

An Introduction to Regression Analysis

By John Tyssedal

Sir Francis Galton (1822-1911), a well-known British anthropologist and meteorologist, was the first to introduce the word regression. It appears in a publication in Nature in 1885 where it is used to describe a biological phenomenon, namely that the heights of descendants of tall ancestors tend to regress down towards a normal average, a phenomenon also known as regression towards the mean. To Francis Galton regression had a pure biological meaning, but the term regression soon came to be applied to relationships in situations other than the one from which it originally arose. Today the goal of performing a regression analysis is to find a relationship between a response variable and one or more regression variables or to make models for prediction. Regression analysis is one of the most widely used statistical techniques. You may find that for practitioners of statistics, regression analysis is the method they are most familiar with. In the simplest case there is only one regression variable to explain the response variable. This is called simple linear regression.

Simple linear regression

The model for simple linear regression is $Y = \alpha + \beta x + \varepsilon$ where α and β are constants. The response variable Y (also called dependent variable) is a random variable while the regression variable x (also called independent or explanatory variable) can be deterministic or stochastic. ε is a noise variable assumed to have expected value 0 and constant variance, σ^2 . To fit a model, we need to have a certain number, say n , of observation pairs (x_i, y_i) of the regression and the response variable. (y_1, y_2, \dots, y_n) can then be considered to be realizations of (Y_1, Y_2, \dots, Y_n) , where

$$Y_i = \alpha + \beta x_i + \varepsilon_i, \quad i = 1, 2, \dots, n. \quad (1)$$

and $(\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n)$ are assumed to be uncorrelated. It is important to be aware of that by a linear model we mean linear in the coefficients such that x_i can be substituted by any transformation, for instance $\ln(x_i)$, x_i^2 or $\sin(x_i)$, as long as the transformation is allowed.

While the response variables, Y_1, Y_2, \dots, Y_n are random variables, the regression variables x_1, x_2, \dots, x_n can be either deterministic or stochastic. This has some implications in the interpretation of the model. If x_1, x_2, \dots, x_n are deterministic we have

$$E[Y_i] = \alpha + \beta x_i.$$

If x_1, x_2, \dots, x_n are stochastic variables we have

$$E[Y_i | x = x_i] = \alpha + \beta x_i.$$

Hence, the parameter α is the expected response value when $x=0$. Depending on whether x_1, x_2, \dots, x_n are deterministic or stochastic the parameter β is the expected change in the mean or the conditional mean respectively of the response when the regression variable changes with one unit.

Now let us assume we have n observation pairs: $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$. If these satisfy (1), we should have:

$$\begin{aligned} y_1 &= \alpha + \beta x_1 + \varepsilon_1^* \\ y_2 &= \alpha + \beta x_2 + \varepsilon_2^* \\ &\vdots \\ y_n &= \alpha + \beta x_n + \varepsilon_n^* \end{aligned}$$

or

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \end{bmatrix} + \begin{bmatrix} \varepsilon_1^* \\ \varepsilon_2^* \\ \vdots \\ \varepsilon_n^* \end{bmatrix},$$

with ε_i^* being realizations of ε_i , $i=1,2,\dots,n$.

$$\text{With } \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \mathbf{X} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}, \boldsymbol{\beta} = \begin{bmatrix} \alpha \\ \beta \end{bmatrix} \text{ and } \boldsymbol{\varepsilon}^* = \begin{bmatrix} \varepsilon_1^* \\ \varepsilon_2^* \\ \vdots \\ \varepsilon_n^* \end{bmatrix}, \text{ we get the shorter notation}$$

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}^*.$$

The least square estimates of (α, β) are the values (a, b) that minimizes

$$Q = \sum_{i=1}^n (y_i - a - bx_i)^2$$

with respect to a and b .

These are known to be the ones that satisfies

$$\frac{dQ}{da} = 0 \Leftrightarrow -2 \sum_{i=1}^n (y_i - a - bx_i) = 0 \quad (2)$$

$$\frac{dQ}{db} = 0 \Leftrightarrow -2 \sum_{i=1}^n (y_i - a - bx_i) x_i = 0. \quad (3)$$

From (2) and (3) we obtain the two normal equations

$$na + b \sum_{i=1}^n x_i = \sum_{i=1}^n y_i \quad (4)$$

$$a \sum_{i=1}^n x_i + b \sum_{i=1}^n x_i^2 = \sum_{i=1}^n x_i y_i \quad (5)$$

which gives

$$b = \frac{\sum_{i=1}^n (x_i - \bar{x}) y_i}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad \text{and} \quad a = \bar{y} - b\bar{x} \quad (6)$$

with corresponding least squares estimators

$$\hat{\beta} = \frac{\sum_{i=1}^n (x_i - \bar{x}) Y_i}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad \text{and} \quad \hat{\alpha} = \bar{Y} - \hat{\beta} \bar{x}. \quad (7)$$

Estimated model (estimated expected regression line) is then:

$$\hat{y} = a + bx \quad \text{or} \quad \hat{y}_i = a + bx_i.$$

The deviations $e_i = y_i - \hat{y}_i$, $i=1,2,\dots,n$ are called residuals. These are estimates of $(\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n)$ and should have approximately the same properties, i.e. uncorrelated with mean equal to zero and constant variance. Figure 1 shows some pattern that may be observed when plotting the residuals.

