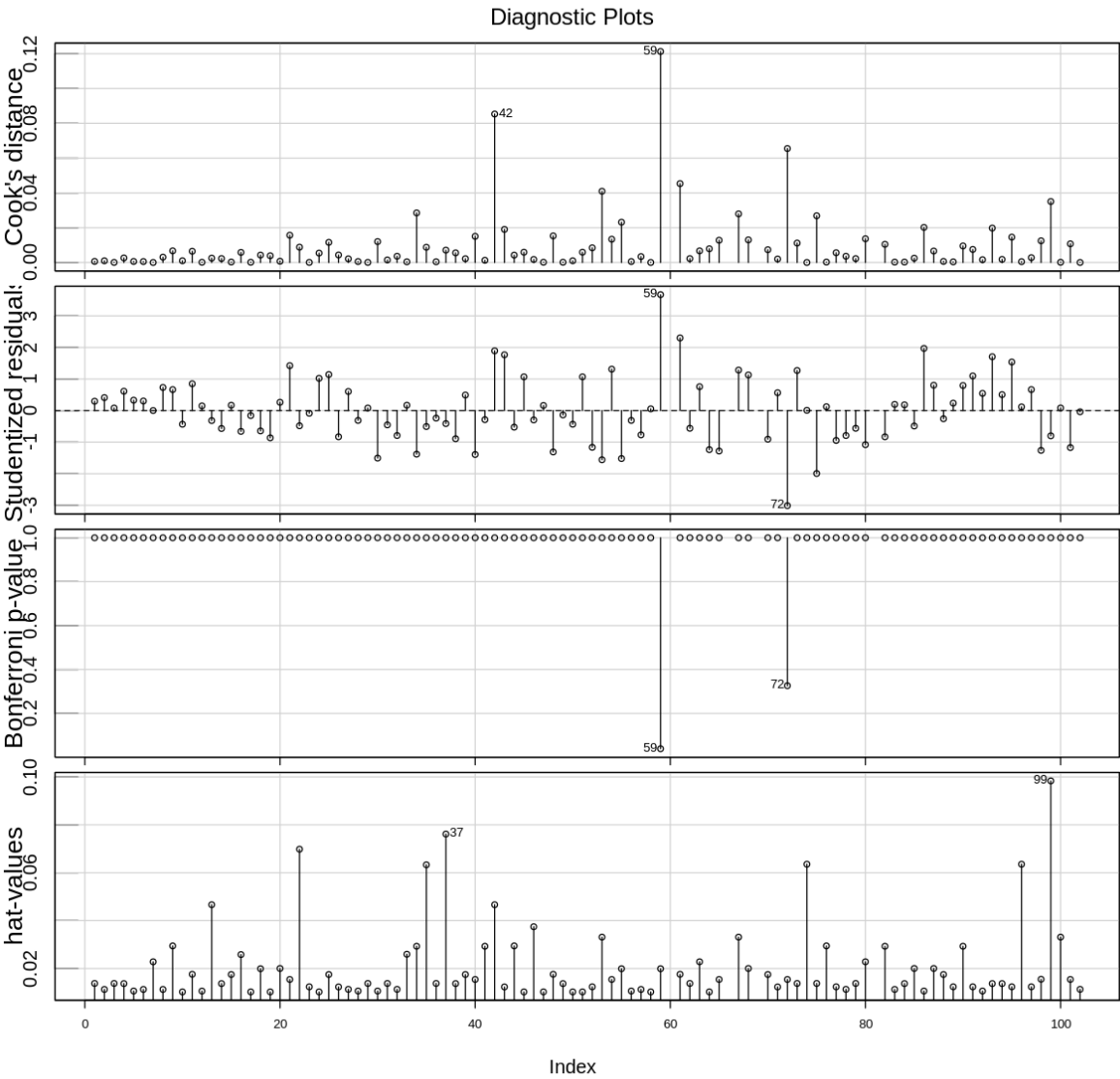
1

The intercept,slope,and R^2 do not change significantly.

```
In [25]: 1 options(repr.plot.width=10, repr.plot.height=10)
2 par(mfrow = c(2,2))
3 plot(model2,cex.axis=1.5,cex.lab=1.8,ps=10)
4 par(mfrow = c(1,1))
```



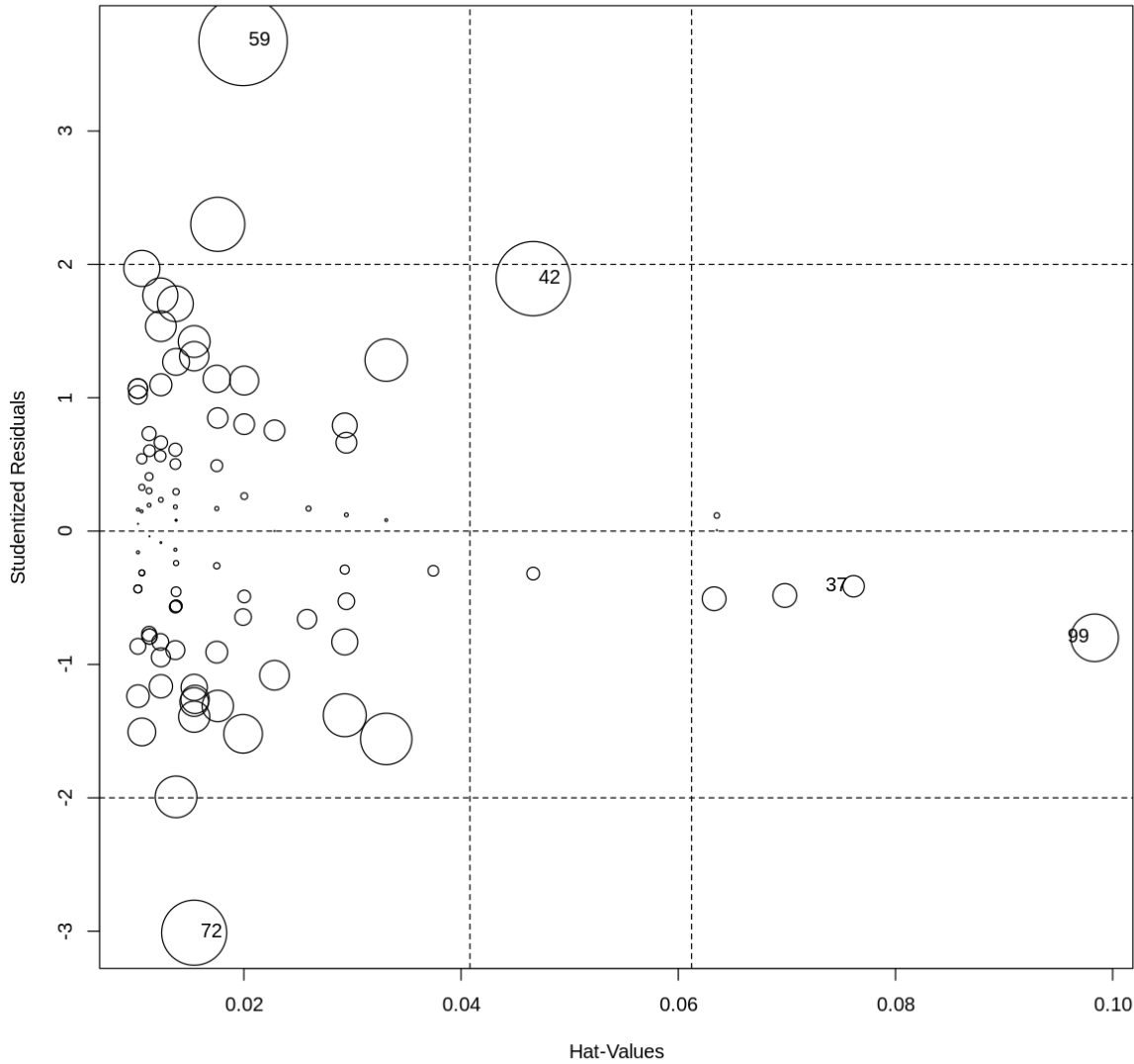
```
In [26]: 1 influenceIndexPlot(model2,cex.lab=2,cex.axis=1.5)
```



```
In [27]: 1 influencePlot(model2)
```

A data.frame: 5 × 3

	StudRes	Hat	CookD
	<dbl>	<dbl>	<dbl>
37	-0.4134772	0.07614503	0.007106858
42	1.8930925	0.04663819	0.085361742
59	3.6718182	0.01992598	0.121284644
72	-3.0115733	0.01542194	0.065522891
99	-0.8001545	0.09833900	0.035045381



```
1 Now the diagnostic plot looks better
```

No documentation for 'cex' in specified packages and libraries: you could try '??cex'

Question 3

For this question use the data set data2.

