

Linear models, assumptions and outputs

UA Summer R Workshop: Week 2

Nicholas M. Caruso
Christina L. Staudhammer

7 June 2016

Create data

Before we start, we'll create a dataframe that includes a continuous response and a continuous explanatory variable using a `function()`. Functions are helpful to create in situations where you might otherwise be re-writing the same code multiple times. Within a function, we can establish the arguments (inside the parentheses), the argument defaults (not necessary, however, you must supply values for all arguments that do not have defaults), the function statements (what it does, within brackets `{}`), and the function values (what the function returns). We can then store the results of our function into a name. Alternatively, we would read in our dataframe using `read.csv()`.

```
library(tidyverse)
library(dplyr)

lmdat.fn <- function(n, alpha=1.5576, beta.var=-2.6513,
                     sd.noise=4.5){
  # n: Number of individuals
  # alpha, beta.var: coefficients (intercept and slope)

  # Generate explanatory variable
  variable <- rnorm(n, mean=0, sd=1)

  # Signal: Build up systematic part of the LM
  expected.response <- alpha + beta.var*variable

  # Noise: generate error (normally distributed) around expected response
  response <- rnorm(n, mean=expected.response, sd=sd.noise)
  response2 <- expected.response + rnorm(n, mean=0, sd=sd.noise)

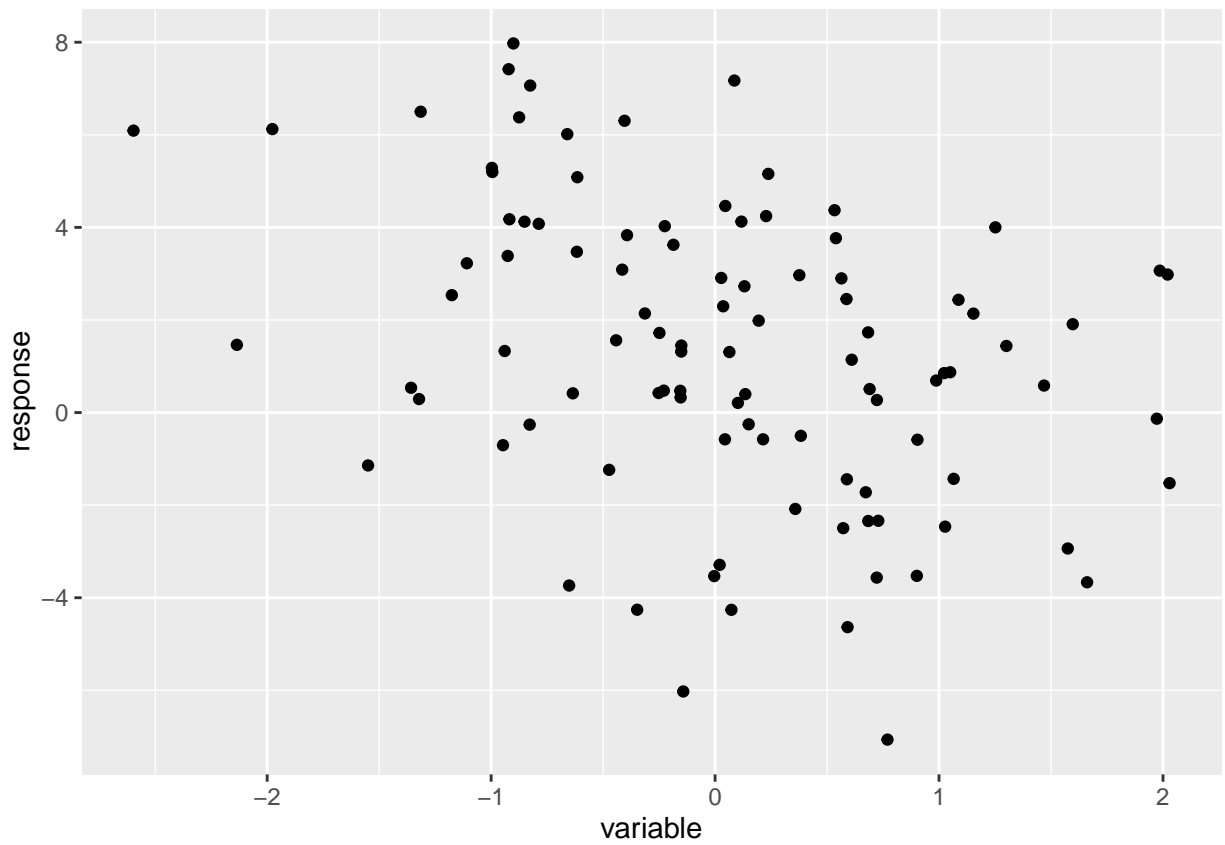
  # Return dataframe
  return(data_frame(variable=variable, response=response,
                    expect.response=expected.response,
                    response2=response2))
}

set.seed(8675309)
dat <- lmdat.fn(n=100, alpha=1.5576, beta.var=-1.5, sd.noise=3)

# write.csv(dat, 'dat_2june2016.csv')
# dat <- read.csv('dat_2june2016.csv')

library(ggplot2)

ggplot(dat, aes(variable, response)) +
  geom_point()
```



Linear model

We use the `lm()` function to call a linear model. The main arguments in this function are the formula, which is written as the response followed by a tilde (\sim), then the explanatory variable(s), and then the data that we are using (typically a data frame). This is typically stored into a name which can be used to inspect the model further.