The plot of e_1 against \hat{y} ("yhat") shows an increase in the residual variance when \hat{y} is increasing. The plot of e_2 against \hat{y} is u-shaped indicating the lack of a quadratic term. This is supported by the plot of e_2 against x which has the same form. Finally, the plot of e against \hat{y} shows a residual plot, where no specific pattern is observed. Such a plot will support our assumption about uncorrelated residuals with a constant variance. In addition, a normal plot should be used to check if a normal distribution is an appropriate assumption. Most procedures for statistical inference in regression analysis rely on such an assumption.

Transformations for constant variance.

When a pattern as observed when e1 is plotted against \hat{y} occurs, the response should be transformed to obtain constant variance, otherwise our least squares estimators will not be optimal. The following calculations give a motivation for how the response should be transformed. Let $z = g(y)$ be a transformation of y . A first order approximation with a Taylor series gives us:

$$z \approx g(\mu) + g'(\mu)(y - \mu) \text{ where } \mu \text{ is some constant.}$$

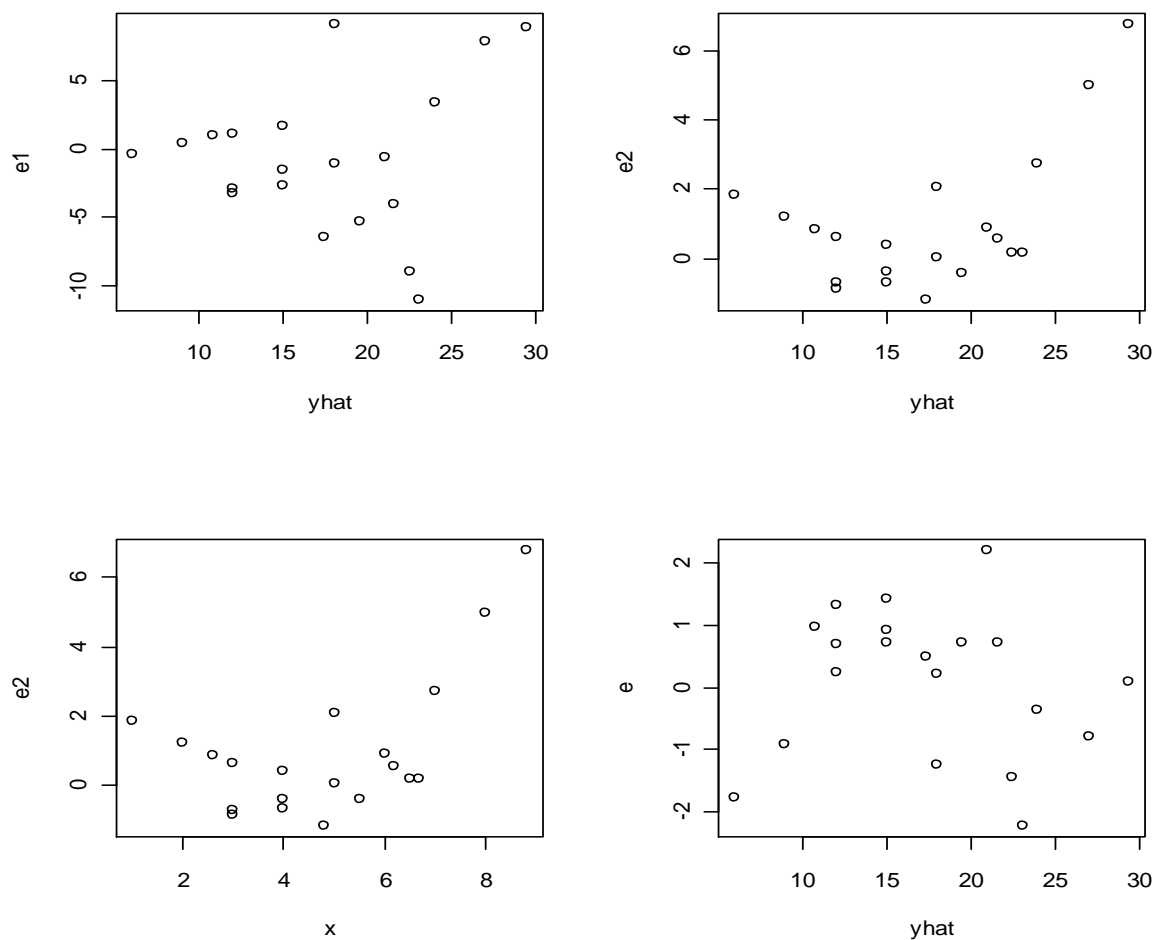


Figure 1. Four residuals plots illustrating: increase in variance with increasing \hat{y} , two u-shaped residual plots illustrating lack of a second order term in the model and finally a plot where the residuals have no particular pattern and thereby support our choice of model.