(i) Read data2 and plot yval as a function of xval. Fit a simple linear regression and add the regression line to the plot. Comment. Obtain a summary for the regression and draw the diagnostic plots. Comment on the results

In [28]:

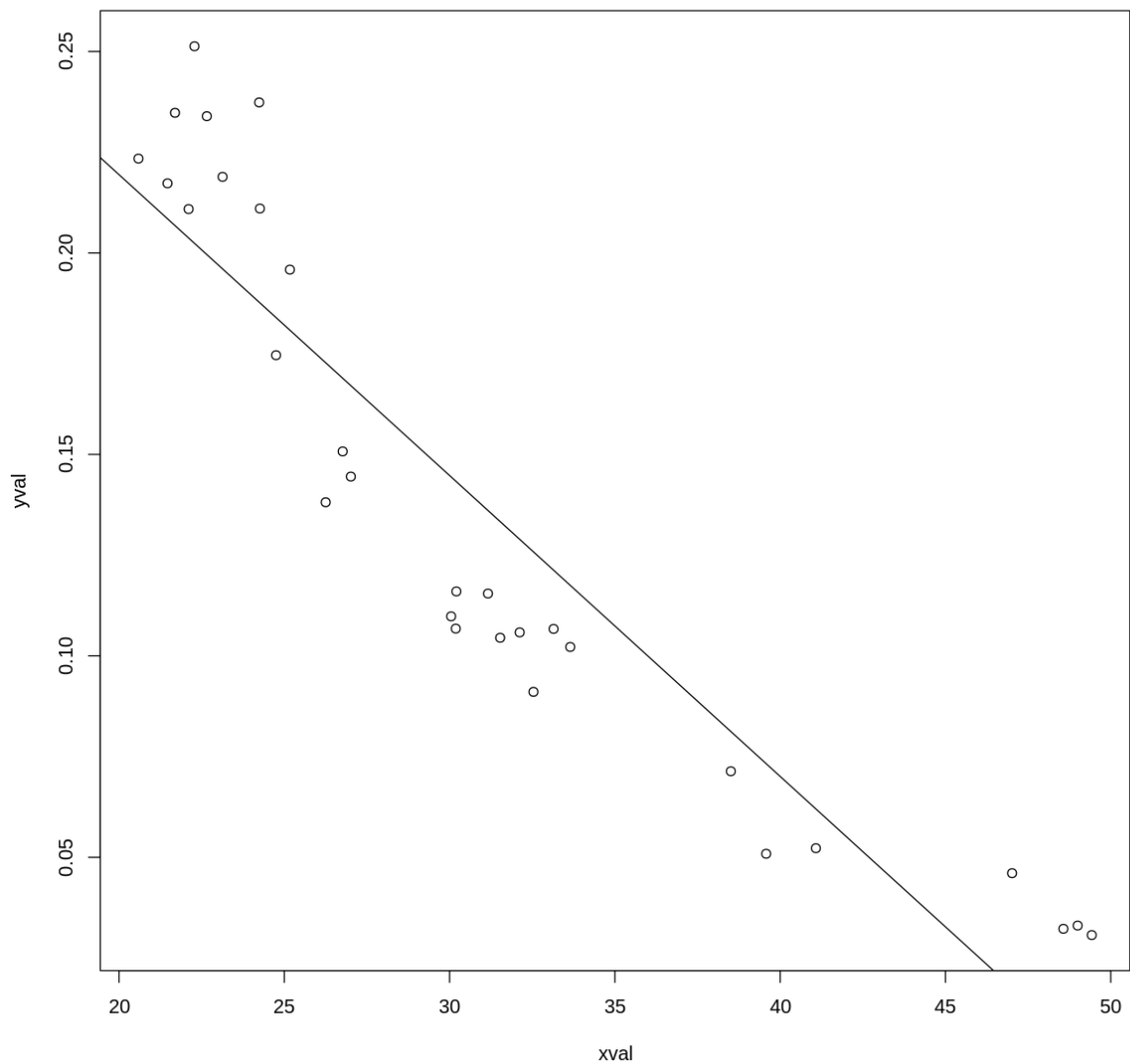
```
1 data2 = read.table('data2')
2 head(data2)
```

A data.frame: 6 × 2

	xval	yval
	<dbl>	<dbl>
1	31.52991	0.1044869
2	33.14513	0.1066878
3	32.11964	0.1058220
4	31.16068	0.1154637
5	32.53664	0.0910629
6	23.13162	0.2188879

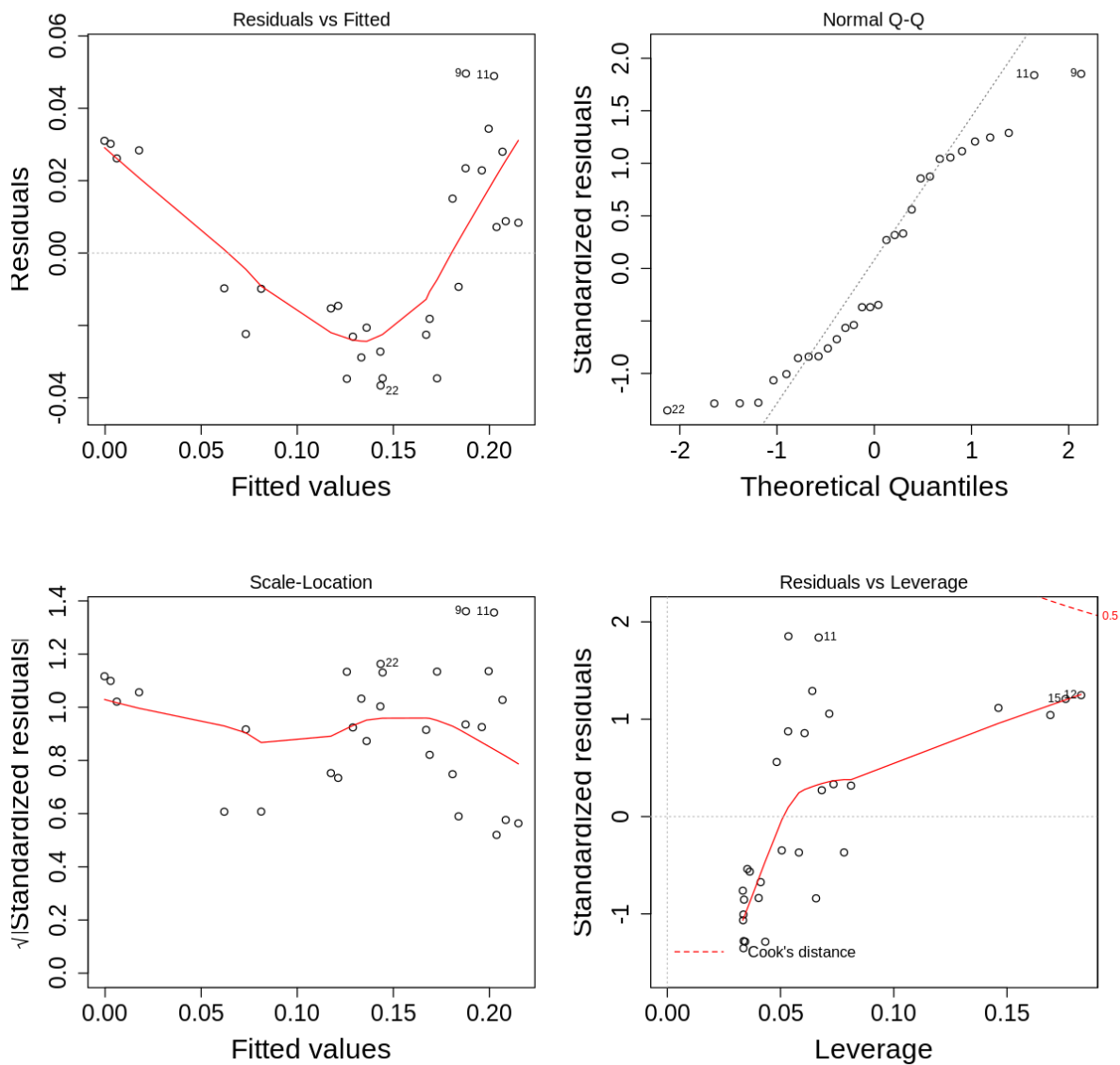
In [29]:

```
1 options(repr.plot.width=10, repr.plot.height=10)
2 plot(yval ~ xval, data = data2)
3 model1 <- lm(yval ~ xval, data = data2)
4 abline(model1)
5 summary(model1)
```



1	The p value is significant and R^2 is 0.85, it seems good.
---	--

```
In [30]: 1 options(repr.plot.width=10, repr.plot.height=10)
2 par(mfrow = c(2,2))
3 plot(modell,cex.axis=1.5,cex.lab=1.8,ps=10)
4 par(mfrow = c(1,1))
```



1 From residuals plot, we know that there is a quadratic pattern but not for the standardized residuals plot and normality do not fit the model

```
In [31]: 1 ncvTest(modell)
2 ##### The variances are uniform
```

Non-constant Variance Score Test
Variance formula: ~ fitted.values
Chisquare = 0.03720501, Df = 1, p = 0.84705

```
In [32]: 1 shapiro.test(modell$residuals)
2 ##### The normality assumption is not satisfied
```