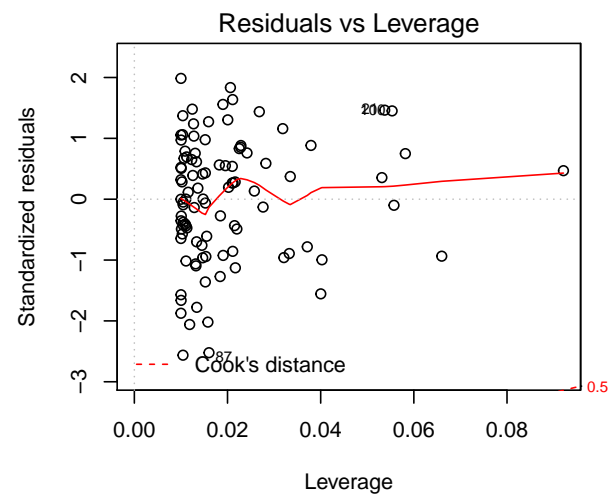
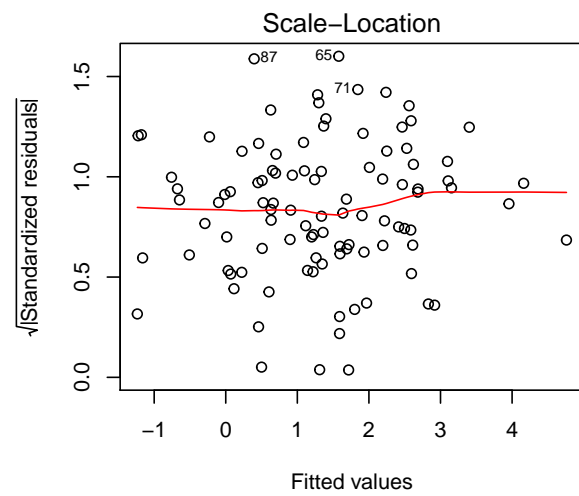
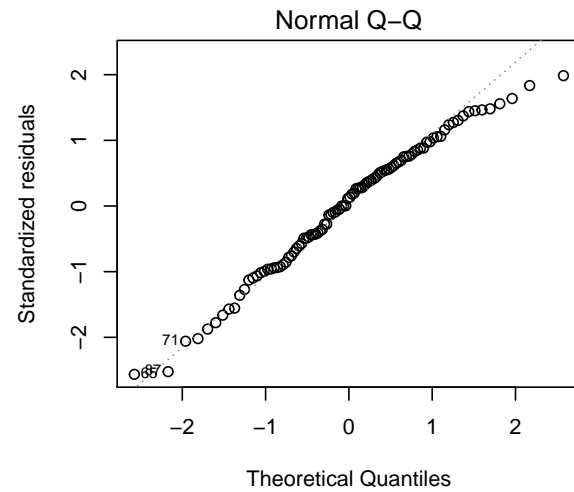
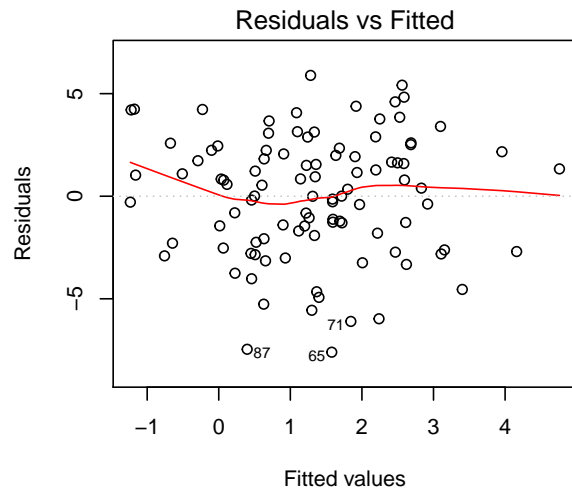
```
mod.lm <- lm(response~variable, data=dat) # store the model
```

Linear model assumptions

- 1) The model is correctly specified (linearity)
- 2) Each $y_i|x_i$ is normally distributed (given is important!)
- 3) The y_i are homoscedastic
- 4) Any observation of y_i is independent of all other y_i
- 5) For a given x_i , the observations of y_i are randomly selected
- 6) The values of x_i are fixed and/or measured without error

Assumptions 1-3 can be assessed using R, while assumptions 4-6 should be addressed in the study design or by using an alternative model; thus we will focus on the first three assumptions. These can be assessed graphically in R using the `plot()` function. This produces 4 figures, if you just use `plot()`, R will cycle through the 4 figures, asking you to hit return to scroll through them. Alternatively you can adjust the `mfrow` graphic parameter `par()` so that we create a plot with 4 panels (2 rows and 2 columns) to view all of them at once, but we will go over each plot individually.

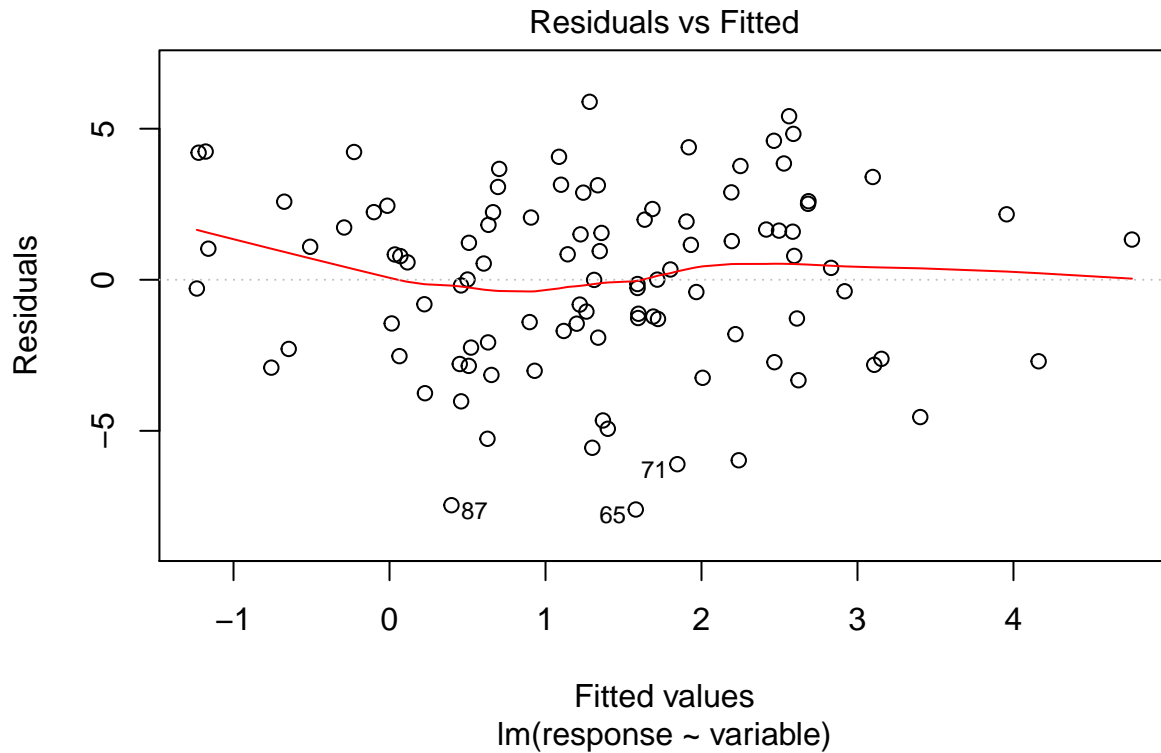
```
par(mfrow=c(2,2))
plot(mod.lm)
```



Plot 1: Linearity and homoscedasticity

The first plot of the residuals (error, y-axis) vs the fitted values (\hat{y} , x-axis) assess both the linearity and homoscedasticity assumptions. If our data are linear, the red line will be relatively flat. We can also use this plot to assess if the residual standard error is constant (Homoscedasticity), in which the data points would be scattered with no obvious pattern. Our model meets both of the assumptions based on this figure.

```
plot(mod.lm, 1)
```