Hence as an approximation it should be possible to write a random variable $Z = g(Y)$ as:

$$Z \approx g(\mu) + g'(\mu)(Y - \mu) \text{ where } \mu = E(Y).$$

Thereby we get

$$\text{Var}(Z) \approx g'(\mu)^2 \text{Var}(Y).$$

For the variance of $Z = g(Y)$ to be (approximately) constant we must have:

$$g'(\mu)^2 \text{Var}(Y) = k \text{ or } g'(\mu) = \frac{\sqrt{k}}{SD(Y)} \text{ where } k \text{ is a constant.}$$

The following situations give raise to the most common transformations:

$$SD(Y) \propto \sqrt{\mu} \Rightarrow g'(\mu) = \frac{\sqrt{k}}{c\sqrt{\mu}} \text{ and an appropriate transformation will be } g(Y) = \sqrt{Y}$$

$$SD(Y) \propto \mu \Rightarrow g'(\mu) = \frac{\sqrt{k}}{c\mu} \text{ and an appropriate transformation will be } g(Y) = \ln(Y)$$

$$SD(Y) \propto \mu^2 \Rightarrow g'(\mu) = \frac{\sqrt{k}}{c\mu^2} \text{ and an appropriate transformation will be } g(Y) = \frac{1}{Y}$$

There is a constant missing in the transformations but that is irrelevant since the goal is to obtain constant variance. In general the transformation can be of the form Y^λ with λ arbitrary, where $\lambda = 0$ means the logarithm. If none of the transformations above gives a satisfactory results, one can try the Box-Cox transformation. The steps are as follows:

Suppose $Y_i > 0, \forall i$.

1. Consider a value of λ from a grid on a selected range for instance $[-2,2], [-1,1]$.
2. For each value of λ compute new responses:

$$V_i = \begin{cases} \frac{Y_i^\lambda - 1}{\lambda(Y^*)^{\lambda-1}}, & \lambda \neq 0 \\ Y^* \ln(Y_i), & \lambda = 0 \end{cases}$$

where $Y^* = \left(\prod_{i=1}^n Y_i \right)^{\frac{1}{n}}$ i.e. the geometric mean of Y_1, \dots, Y_n . Hence

$$\mathbf{Y} = (Y_1, \dots, Y_n)^T \rightarrow \mathbf{V} = (V_1, \dots, V_n)^T.$$

3. Regress \mathbf{V} on \mathbf{X} and find $SS_E(\lambda) = \sum_{i=1}^n (v_i - \hat{v}_i)^2$ for each λ .

4. Plot $SS_E(\lambda)$ versus λ . Draw a smooth curve and find $\lambda = \lambda^*$ that gives the smallest SS_E . If λ^* is close to values like $\dots, -2, -1.5, -1, -0.5, 0, 0.5, 1, 1.5, 2, \dots$, use these instead for easier interpretation.
5. Use the transformation Y^{λ^*} .

Transformations can also sometimes transfer a nonlinear model into a linear one as shown for the following three non-linear models below:

$$Y_i = \alpha e^{\beta x_i + \varepsilon_i} \Rightarrow \ln(Y_i) = \ln \alpha + \beta x_i + \varepsilon_i,$$

$$Y_i = \alpha x_i^\beta \varepsilon_i \Rightarrow \ln(Y_i) = \ln \alpha + \beta \ln(x_i) + \ln(\varepsilon_i),$$

$$Y_i = \frac{x_i}{\alpha + (\beta + \varepsilon_i)x_i} \Rightarrow \frac{1}{Y_i} = \frac{\alpha}{x_i} + \beta + \varepsilon_i.$$

After the transformation the model has become linear in the parameters $(\ln(\alpha), \beta)$, $(\ln(\alpha), \beta)$ and (α, β) respectively. These cases show that it may also be necessary to transform regression variables to obtain a linear model. Variance stabilized data will also normally have a distribution that is close to the normal distribution.

An example

The dataset below consists of 30 observation pairs of density and stiffness for a particular wood product. Since data for density are more easily obtainable than stiffness, it is of interest to find a relationship between stiffness and density such that given a density one is able to estimate or predict the stiffness. Therefore, stiffness will be taken to be the response variable and density to be the regression variable.

i	(x_i, y_i)	i	(x_i, y_i)	i	(x_i, y_i)
1	(9.5, 14184)	11	(17.4, 43243)	21	(25.6, 96305)
2	(8.4, 17502)	12	(15.0, 25319)	22	(23.4, 104170)
3	(9.8, 14007)	13	(15.2, 28028)	23	(24.4, 72594)
4	(11, 19443)	14	(16.4, 41792)	24	(23.3, 49512)
5	(8.3, 7573)	15	(16.7, 49499)	25	(19.5, 32207)
6	(9.9, 14194)	16	(15.4, 25312)	26	(21.2, 48218)
7	(8.6, 9714)	17	(15.0, 26222)	27	(22.8, 70453)
8	(6.4, 8076)	18	(14.5, 22148)	28	(21.7, 47661)
9	(7.0, 5304)	19	(14.8, 26751)	29	(19.8, 38138)
10	(8.2, 10728)	20	(13.6, 18036)	30	(21.3, 53045)

We always start by plotting the response variable against the regression variable. The plot is shown in Figure 2. It seems like stiffness increases with increasing density.