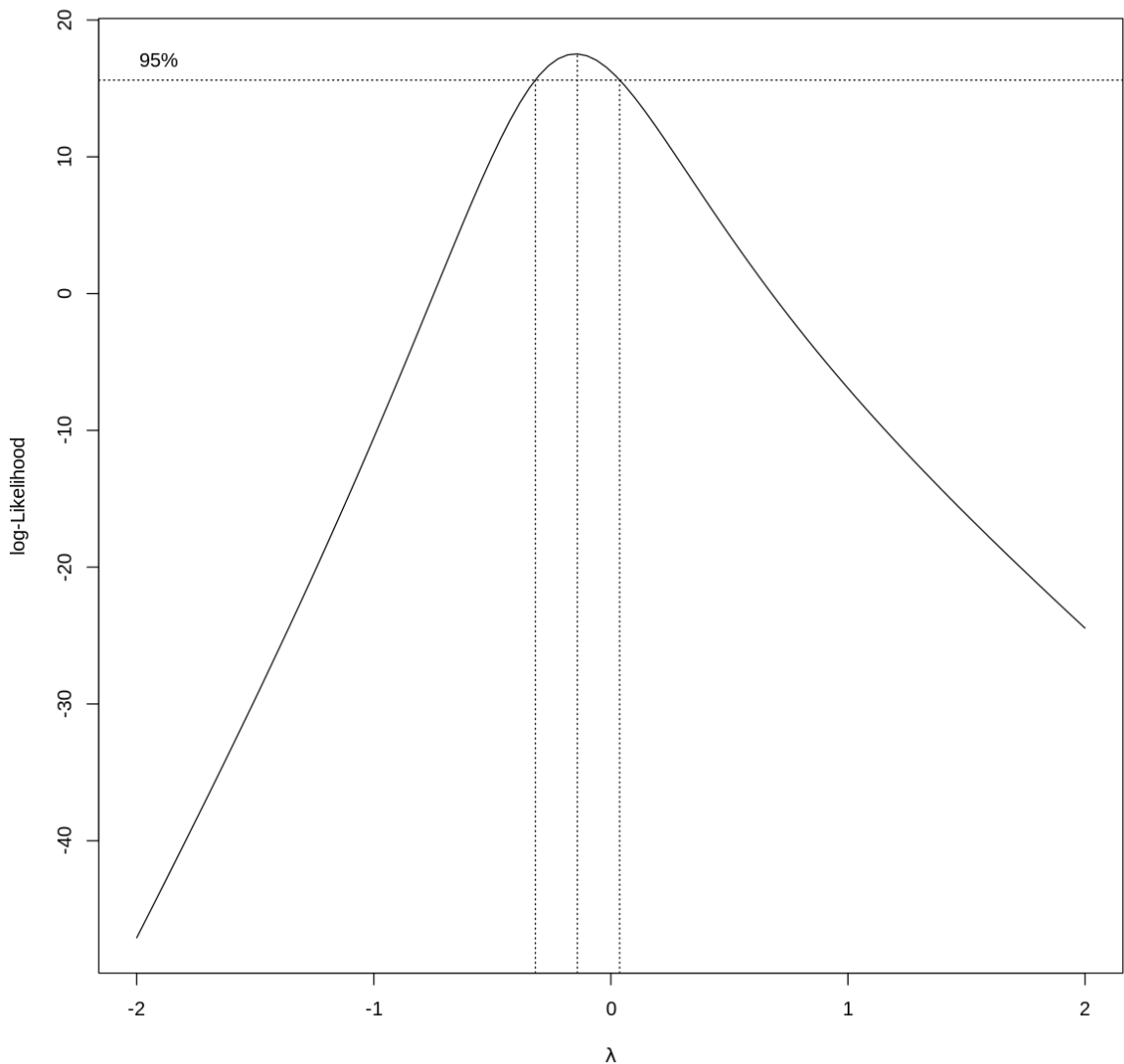
No documentation for 'cex' in specified packages and libraries: you could try '?cex'

Shapiro-Wilk normality test

```
data:  model1$residuals
W = 0.91828, p-value = 0.0242
```

(ii) Use the function `boxcox` on the package `MASS` with the argument set to the model you fitted in (i).

```
In [33]: 1 library(MASS)
          2 boxcox(model1)
```



```
1 Becuase 0 is i the interval, we can chose lamda = 0 and use
  log(transform)
```

use a logarithmic transformation for yval and fit a new model. Obtain a summary of the new regression and compare with the previous one. Draw the diagnostic plots and compare with the previous results.

```
In [34]: 1 data2$logyval = log(data2$yval)
          2 data2$logxval = log(data2$xval)
```

```
In [35]: 1 plot(logyval ~ logxval, data = data2)
2 model2 <- lm(logyval ~ logxval, data = data2)
3 abline(model2)
4 summary(model2)
```

Call:
lm(formula = logyval ~ logxval, data = data2)

Residuals:

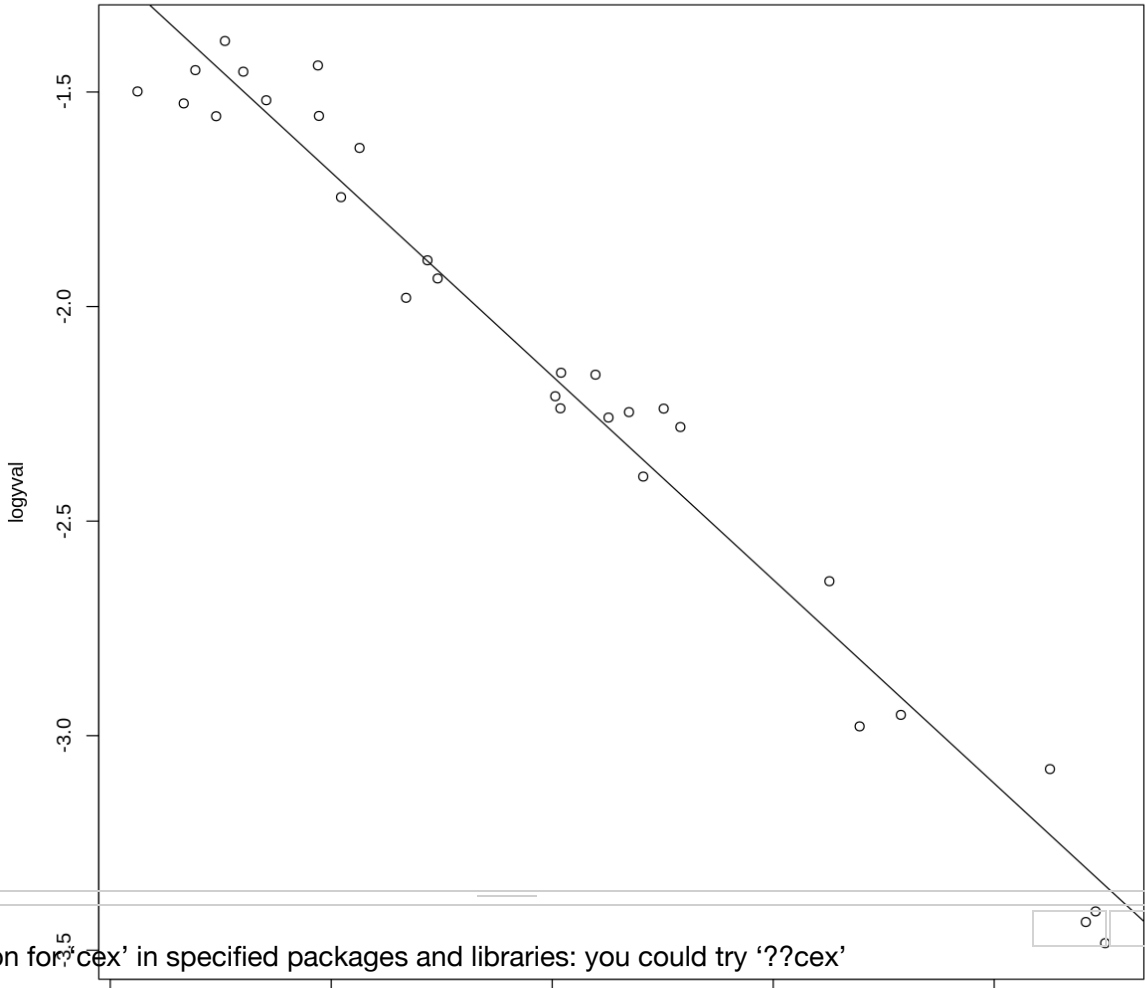
	Min	1Q	Median	3Q	Max
	-0.227479	-0.074795	-0.008457	0.092254	0.220557

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	5.90701	0.27207	21.71	<2e-16 ***
logxval	-2.37322	0.07983	-29.73	<2e-16 ***