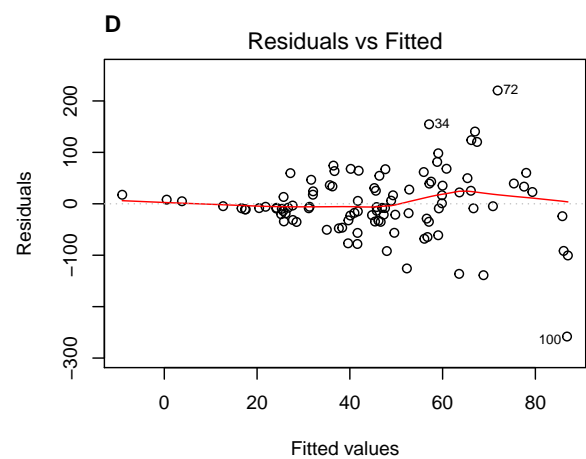
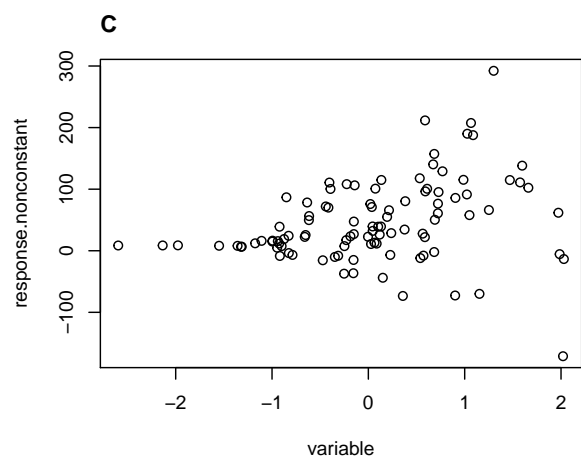
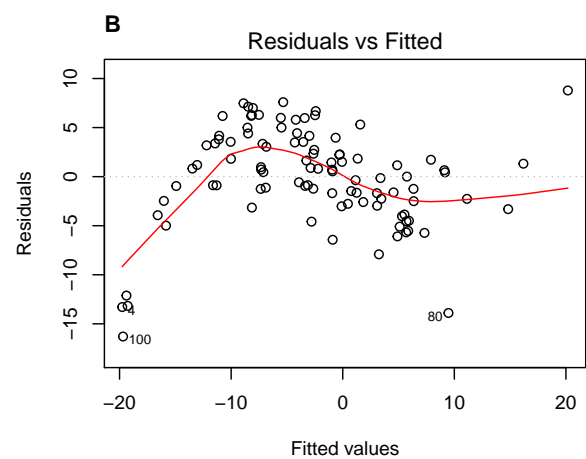
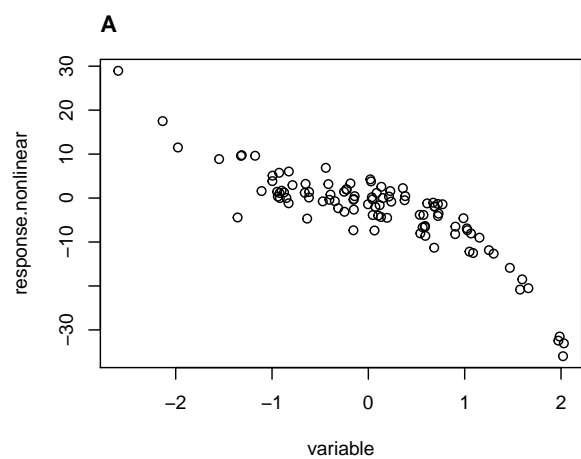