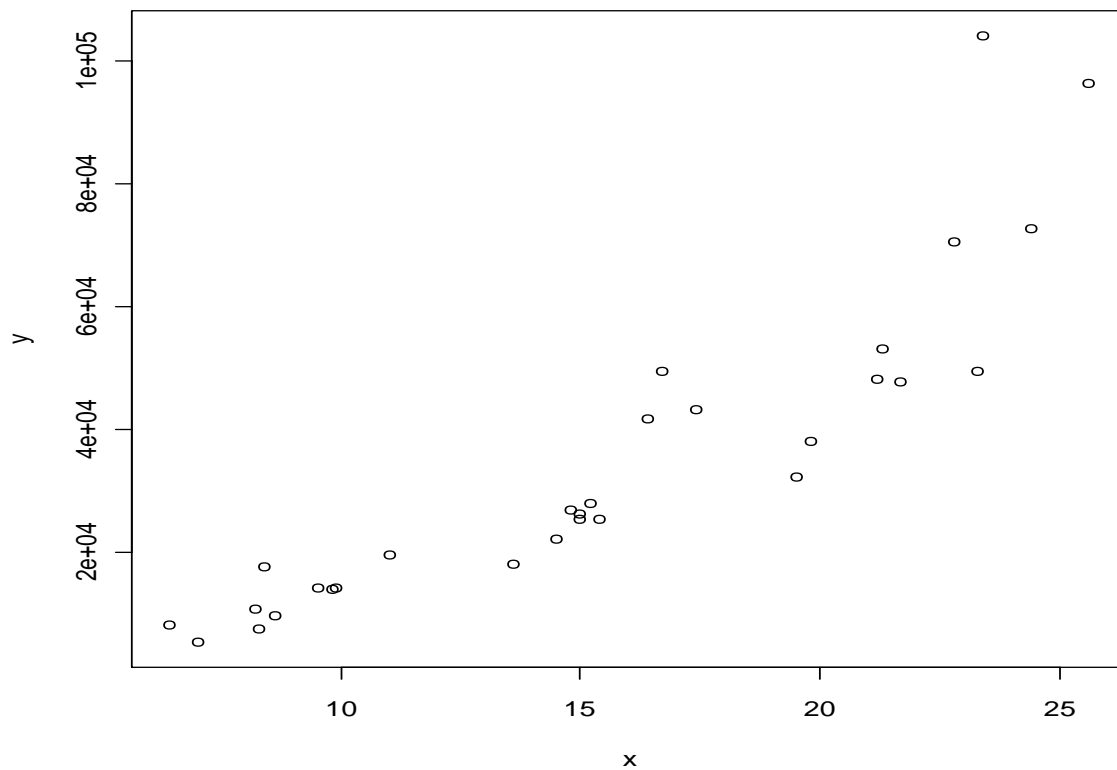


Figure 2. A plot of stiffness against density

Below is the result of a least square analysis performed with R. We shall comment on this analysis in greater detail later. Now we notice that α is estimated to -25513.3 and β to 3888.8 and that both are highly significant. The estimated value for σ is 11620.

```

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept) -25513.3    6101.6   -4.181  0.000258 ***
x              3888.8     369.8   10.515  3.14e-11 ***
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1 ---
Residual standard error: 11620 on 28 degrees of freedom