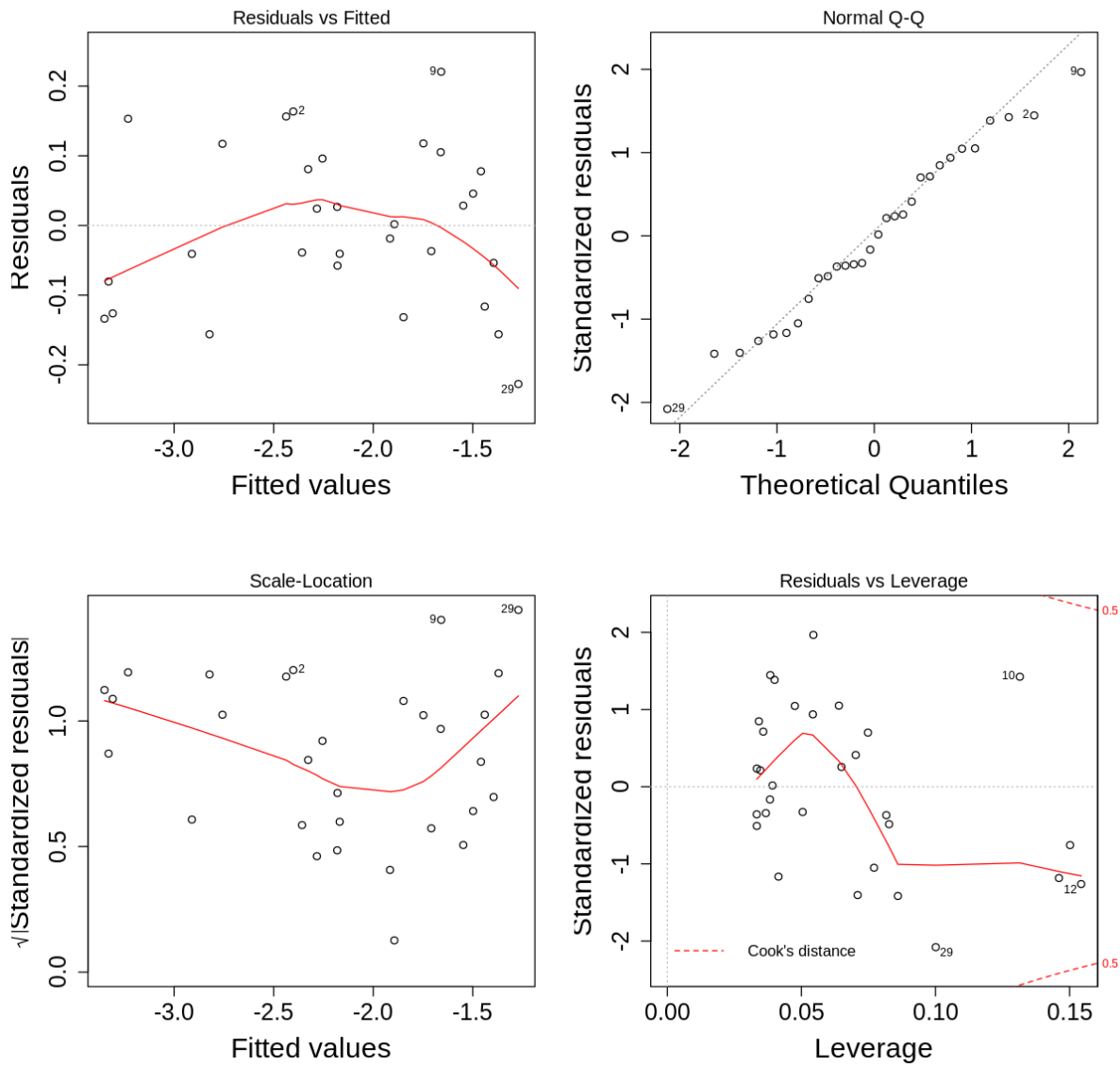
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.1153 on 28 degrees of freedom
Multiple R-squared: 0.9693, Adjusted R-squared: 0.9682
F-statistic: 883.8 on 1 and 28 DF, p-value: < 2.2e-16



```
1 It is much better than the previous one! R^2 reach to 0.97!
```

```
In [36]: 1 par(mfrow = c(2,2))
2 plot(model2,cex.axis=1.5,cex.lab=1.8,ps=10)
3 par(mfrow = c(1,1))
```



```
1 The plot shows that varivance is uniform and normality is stasfied
```

```
In [37]: 1 ncvTest(model2)
2 ##### The variances are uniform
```

Non-constant Variance Score Test
Variance formula: ~ fitted.values
Chisquare = 0.03014715, Df = 1, p = 0.86216

```
In [38]: 1 shapiro.test(model2$residuals)
2 ##### specified packages and libraries you could try: 'cex'
```

Shapiro-Wilk normality test

```
data: model2$residuals  
W = 0.97749, p-value = 0.7557
```

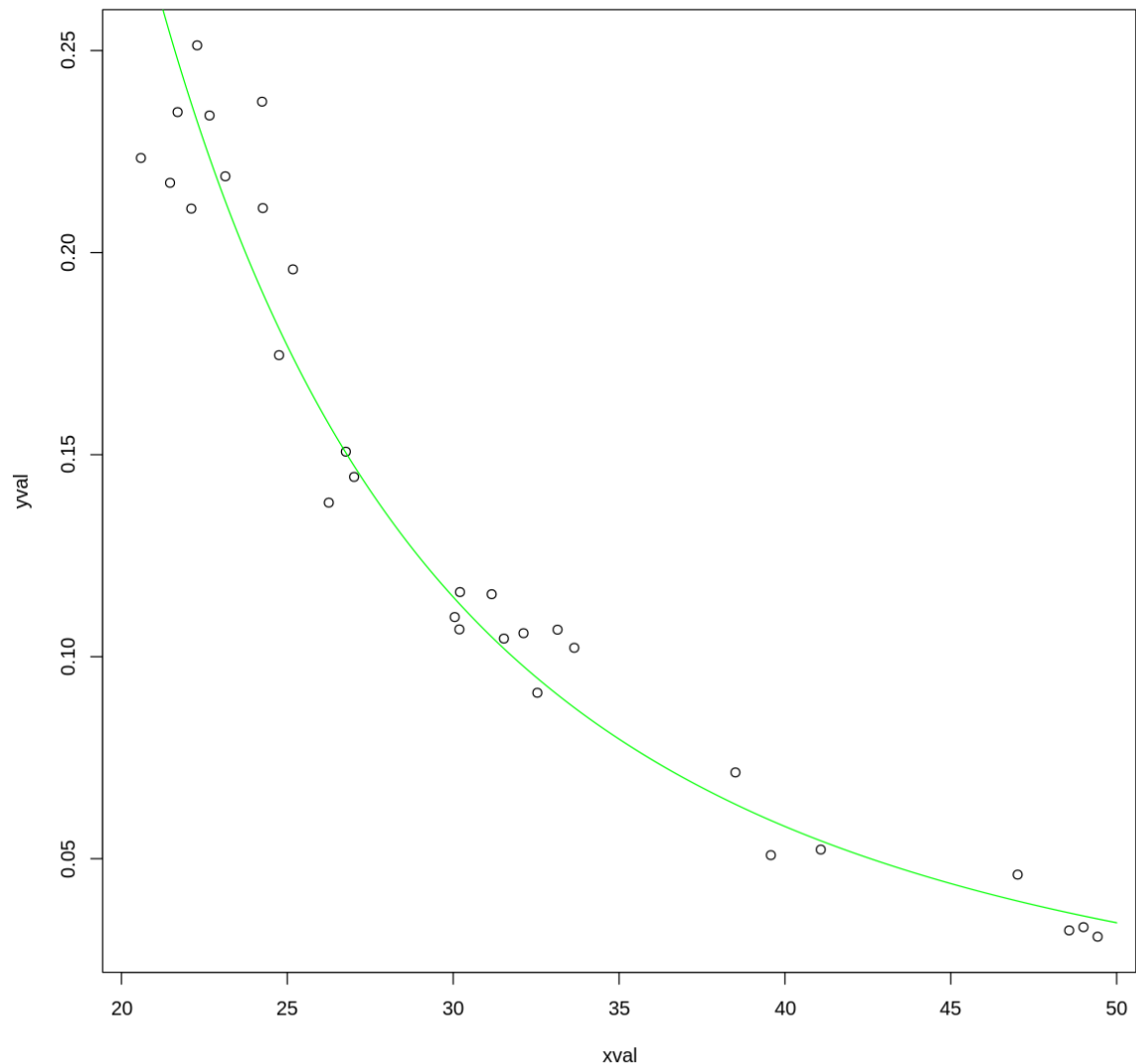
(iv) Write down the final model in terms of the original variables. Draw a scatterplot of y_{val} against x_{val} and add the regression line for the first model and the curve you obtained with the second regression.

```
1 log(yval) = 5.90701-2.37322*log(xval) => yval =  
  exp(5.90701)*xval^(-2.37322)
```

```
In [39]: 1 xCurve <- seq(20, 50, 0.001)  
          2 yCurve <- exp(5.90701)*xCurve^(-2.37322)
```



```
In [40]: 1 plot(yval ~ xval, data = data2)
2 lines(xCurve, yCurve, col = 'green', lty = 1) ## Plot the curve
```



Question 4

For this question use the data set data3.