Violations in plot 1

Here is an example of two responses that do not meet these assumptions. For the first response, I created a non linear response, in which the response is a polynomial function (3rd order) of our explanatory variable with added normally distributed error. The second response was created with a trend in the normally distributed error, in which size and standard deviation of the error are positively related to the explanatory variable. The top graphics show the data (A) and a violation of the assumption of linearity (B). The bottom graphics show the data (C) and a violation of the assumption of homoscedasticity (D).

```
dat <- dat %>%
  mutate(response.nonlinear=-3.453*variable + -2.345*(variable^2) + -1.97*(variable^3) +
          rnorm(length(variable), 0, 3),
          response.nonconstant=-3.453*variable + rnorm(length(variable),
                                                         rank(variable), rank(variable)))

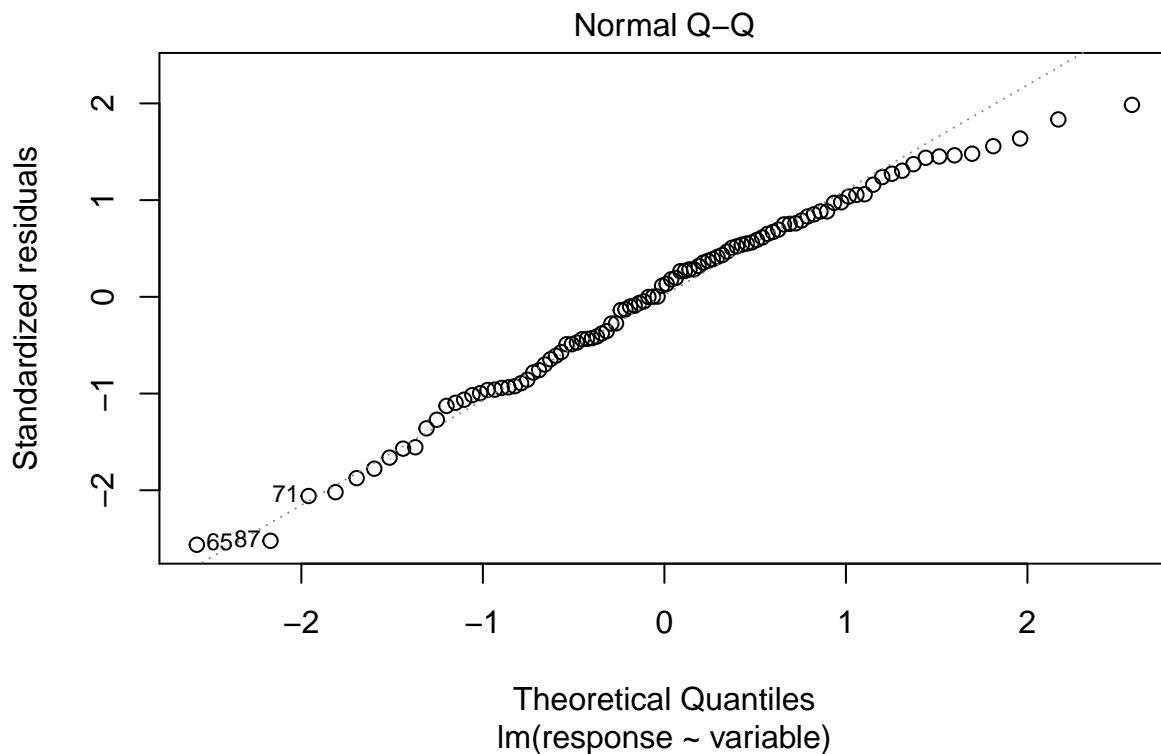
par(mfrow=c(2,2))
plot(response.nonlinear~variable, dat)
mtext('A', side=3, line=1, adj=0, font=2)
plot(lm(response.nonlinear~variable, dat), 1) # non linear, but constant variance
mtext('B', side=3, line=1, adj=0, font=2)
plot(response.nonconstant~variable, dat)
mtext('C', side=3, line=1, adj=0, font=2)
plot(lm(response.nonconstant~variable, dat), 1) # variance is not constant, but is linear
mtext('D', side=3, line=1, adj=0, font=2)
```



Plot 2: Errors are normally distributed

The second plot assesses whether the errors are normally distributed. This can also be tested using a Shapiro Wilks test (null hypothesis is that the errors are normally distributed), a common mistake is to test the normality of the response rather than of the errors (residuals). We can call the residuals of the model using the `residuals()` function. Additionally, most models are robust to moderate violations of the normality assumption, often this test would have to show extreme deviations from normality (i.e., $p \ll 0.05$) to preclude use of a linear model. Alternatively, we can use a qqplot to graphically inspect our model. For the graph, the data should follow the dotted line.

```
plot(mod.lm, 2)
```



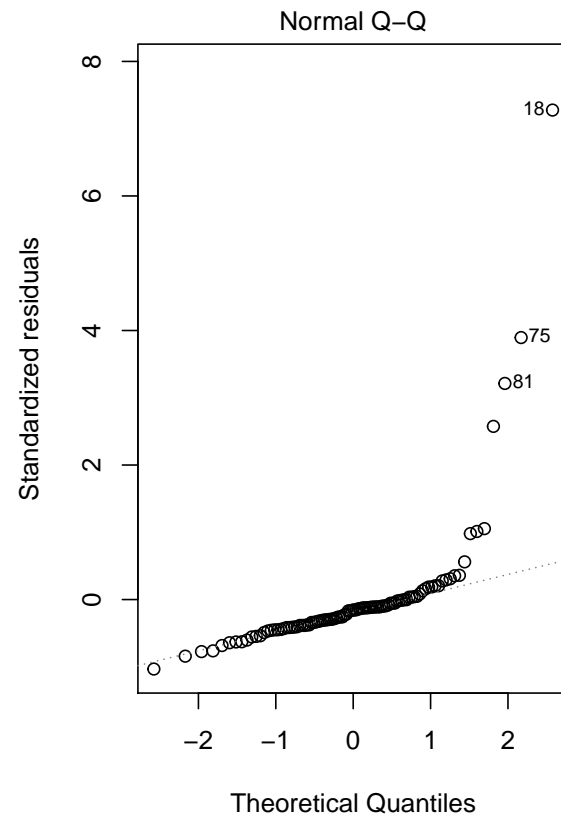
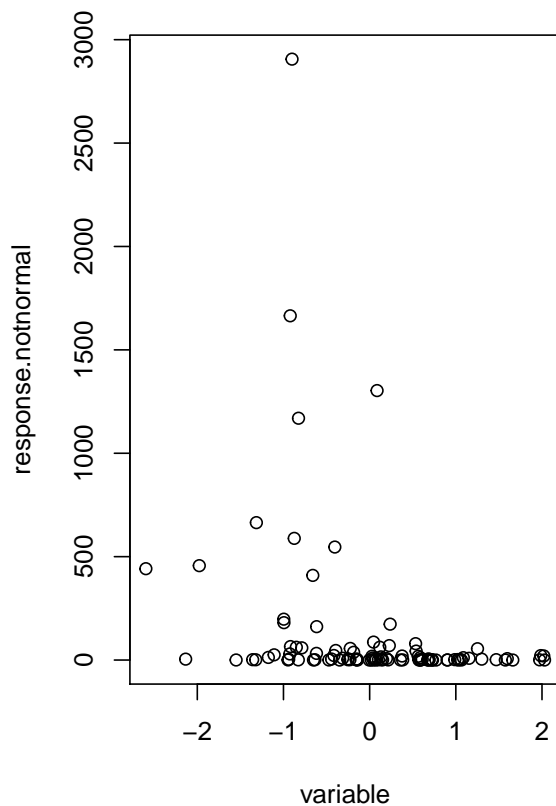
```
shapiro.test(residuals(mod.lm))
```

```
##  
##  Shapiro-Wilk normality test  
##  
## data:  residuals(mod.lm)  
## W = 0.9826, p-value = 0.2107
```


Violations in plot 2

Here is an example of a response that violates the normality assumption.

```
dat <- dat %>%  
  mutate(response.notnormal=exp(response))  
  
shapiro.test(residuals(lm(response.notnormal~variable, dat)))  
  
##  
## Shapiro-Wilk normality test  
##  
## data: residuals(lm(response.notnormal ~ variable, dat))  
## W = 0.48663, p-value < 2.2e-16  
  
par(mfrow=c(1,2))  
plot(response.notnormal~variable, dat)  
plot(lm(response.notnormal~variable, dat), 2)
```