```

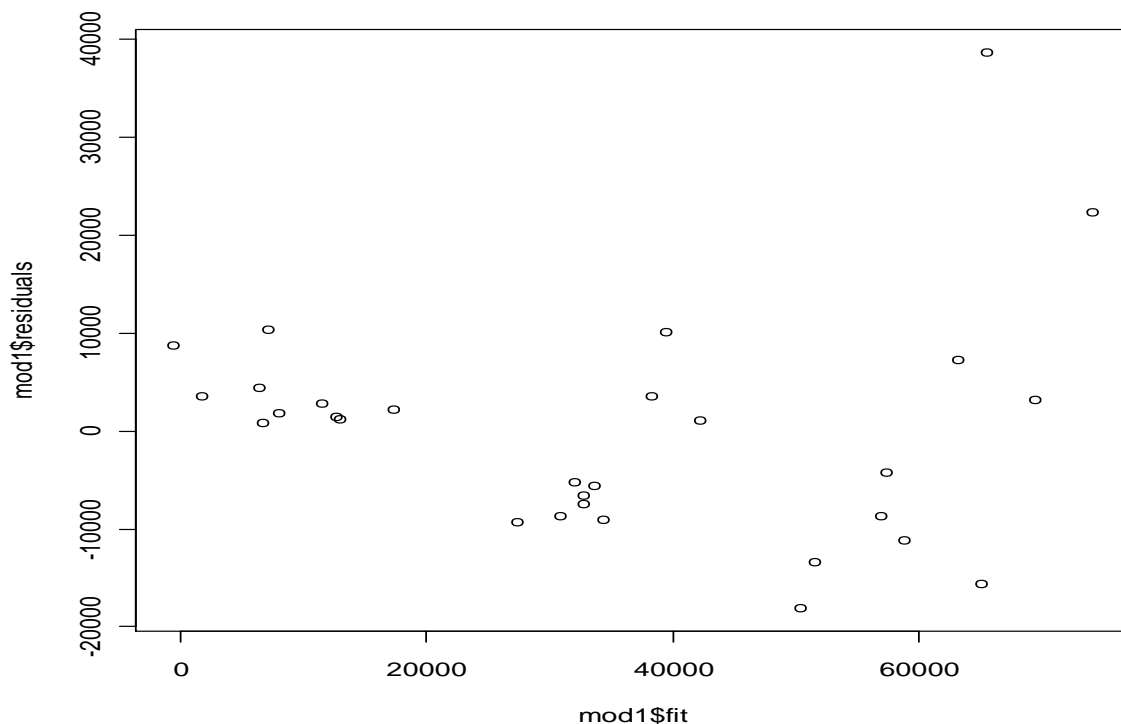


Figure 3. Plot of residuals against \hat{y} .

The plot of the residuals shows that the variance is increasing when \hat{y} or the estimated mean is increasing. This calls upon a transformation. A plot of the log transformed stiffness against density is given in Figure 4.

With $\log(y)$ as the response we get the following least squares estimates

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	8.25193	0.12819	64.37	< 2e-16 ***
x	0.12518	0.00777	16.11	1.08e-15 ***

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

Residual standard error: 0.2441 on 28 degrees of freedom

The plot of the residuals against \hat{y} with the log transformed response indicates a far better fit as can be seen in Figure 5, and the normal plot indicates that their distribution is not too far away from a normal distribution.

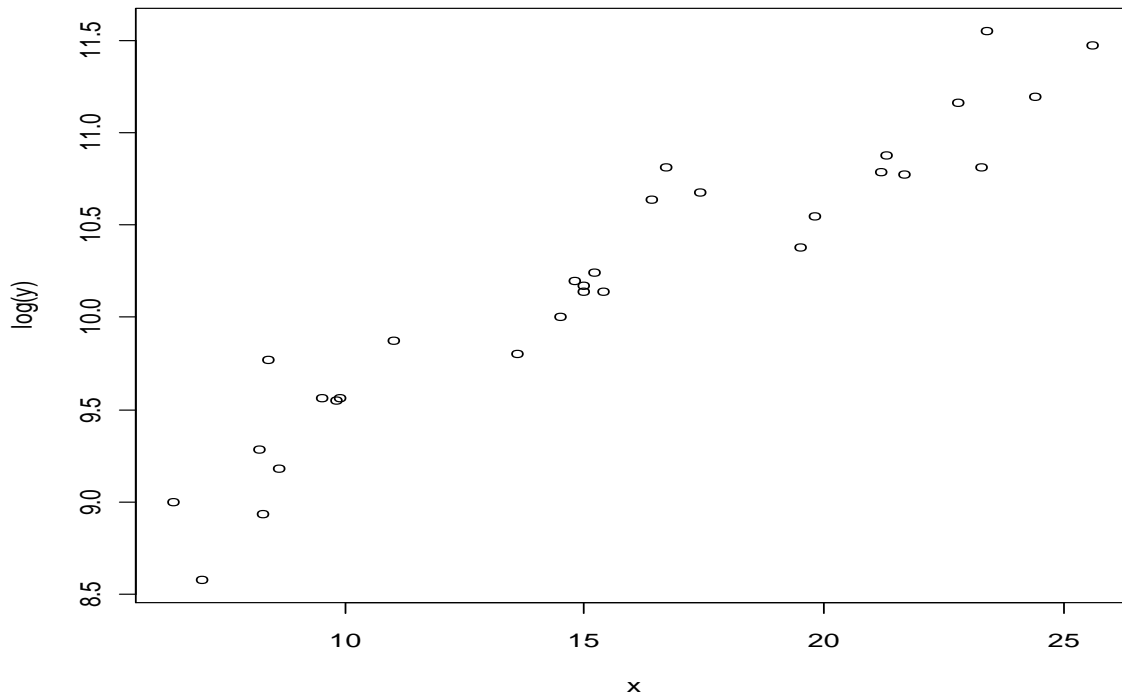


Figure 4. Plot of the logarithm of stiffness against density.