(i) Read data3 and plot vary as a function of varx. Fit a simple linear regression and add the regression line to the plot. Comment. Obtain a summary for the regression and draw the diagnostic plots. Comment on the results.

```
In [41]: 1 data3 = read.table('data3')
          2 head(data3)
```

A data.frame: 6 × 2

	varx	vary
	<dbl>	<dbl>
1	8.121116	4.105132
2	2.800989	3.313017
3	2.977722	3.320299
4	5.797381	3.755601
5	4.732074	3.677448
6	9.417352	4.366827

```
In [42]: 1 plot(vary ~ varx, data = data3)
        2 modell <- lm(vary ~ varx, data = data3)
        3 abline(modell)
        4 summary(modell)
```

Call:
lm(formula = vary ~ varx, data = data3)

Residuals:

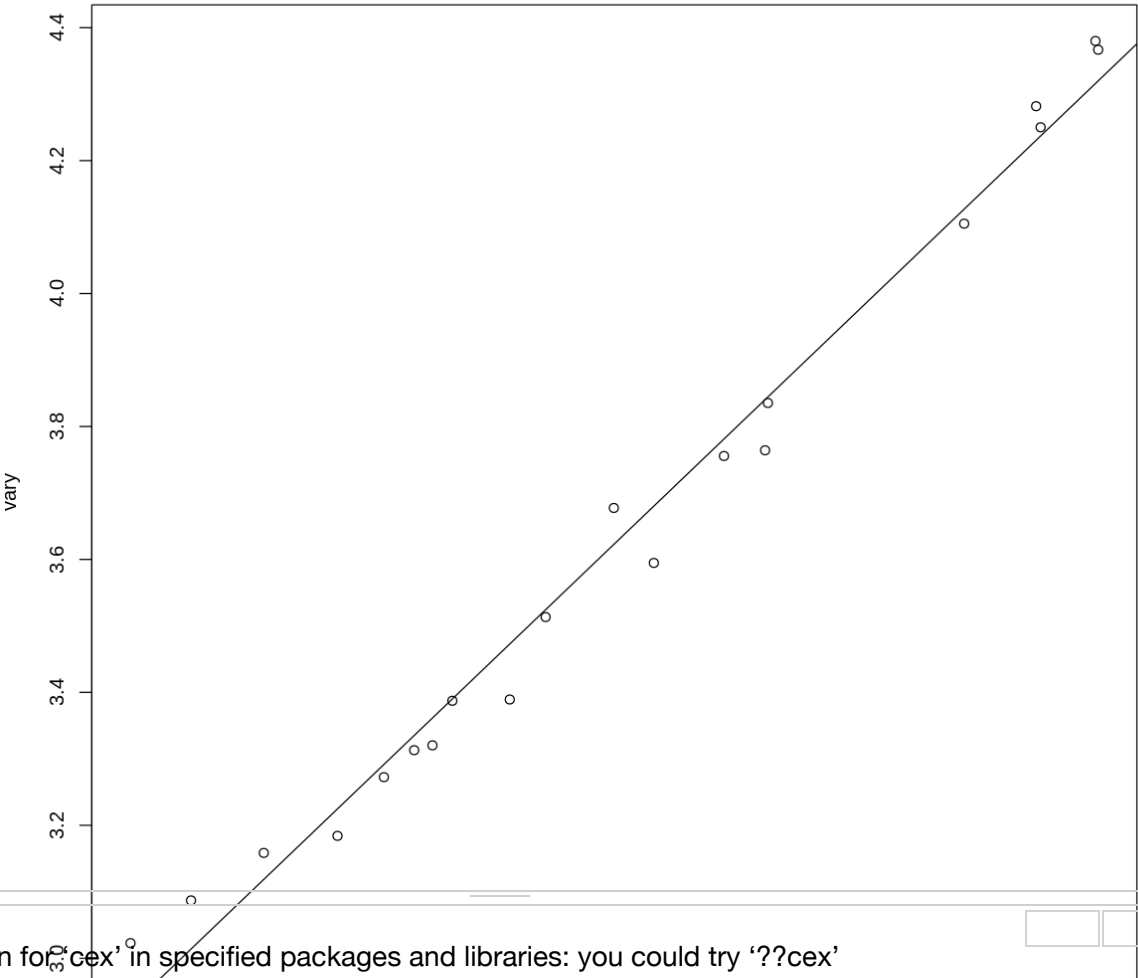
Min	1Q	Median	3Q	Max
-0.08550	-0.02935	-0.01009	0.04794	0.09598

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	2.918281	0.023713	123.06	<2e-16 ***
varx	0.148842	0.004231	35.18	<2e-16 ***

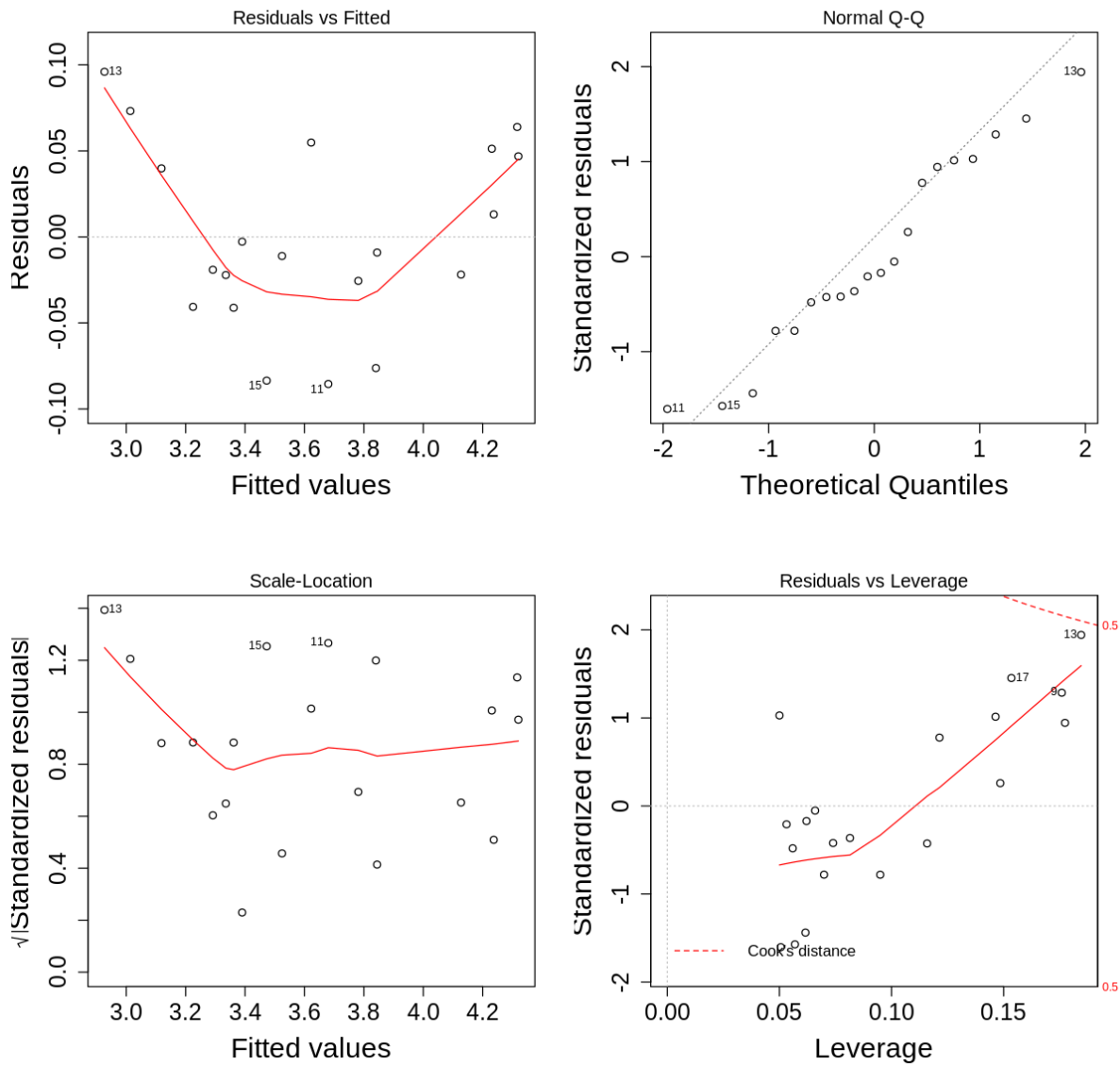
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.05471 on 18 degrees of freedom
Multiple R-squared: 0.9857, Adjusted R-squared: 0.9849
F-statistic: 1238 on 1 and 18 DF, p-value: < 2.2e-16



```
1 The p value is significant and R^2 is 0.0849, it seems very good.
```

```
In [43]: 1 par(mfrow = c(2,2))
2 plot(model1,cex.axis=1.5,cex.lab=1.8,ps=10)
3 par(mfrow = c(1,1))
```



```
1 From residuals plot, we know that there is a quadratic pattern but
2 not for the standardized residuals plot
3 and normality seems not satisfied.
4 We need further check.
5 Point 13 is very close to cook's distance 0.5 and it seems that
6 larger residuals have largerr leverage
```

```
In [44]: 1 ncvTest(model1)
2 ##### The variances are uniform
```