Shapiro Wilks test vs. graphic diagnostics

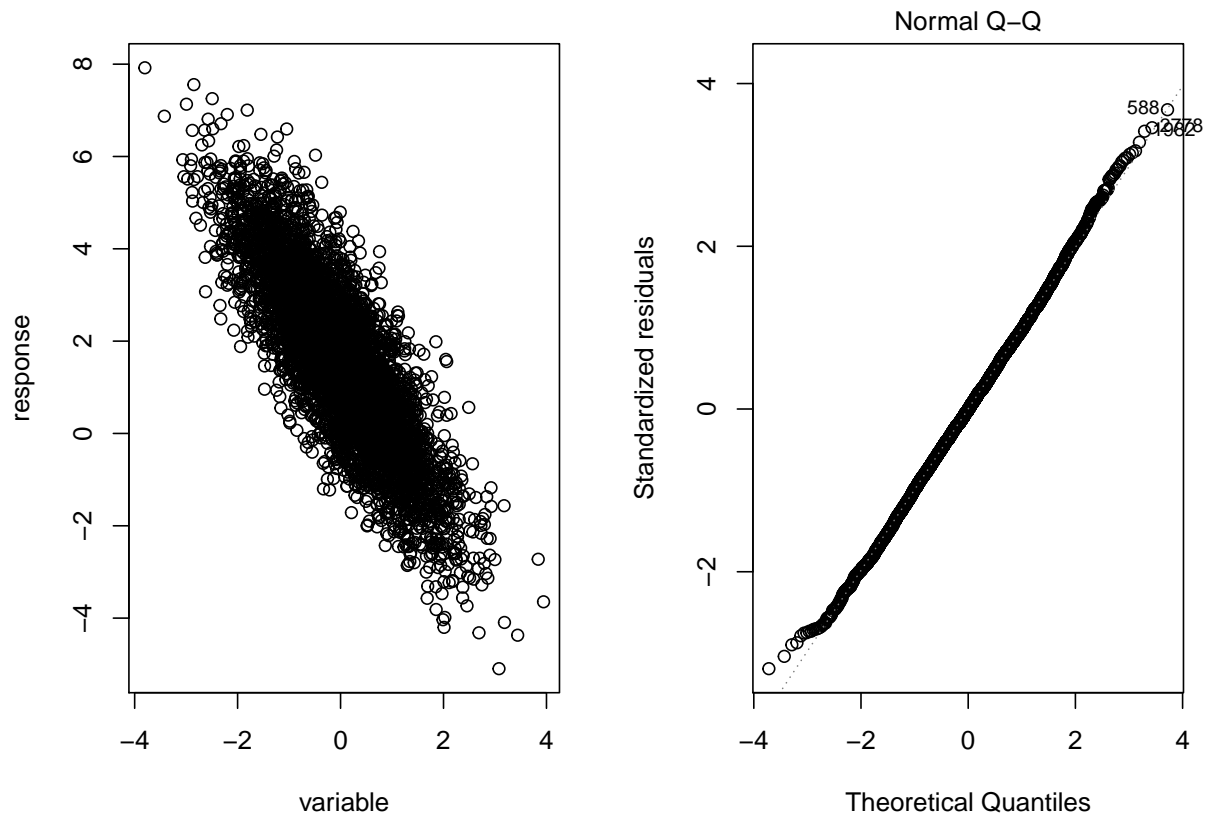
A good reason to inspect the model using plots rather than Shapiro Wilks test is that the Shapiro Wilks test is sensitive to large sample sizes and will often show non-normally distributed errors with larger datasets (Type I error). Obviously, our data was created using normally distributed errors, and the diagnostic plot shows no deviations from normality in the errors.

```
set.seed(7375)
dat.large <- lmdat.fn(n=5000, alpha=1.5576, beta.var=-1.5, sd.noise=1)

shapiro.test(residuals(lm(response~variable, dat.large)))

##
##  Shapiro-Wilk normality test
##
## data:  residuals(lm(response ~ variable, dat.large))
## W = 0.99922, p-value = 0.02476

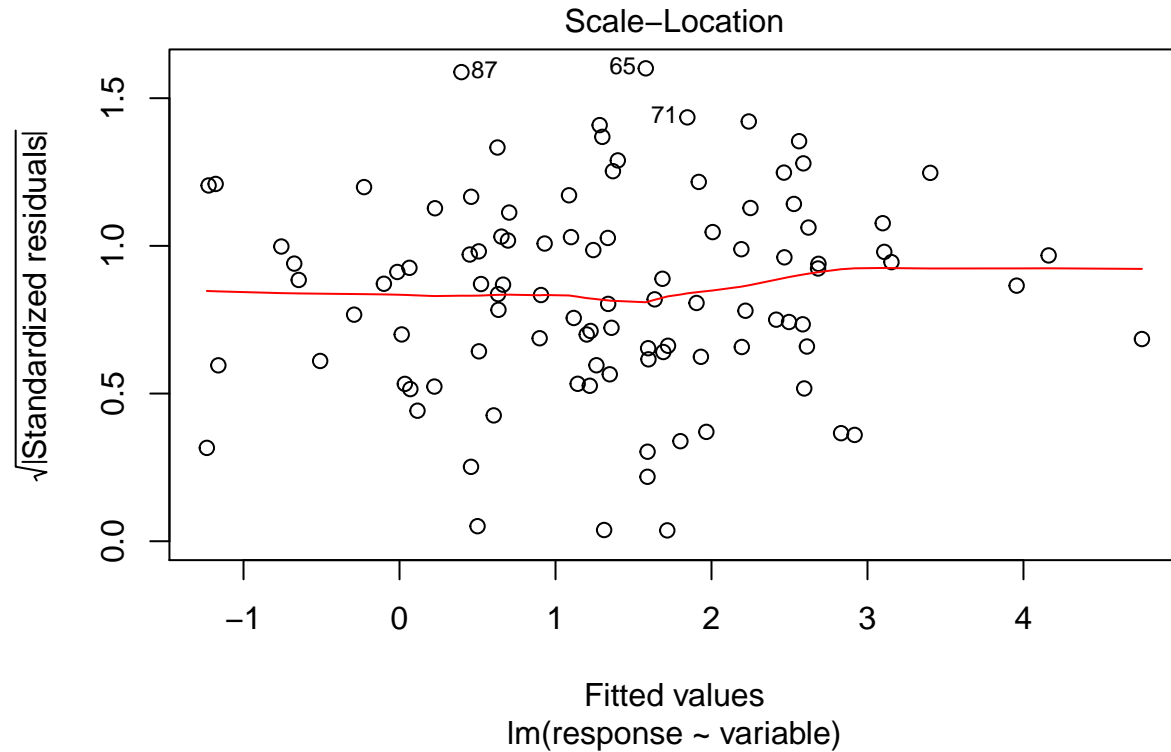
par(mfrow=c(1,2))
plot(response~variable, dat.large)
plot(lm(response~variable, dat.large), 2)
```