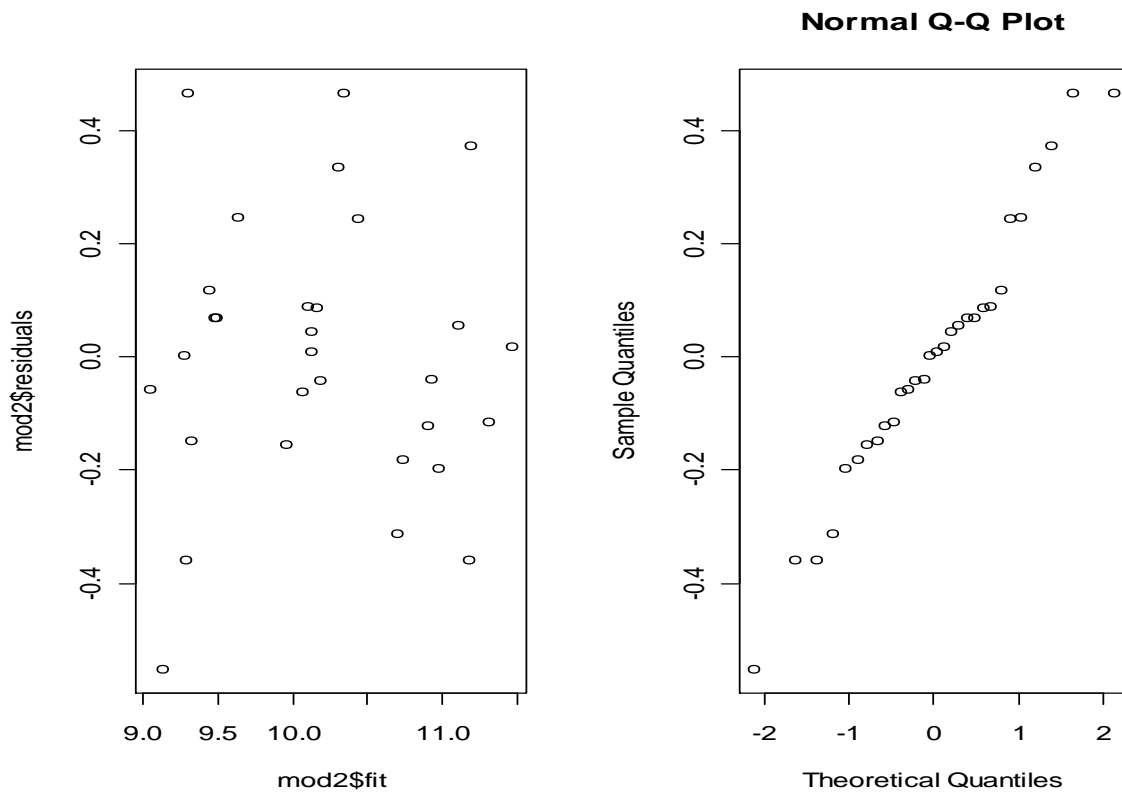


Figure 5. Plot of residuals against \hat{y} for the model with the log transformed response together with a Normal Q-Q plot of the residuals.

Approximation of the expectation and variance for functions of random variables.

The technique used to find appropriate transformations can be used to find an approximation of expectation and variance for non-linear functions of random variable.

An example.

The estimated simple regression model can be written as $\hat{y} = \bar{y} + b(x - \bar{x})$. Suppose we are interested in the x -value for which the regression line crosses the x -axis. Simple calculations will suggest this value to be estimated with $x = \bar{x} - \frac{\bar{y}}{b}$. Now suppose we are interested in the uncertainty in this estimate. An estimator for the crossing value is $\bar{x} - \frac{\bar{Y}}{\hat{\beta}}$. The random variables here are \bar{Y} and $\hat{\beta}$.

Both these have known variances given as $\frac{\sigma^2}{n}$ and $\frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$ respectively. They can also

be shown to be uncorrelated random variables. However, we have no formal way to find the exact value of $Var\left(\frac{\bar{Y}}{\hat{\beta}}\right)$ by calculations.

For a function $g(x, y)$ a first order Taylor expansion gives:

$$g(x, y) \approx g(a, b) + \frac{dg}{dx}(a, b)(x - a) + \frac{dg}{dy}(a, b)(y - b).$$

Then expanding a random variable $Z = g(X, Y)$ around $\mu_1 = E(X)$ and $\mu_2 = E(Y)$ gives

$$g(X, Y) \approx g(\mu_1, \mu_2) + \frac{dg}{dx}(\mu_1, \mu_2)(X - \mu_1) + \frac{dg}{dy}(\mu_1, \mu_2)(Y - \mu_2) \quad (8)$$

From this expression we have:

$$Var(g(X, Y)) \approx \frac{dg}{dx}(\mu_1, \mu_2)^2 Var(X) + \frac{dg}{dy}(\mu_1, \mu_2)^2 Var(Y) + 2\left(\frac{dg}{dx}(\mu_1, \mu_2) \frac{dg}{dy}(\mu_1, \mu_2) Cov(X, Y)\right)$$

If the two random variable X and Y are uncorrelated the last term vanishes.

Applied to our crossing problem we get:

$$Var\left(\frac{\bar{Y}}{\hat{\beta}}\right) \approx \frac{1}{\beta^2} \cdot \frac{\sigma^2}{n} + \frac{(\alpha + \beta \bar{x})^2}{\beta^4} \cdot \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sigma^2}{\beta^2} \left(\frac{1}{n} + \frac{(\alpha + \beta \bar{x})^2}{\beta^2 \sum_{i=1}^n (x_i - \bar{x})^2} \right).$$

In practice estimates of α and β will be used to get an approximate value for the variance.