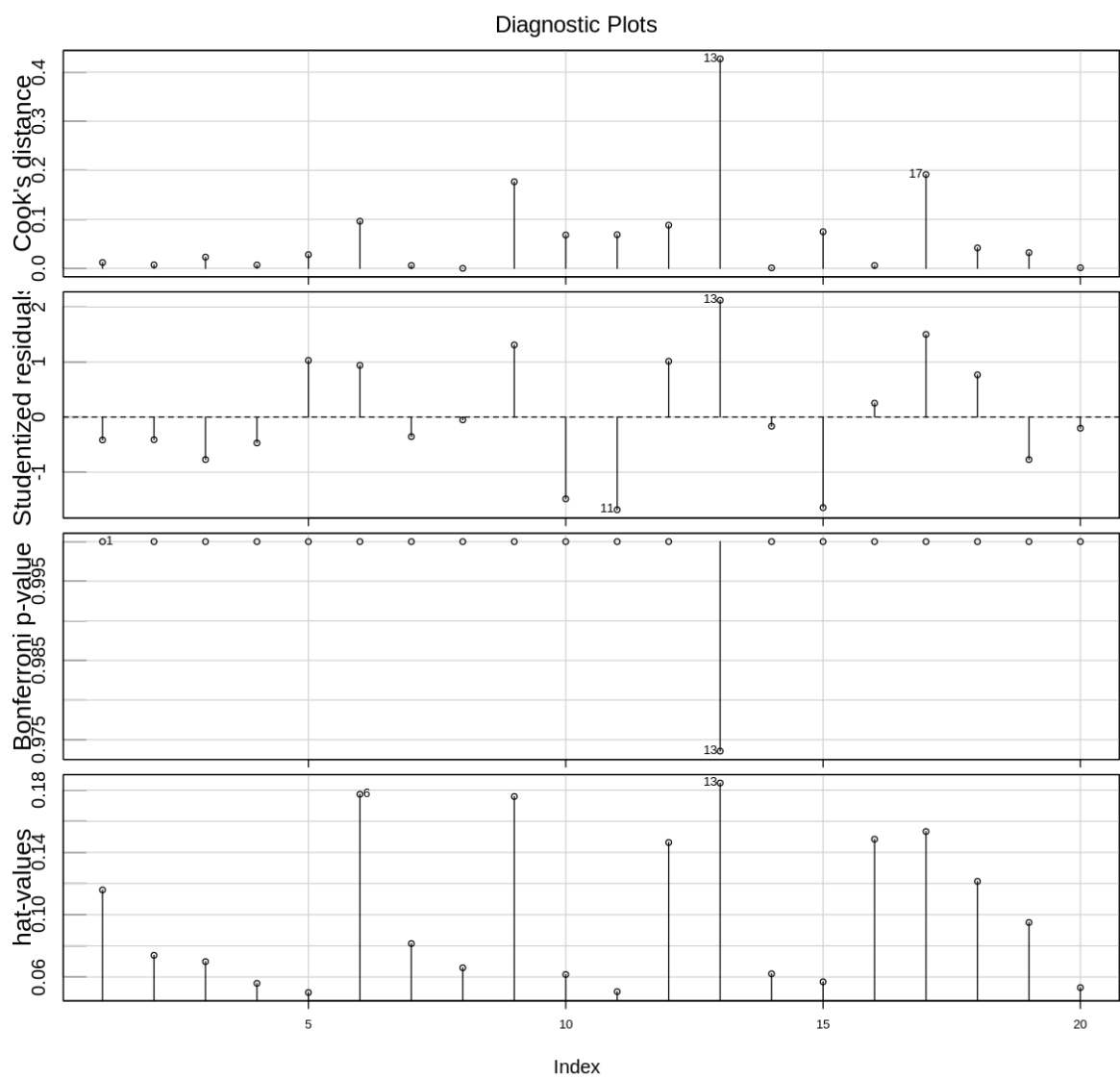
Non-constant Variance Score Test

```
In [45]: 1 shapiro.test(model1$residuals)
2 #####normality assumption is satisfied.
```

Shapiro-Wilk normality test

data: model1\$residuals
W = 0.95481, p-value = 0.446

```
In [46]: 1 influenceIndexPlot(model1,cex.lab=2,cex.axis=1.5)
```

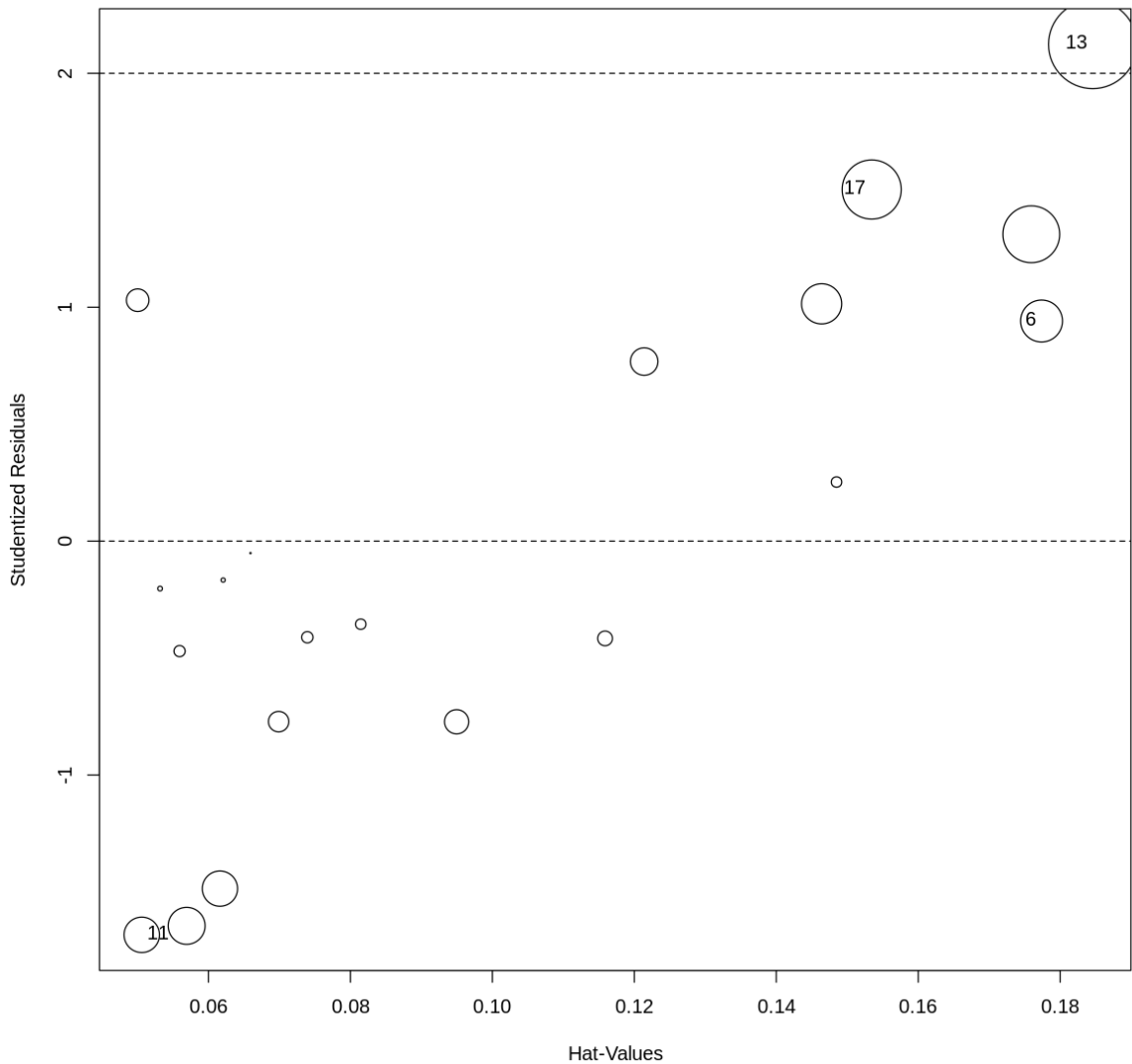


```
In [47]: 1 influencePlot(model1)
```

A data.frame: 4 × 3

StudRes	Hat	CookD
No documentation for 'cex' in specified packages and libraries: you could try '??cex'		
<dbl>	<dbl>	<dbl>

	StudRes	Hat	CookD
	<dbl>	<dbl>	<dbl>
6	0.941004	0.1773770	0.09607734
11	-1.683430	0.0506035	0.06854195