Plot 3: Linearity and homoscedasticity (standardized residuals)

The third diagnostic plot can also assess non-linearity and homoscedasticity. The third plot is very similar to the first, except the scale of the residuals is standardized such that they are all positive (i.e., large negative residuals become large positive residuals).

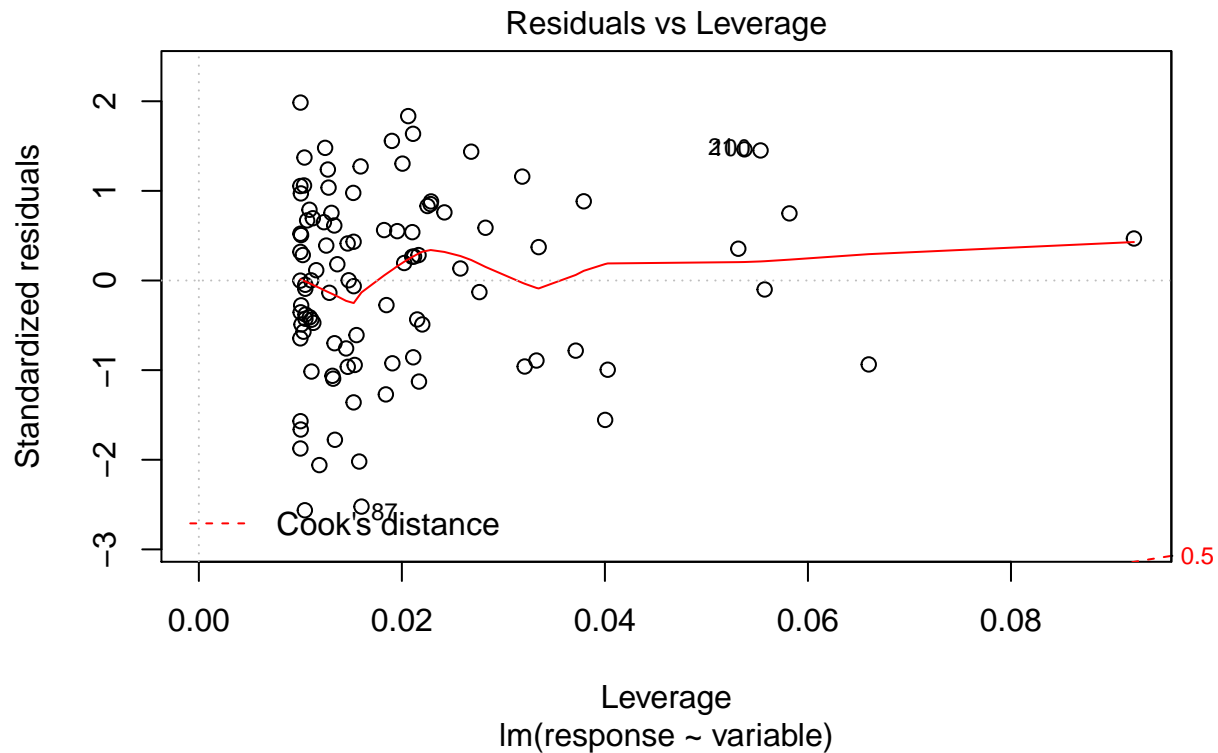
```
plot(mod.lm, 3)
```



Plot 4: Leverage and Cook's distance

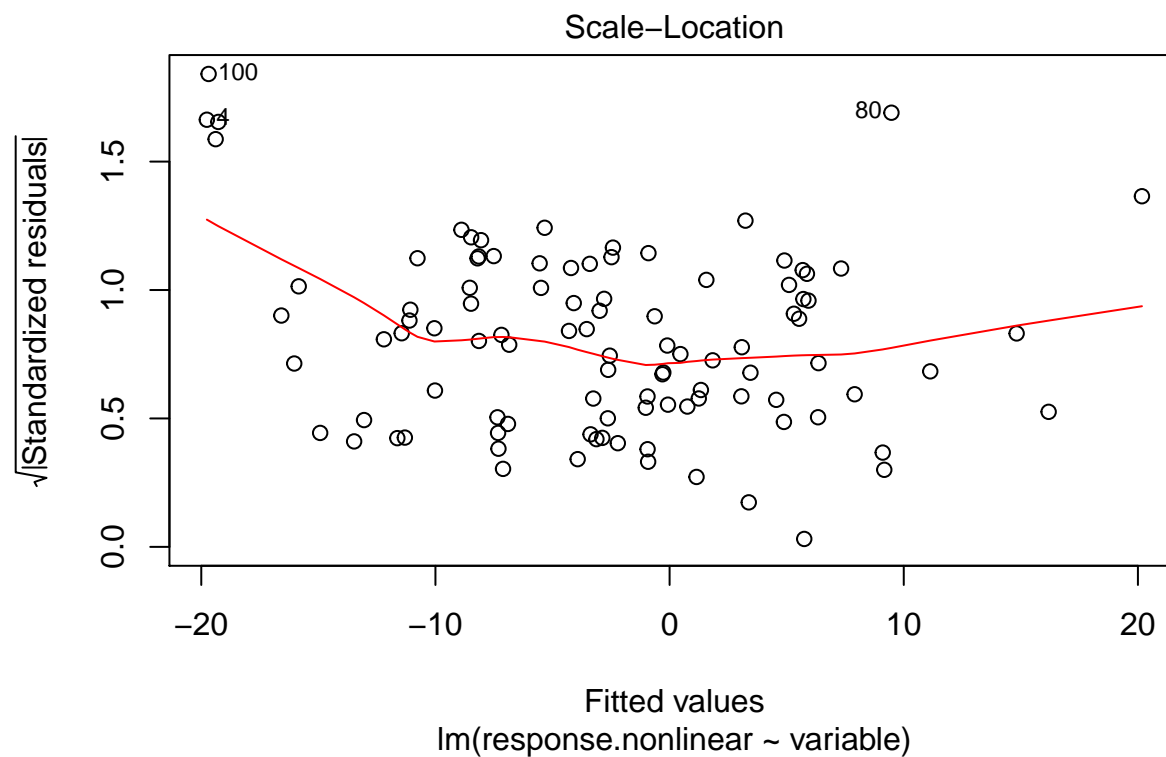
The last plot (actually 5th of the diagnostic plots, there's six total), shows the standardized residuals against leverage. Leverage is a measure of how much impact that a given y_i has on \hat{y}_i , such that the further a given x_i is from the \bar{x} , the more influence it has on the estimate of regression coefficients (higher leverage). For this plot, the red smoothed line should stay close to the gray dotted line and no points should have large Cook's distance (all should be < 0.5 and inside the dotted red lines). Cook's distance is a measure of the influence of a data point, typically these high influence data points are those with large residuals and/or high leverage.

```
plot(mod.lm, 5)
```