The general approximation formula for the variance of a random variable $g(X_1, X_2, \dots, X_n)$ is:

$$Var(g(X_1, X_2, \dots, X_n)) \approx \sum_{i=1}^n \frac{dg}{dx_i}(\mu_1, \dots, \mu_n)^2 Var(X_i) + 2 \sum_{i < j} \frac{dg}{dx_i}(\mu_1, \dots, \mu_n) \cdot \frac{dg}{dx_j}(\mu_1, \dots, \mu_n) Cov(X_i, X_j)$$

Using the Taylor expansion, it is also possible to find an approximation to the expected value of a function of random variables. By taking the expected value of the linearized expression in (8) we get:

$$E[g(X, Y)] \approx g(\mu_1, \mu_2)$$

and in general

$$E[g(\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n)] \approx g(\mu_1, \mu_2, \dots, \mu_n).$$

For the approximation formulas to work well the variances of the respective variables should be small. However, simulation shows that the approximation formula for the variance works well in most situations.

The reader should be aware of how many times these approximation formulas are needed when working with measured data. For instance, we measure length, width, height, radius, current, resistance etc. and may want to calculate, area, volume or voltage. In each case we have expressions that are not linear in the variables we measure. It is also often the case that our measurements of the different variables may be assumed to be uncorrelated and hence the covariance terms in the expression for the variance vanish.

One more example

Sometimes measurements are given with a relative uncertainty. Let us assumed that the height of a cylinder is measured with a relative uncertainty of 2% and that the radius is measured with a relative uncertainty of 1%.

As an approximation:

$$\ln(X) \approx \ln(\mu_x) + \frac{1}{\mu_x}(X - \mu_x).$$

And hence

$$\sigma_{\ln(X)} \approx \frac{\sigma_X}{\mu_X}, \text{ the coefficient of variation or relative uncertainty.}$$

The volume of a cylinder is given by $V = \pi r^2 h$. Taking the logarithm, we get:

$$\ln(V) = \ln(\pi) + 2\ln(r) + \ln(h),$$

and a direct calculation gives $\sigma_{\ln(V)} = \sqrt{4\sigma_{\ln(r)}^2 + \sigma_{\ln(h)}^2}$.

The relative uncertainty in V is $\frac{\sigma_V}{\mu_V} \approx \sigma_{\ln(V)} = \sqrt{4(0.01)^2 + (0.02)^2} = 0.028$ or 2.8%.

Multiple linear regression. Least squares estimators and their covariance matrix.

An example. Sour precipitation.

To study the influence of sour precipitation in Norwegian lakes, SFT in 1986 carried through an investigation where data from 1005 lakes were collected. In this example 26 random lakes, 16 from Telemark and 10 from Trøndelag are chosen out. The variables and the data are given below.

y = Measured pH-value
x1 = Content of SO_4
x2 = Content of NO_3
x3 = Content of Ca
x4 = Content of latent aluminum
x5 = Content of organic material
x6 = Area of water
x7 = Location (0 – Telemark, 1-Trøndelag)

	y	x1	x2	x3	x4	x5	x6	x7
1	5.38	4.9	39	1.54	78	2.02	0.30	0
2	5.68	4.1	75	1.55	17	2.98	1.85	0
3	5.04	3.5	80	0.83	157	3.40	0.25	0
4	4.81	3.8	75	0.53	163	3.42	0.30	0
5	4.92	3.8	90	0.82	105	3.41	0.25	0
6	5.34	2.6	49	0.62	114	1.90	0.65	0
7	5.74	2.7	79	1.08	15	2.53	1.15	0
8	5.17	2.6	90	0.67	90	1.89	0.91	0
9	5.02	2.4	64	0.41	107	0.97	0.60	0
10	5.88	2.8	27	1.15	12	3.04	0.37	0
11	5.36	3.4	13	0.89	93	2.95	0.58	0
12	5.26	2.7	14	0.74	72	2.75	0.15	0
13	5.69	3.2	13	1.03	99	3.34	0.53	0
14	5.51	2.5	79	0.67	80	0.78	0.56	0
15	5.25	1.5	77	0.33	72	0.11	1.08	0
16	6.06	3.7	15	1.94	4	5.53	0.72	0
17	6.08	1.9	16	1.05	0	5.63	0.48	1
18	6.08	1.3	2	0.81	0	3.84	0.25	1
19	6.20	2.2	32	1.40	7	3.77	1.87	1
20	5.64	2.2	21	0.75	13	5.62	0.53	1

21	5.75	1.6	7	0.79	2	3.95	0.35	1
22	5.43	2.0	27	0.47	16	2.39	1.03	1
23	5.82	1.8	17	0.74	3	3.08	0.32	1
24	5.50	2.0	6	0.49	10	2.82	0.28	1
25	5.62	1.5	3	0.36	7	1.71	0.10	1
26	5.41	1.7	7	0.54	9	6.48	0.21	1

We notice that there are seven potential regression variables that may affect the response. Our modeling goal is to find a functional relationship between the response y and those regression variables that really can explain the variation in y . Let us first find out how we can obtain the least squares estimators in a multiple regression model.