1

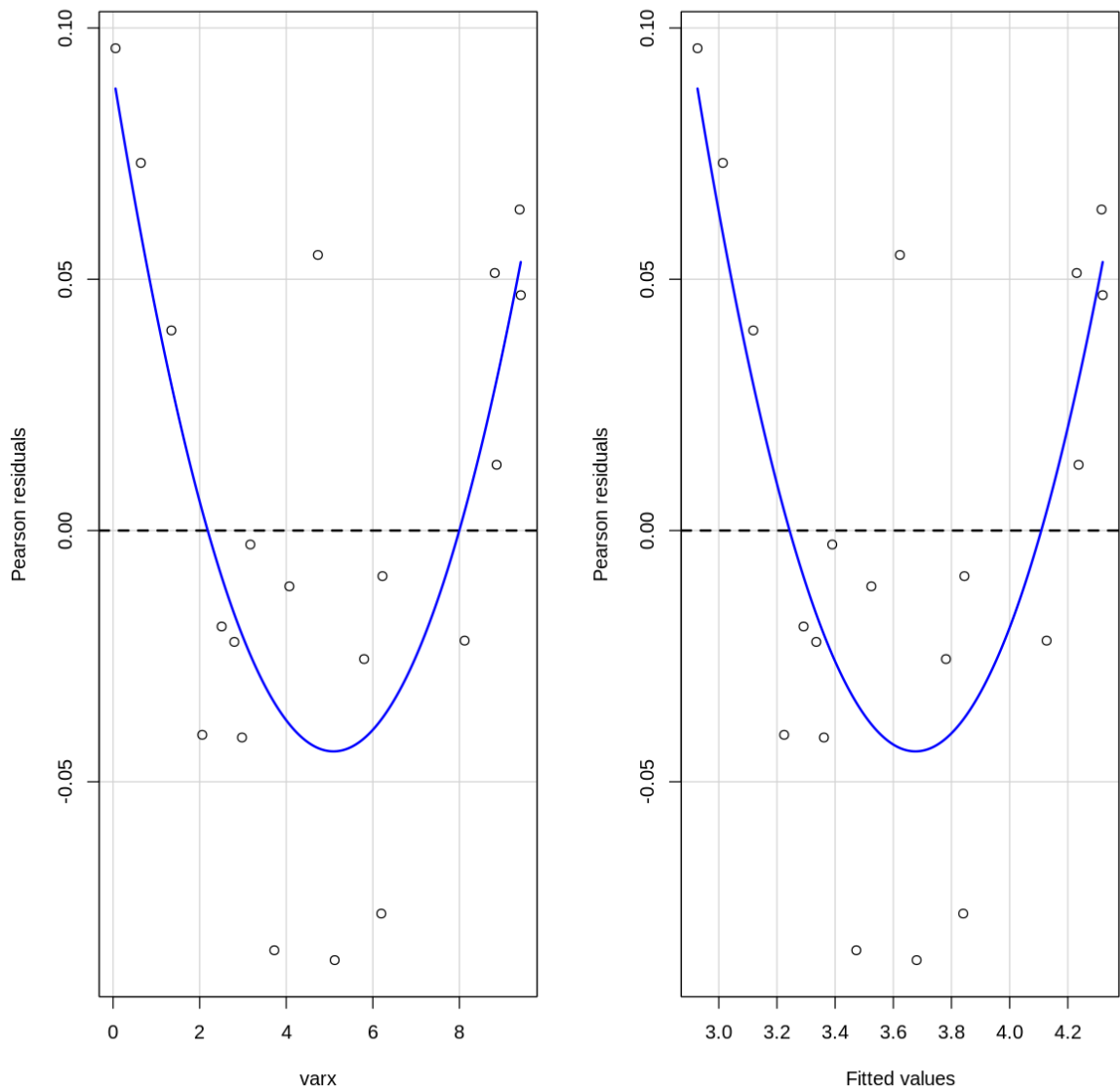
The influence plots we see that 13 have large standard residuals

(ii) Use the function `residualPlots` in package `car` and interpret the test produced by the function. What is your conclusion?

In [48]:

1 residualPlots(model1)

```
Test stat Pr(>|Test stat|)
varx      4.9408      0.0001241 ***
Tinker test 4.9408      7.781e-07 ***
```



1 Very highly to be quadratic term, it means that the model is not sufficient

(iii)Add a quadratic term to the regression model and obtain a summary, draw the diagnostic plots andcomment. Draw a scatterplot of the data and add the lines/curves for both models. Write down youfinal model.

```
In [49]: 1 model2 <- lm(vary ~ poly(varx,2,row=TRUE),
          2           data=data3)
          3 summary(model2)
```

No documentation for 'cex' in specified packages and libraries: you could try '??cex'

```
Call:
lm(formula = vary ~ poly(varx, 2, raw = TRUE), data = data3)

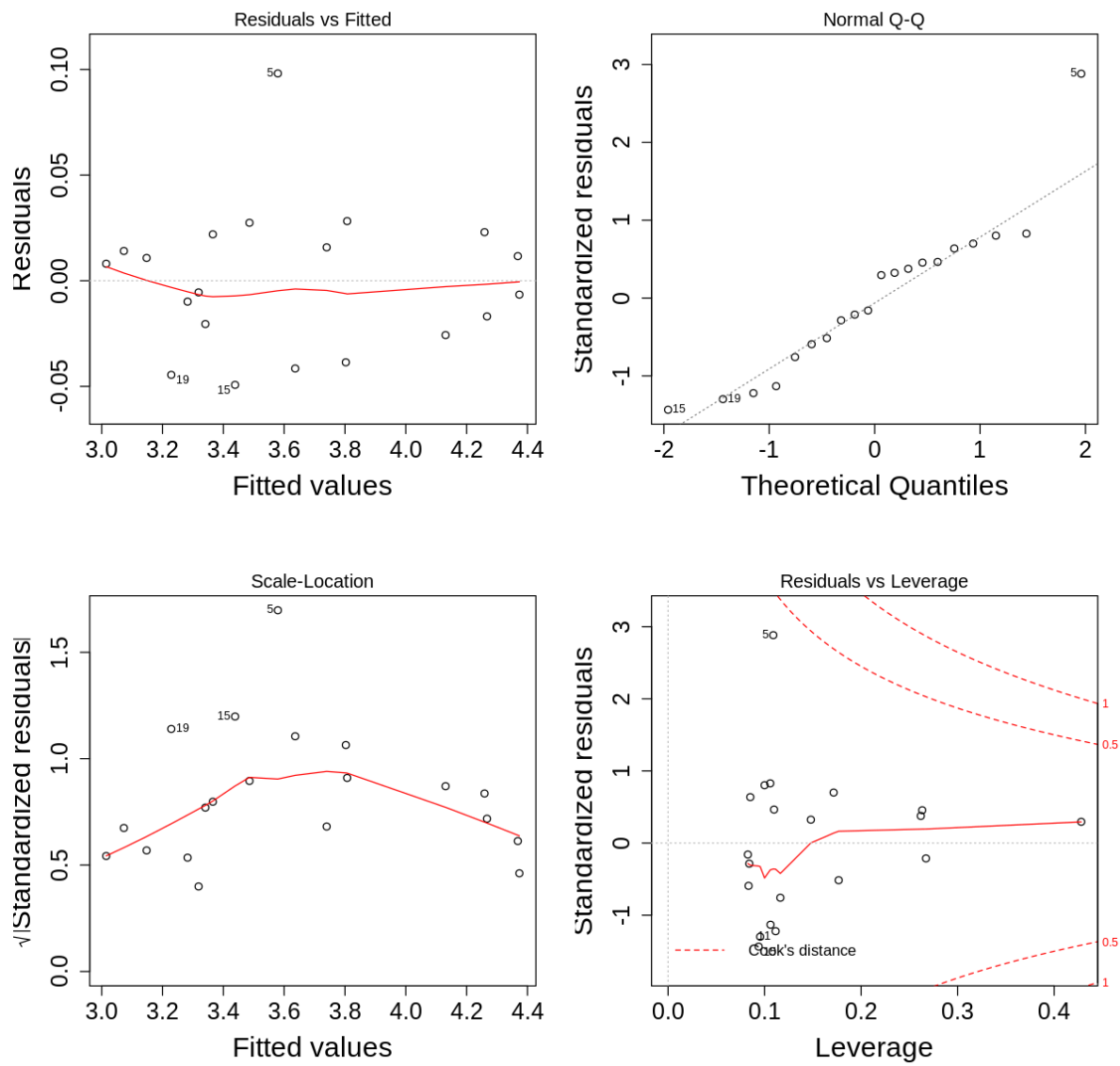
Residuals:
    Min       1Q   Median       3Q      Max
-0.04929 -0.02179  0.00126  0.01733  0.09812

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)      3.009235    0.024152 124.597 < 2e-16 ***
poly(varx, 2, raw = TRUE)1  0.095847    0.011083   8.648 1.24e-07 ***
poly(varx, 2, raw = TRUE)2  0.005204    0.001053   4.941 0.000124 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.03607 on 17 degrees of freedom
Multiple R-squared:  0.9941,    Adjusted R-squared:  0.9934
```

The R^2 square is better than the previous model reach to 0.994, the intercept do not change very much. The slope of varx reduce to less than 0.1 because of adding varx^2


```
In [50]: 1 par(mfrow = c(2,2))
2 plot(model2,cex.axis=1.5,cex.lab=1.8,ps=10)
3 par(mfrow = c(1,1))
```



```
In [51]: 1 ncvTest(model2)
2 shapiro.test(model2$residuals)
```