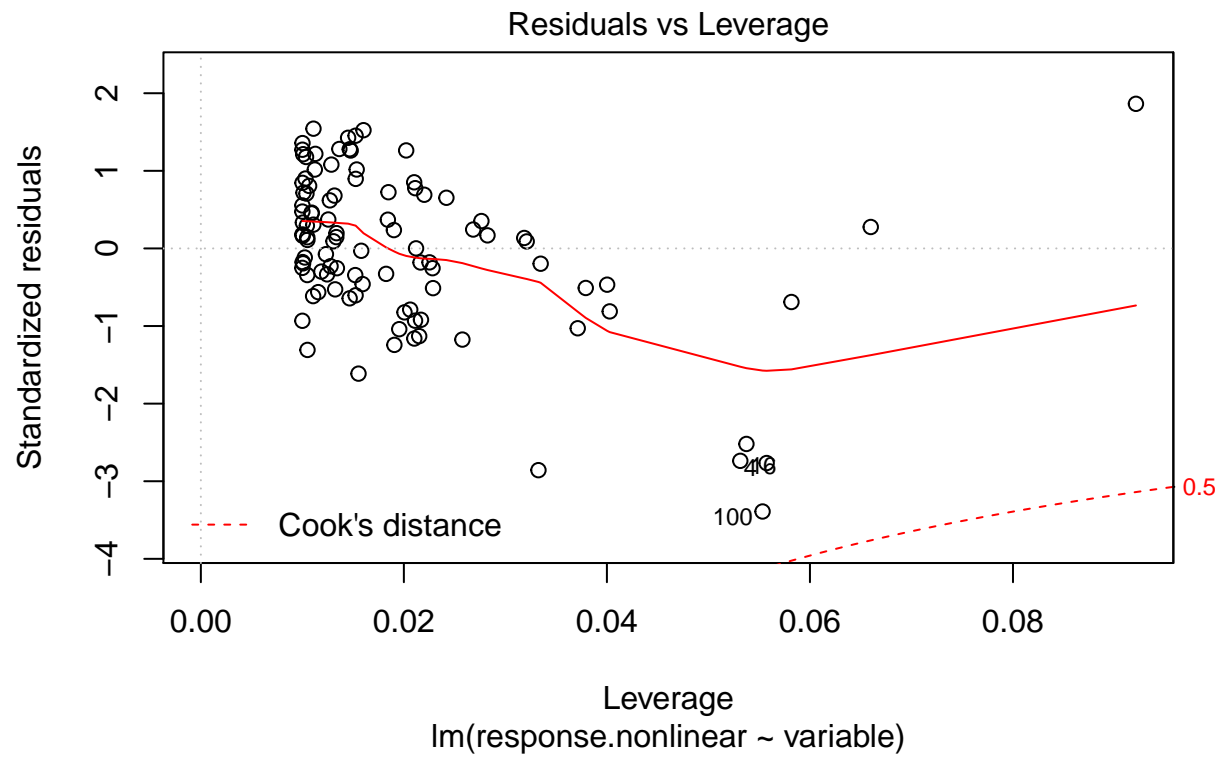


Next, we will show examples where data show violations in plots 3 and 4. We will revisit our nonlinear and non normal data that we created before.

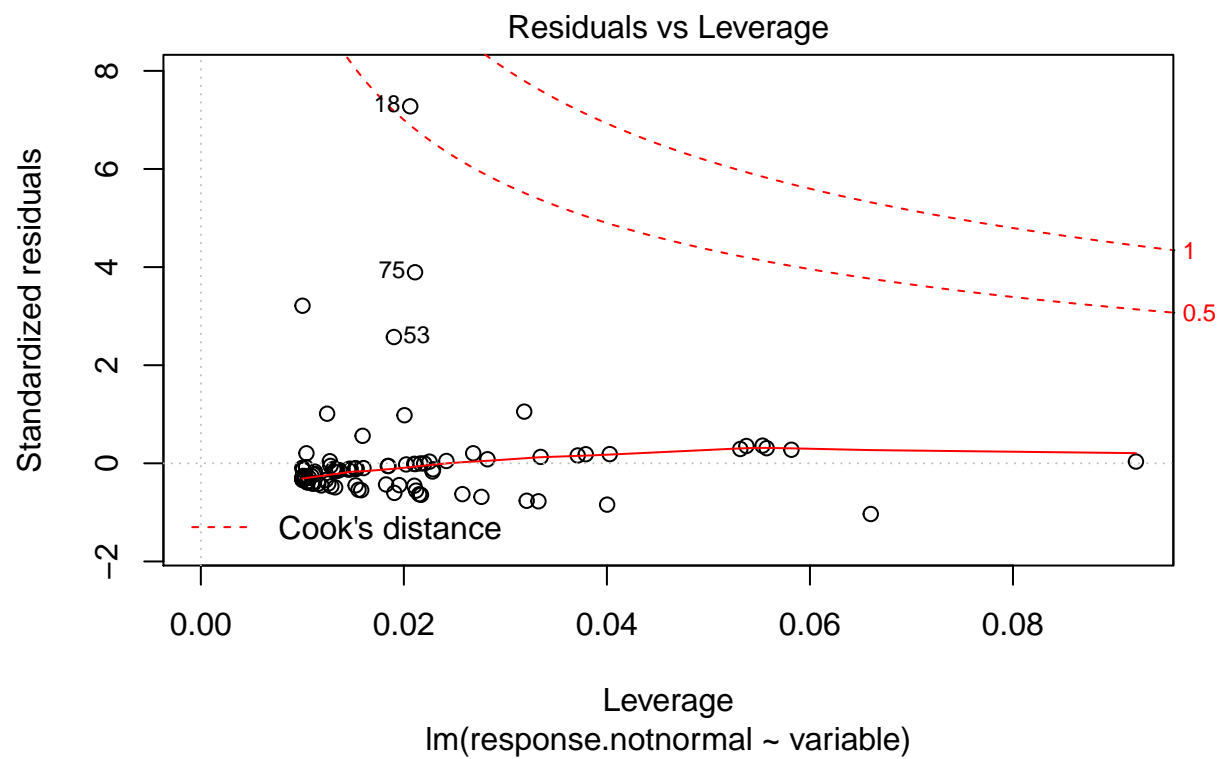
```
# a non linear relationship violates this diagnostic plot  
plot(lm(response.nonlinear~variable, dat), 3)
```



```
# a non linear relationship violates this diagnostic plot  
plot(lm(response.nonlinear~variable, dat), 5)
```



```
# a non normal relationship has points with large Cook's distance  
plot(lm(response.notnormal~variable, dat), 5)
```