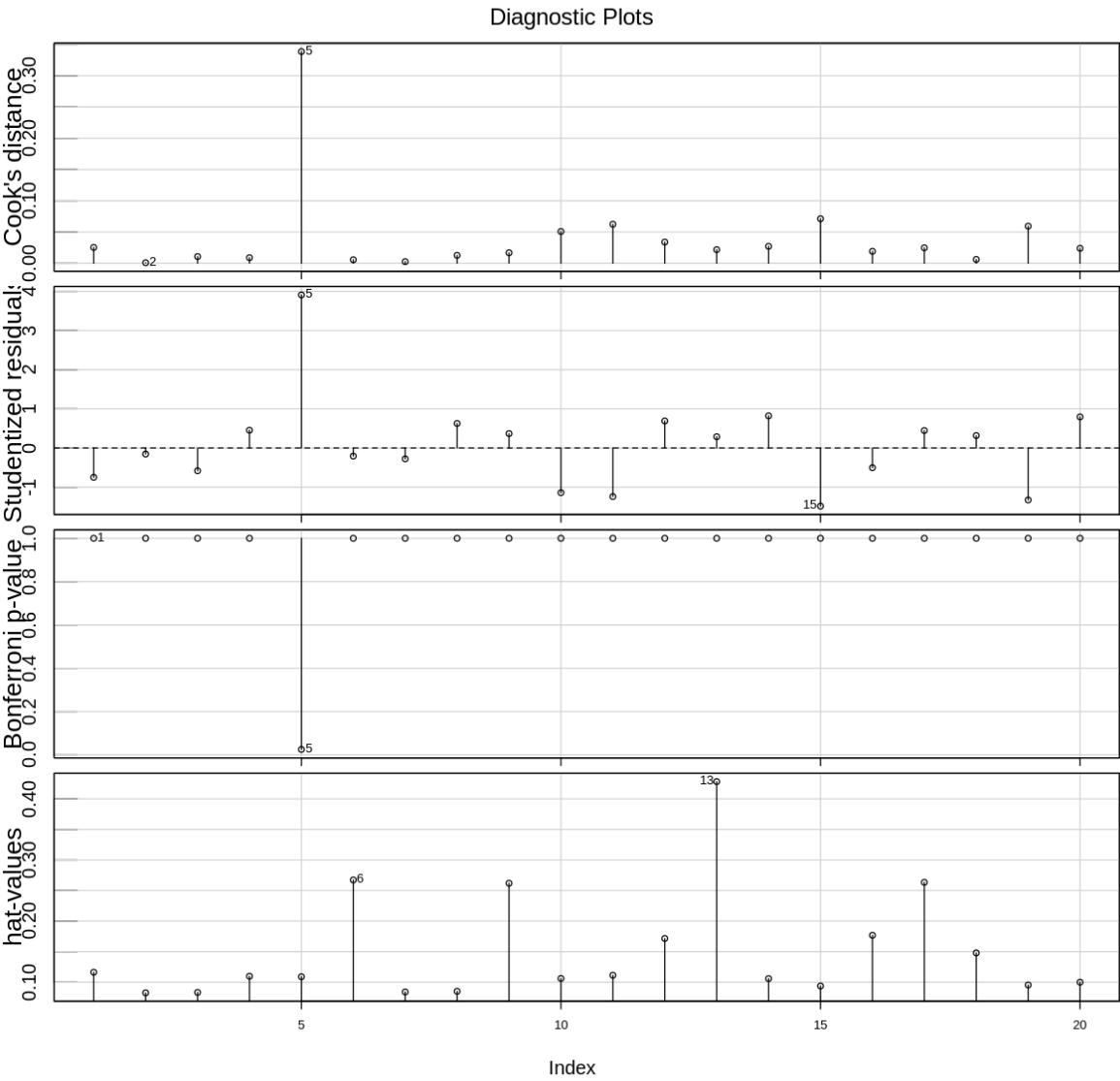
Non-constant Variance Score Test
Variance formula: ~ fitted.values
Chisquare = 0.1112209, Df = 1, p = 0.73876

Shapiro-Wilk normality test

data: model2\$residuals
W = 0.91348, p-value = 0.07427

1	Both homoscedasticity and normality are satisfied. But ther is one outliers 5
---	---

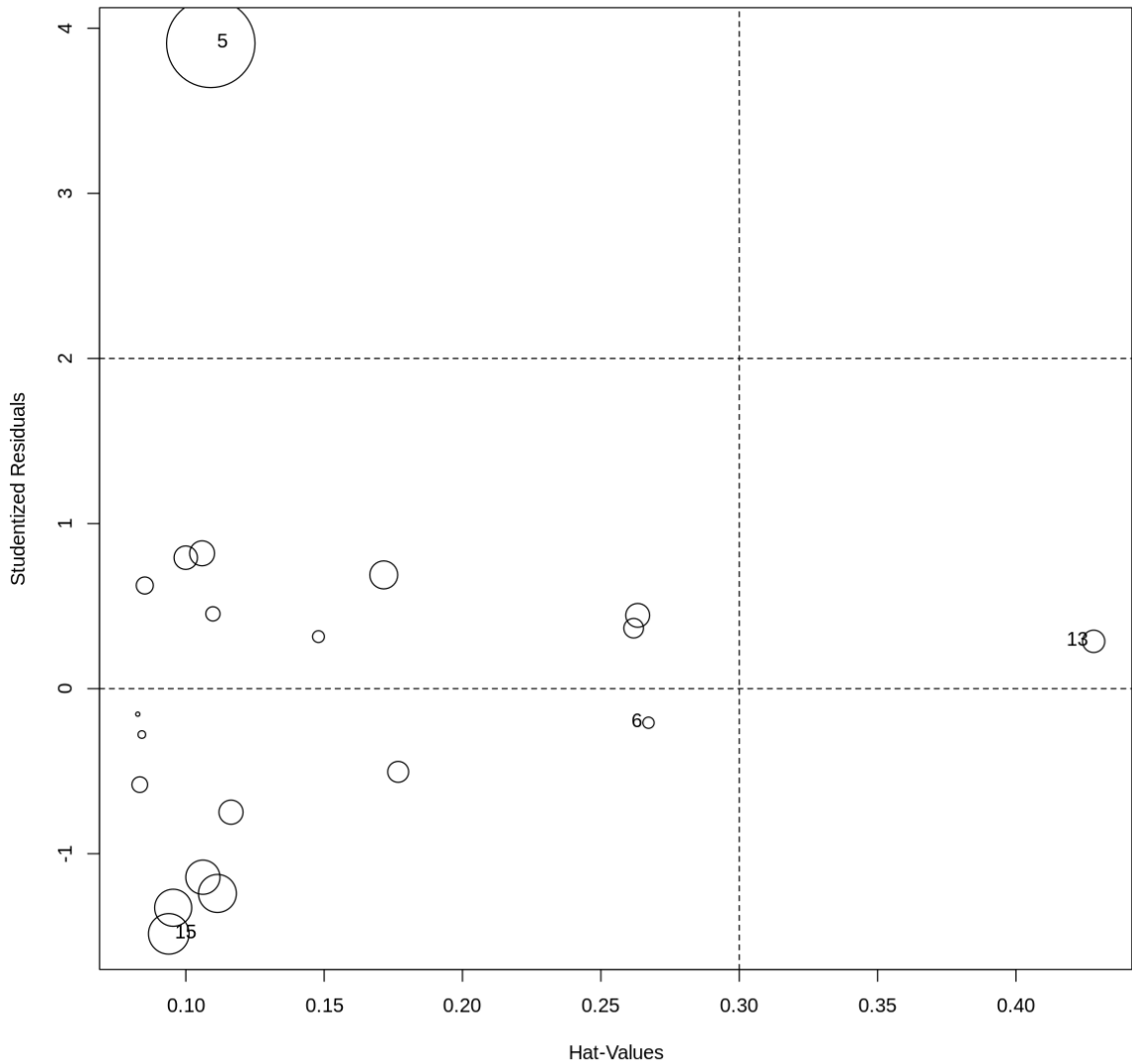
```
In [52]: 1 influenceIndexPlot(model2,cex.lab=2,cex.axis=1.5)
```



```
In [53]: 1 influencePlot(model2)
```

A data.frame: 4 × 3

	StudRes	Hat	CookD
	<dbl>	<dbl>	<dbl>
5	3.9090487	0.10901910	0.33871007
6	-0.2068640	0.26722267	0.00551212
13	0.2864369	0.42808499	0.02163928
15	-1.4854910	0.09384668	0.07113021



```
1 The influence plot show outlier 5
```

```
In [54]: 1 library(alr4)
No documentation for text in specified packages and libraries, you could try '?pex',
```

```
3 partial.residuals = list(cex=2, col=gray(0.5),
4                           lty = 2))
```

Loading required package: effects

Registered S3 methods overwritten by 'lme4':

method	from
cooks.distance.influence.merMod	car
influence.merMod	car
dfbeta.influence.merMod	car
dfbetas.influence.merMod	car

lattice theme set by effectsTheme()
See ?effectsTheme for details.