Summarise model results

Ok now that we know our data meet the assumptions of the model, we can look at the results. We will use `summary` and `anova` to extract the results of this model.

`summary` vs. `anova`

`summary()` and `anova()` test two different hypotheses for linear models, `summary()` (using t -test statistics) tests if a given model coefficient $\beta = 0$ while `anova()` (using F -test) tests if a given variable reduces the residual sum of squares significantly. For a linear model with one variable, these two outputs yield identical results (p values), but this won't be true for more complicated models.

```
mod.lm # calling just the object shows model formula and coefficients
```

```
##
## Call:
## lm(formula = response ~ variable, data = dat)
##
## Coefficients:
## (Intercept)      variable
##      1.394      -1.296
```

```
summary(mod.lm)
```

```
##
## Call:
## lm(formula = response ~ variable, data = dat)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -7.604 -2.116  0.367  2.184  5.889
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   1.3945     0.2987   4.669 9.62e-06 ***
## variable     -1.2958     0.3225  -4.017 0.000115 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.982 on 98 degrees of freedom
## Multiple R-squared:  0.1414, Adjusted R-squared:  0.1326
## F-statistic: 16.14 on 1 and 98 DF,  p-value: 0.0001155
```

```
anova(mod.lm)
```

```
## Analysis of Variance Table
##
## Response: response
##           Df Sum Sq Mean Sq F value    Pr(>F)
## variable   1 143.51  143.508    16.14 0.0001155 ***
## Residuals 98  871.39    8.892
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```


lm vs. aov

Alternatively, we could have used the `aov()` function to fit our model. This function, however, would only allow us to test the hypothesis that the added variable reduces the residual sum of squares (F -test, analysis of variance), regardless of using `summary()` or `anova()`.

```
mod.aov <- aov(response~variable, dat)
summary(mod.aov)
```

```
##           Df Sum Sq Mean Sq F value    Pr(>F)
## variable     1  143.5   143.51   16.14 0.000115 ***
## Residuals   98  871.4     8.89
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
anova(mod.lm)
```

```
## Analysis of Variance Table
##
## Response: response
##           Df Sum Sq Mean Sq F value    Pr(>F)
## variable     1 143.51  143.508   16.14 0.0001155 ***
## Residuals   98 871.39    8.892
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

broom package

We can further inspect our model using the *broom* package, which summarise model objects into a dataframe, which are easier to use for manipulating and plotting. This package has 3 main functions: `glance()`, which shows model statistics in one row, `tidy()`, which shows the model coefficients in a tidy summary, and `augment()`, which provides additional information on the fitted values and residuals. `glance()` and `tidy()` are also useful for creating model summary tables. Note that `glance` provides the F -test statistics, while `tidy` provides the t -test statistics. However, we can produce a tidy sum of squares table using `tidy()` with `anova()`.

```
library(broom)
library(knitr)
glance(mod.lm)
```

```
##   r.squared adj.r.squared   sigma statistic    p.value df    logLik
## 1 0.1414015   0.1326403 2.981904  16.13949 0.0001154607  2 -250.1399
##           AIC       BIC deviance df.residual
## 1 506.2799 514.0954 871.3919           98
```

```
kable(tidy(mod.lm), digits=5)
```

term	estimate	std.error	statistic	p.value
(Intercept)	1.39449	0.29867	4.66905	0.00001
variable	-1.29580	0.32255	-4.01740	0.00012

```
kable(tidy(anova(mod.lm)), digits=5)
```

term	df	sumsq	meansq	statistic	p.value
variable	1	143.5084	143.50839	16.13949	0.00012
Residuals	98	871.3919	8.89175	NA	NA

```
head(augment(mod.lm))
```

```
##   response variable    .fitted   .se.fit   .resid    .hat
## 1  5.2860866 -0.9965824  2.68586388 0.4509706  2.6002227 0.02287226
## 2 -3.5650202  0.7218241  0.45915178 0.3681751 -4.0241720 0.01524479
## 3  3.4738408 -0.6172088  2.19427143 0.3681732  1.2795693 0.01524463
## 4 -1.5246633  2.0293916 -1.23519483 0.7039776 -0.2894685 0.05573528
## 5 -1.4315520  1.0654161  0.01392519 0.4423806 -1.4454772 0.02200923
## 6  0.6915967  0.9872197  0.11525202 0.4240906  0.5763446 0.02022693
##   .sigma   .cooksd  .std.resid
## 1 2.985312 0.0091077290  0.88214740
## 2 2.968819 0.0143152999 -1.35993665
## 3 2.994375 0.0014473398  0.43242015
## 4 2.997083 0.0002945289 -0.09989891
## 5 2.993559 0.0027035867 -0.49017385
## 6 2.996653 0.0003935736  0.19526562
```