The model for multiple linear regression is:

$$Y = \beta_0 + \beta_1 x_1 + \cdots + \beta_k x_k + \varepsilon,$$

where again Y is the response, x_1, \dots, x_k the regression variables and ε the error term.

Alternatively, we have

$$Y_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_k x_{ik} + \varepsilon_i.$$

Here $E(\varepsilon_i) = 0$, $\text{Var}(\varepsilon_i) = \sigma^2$, $i=1,2,\dots,n$ and $E(\varepsilon_i \varepsilon_j) = 0$, $i \neq j$.

We observe

$$(y_i, x_{i1}, x_{i2}, \dots, x_{ik}), \quad i=1,2,\dots,n$$

Which according to the model must satisfy:

$$\begin{aligned} y_1 &= \beta_0 + \beta_1 x_{11} + \beta_2 x_{12} + \cdots + \beta_k x_{1k} + \varepsilon_1^* \\ y_2 &= \beta_0 + \beta_1 x_{21} + \beta_2 x_{22} + \cdots + \beta_k x_{2k} + \varepsilon_2^* \\ &\vdots \\ y_n &= \beta_0 + \beta_1 x_{n1} + \beta_2 x_{n2} + \cdots + \beta_k x_{nk} + \varepsilon_n^* \end{aligned}$$

where ε_i^* is a realization of ε_i . Written in matrix form this becomes:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}^*$$

Estimates for $\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_k)^t$, $(b_0, b_1, \dots, b_k)^t$, are obtained by minimizing

$$Q = \sum_{i=1}^n (y_i - b_0 - b_1 x_{i1} - b_2 x_{i2} - \cdots - b_k x_{ik})^2$$

with respect to b_0, b_1, \dots, b_k . Note that Q can be shown to be a convex function in b_0, b_1, \dots, b_k .

The estimated model becomes:

$$\hat{y}_i = b_0 + b_1 x_{i1} + \dots + b_k x_{ik}.$$

Setting the partial derivatives equal to zero gives the normal equations:

$$\begin{aligned} -2 \sum_{i=1}^n (y_i - b_0 - b_1 x_{i1} - b_2 x_{i2} - \dots - b_k x_{ik}) &= -2 \sum_{i=1}^n (y_i - \hat{y}_i) = 0 \\ -2 \sum_{i=1}^n x_{i1} (y_i - b_0 - b_1 x_{i1} - b_2 x_{i2} - \dots - b_k x_{ik}) &= -2 \sum_{i=1}^n x_{i1} (y_i - \hat{y}_i) = 0 \\ &\vdots \\ -2 \sum_{i=1}^n x_{ik} (y_i - b_0 - b_1 x_{i1} - b_2 x_{i2} - \dots - b_k x_{ik}) &= -2 \sum_{i=1}^n x_{ik} (y_i - \hat{y}_i) = 0 \end{aligned}$$

which can be written as:

$$\begin{bmatrix} 1 & 1 & \dots & 1 \\ x_{11} & x_{21} & \dots & x_{n1} \\ \vdots & \vdots & \ddots & \vdots \\ x_{1k} & x_{2k} & \dots & x_{nk} \end{bmatrix} \begin{bmatrix} y_1 - \hat{y}_1 \\ y_2 - \hat{y}_2 \\ \vdots \\ y_n - \hat{y}_n \end{bmatrix} = \mathbf{0} \text{ or } \mathbf{X}'(\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{X}'(\mathbf{y} - \mathbf{X}\mathbf{b}) = \mathbf{0}$$

or

$$\mathbf{X}'\mathbf{X}\mathbf{b} = \mathbf{X}'\mathbf{y} \text{ which implies } \mathbf{b} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y} \text{ if } (\mathbf{X}'\mathbf{X})^{-1} \text{ exists.}$$

The least squares estimator is therefore: $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y}$. Notice that the first normal equation guarantees that the sum of the residuals is zero. If there is no constant term in the model, we have no such guarantee.

Expectation and Covariance matrix of the estimators.

The least square estimator is given by:

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon})$$

which implies $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}) = \boldsymbol{\beta} + (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\boldsymbol{\varepsilon}$.

Since $E(\boldsymbol{\varepsilon}) = \mathbf{0}$ we have that $E(\hat{\boldsymbol{\beta}}) = \boldsymbol{\beta}$ i.e. the least square estimator is unbiased.

The covariance matrix of $\hat{\boldsymbol{\beta}}$ is:

$$\begin{aligned}
Cov(\hat{\boldsymbol{\beta}}) &= E \left\{ \begin{bmatrix} \hat{\beta}_0 - \beta_0 \\ \hat{\beta}_1 - \beta_1 \\ \vdots \\ \hat{\beta}_k - \beta_k \end{bmatrix} \begin{bmatrix} \hat{\beta}_0 - \beta_0, & \hat{\beta}_1 - \beta_1, & \dots, & \hat{\beta}_k - \beta_k \end{bmatrix} \right\} \\
&= E \begin{bmatrix} (\hat{\beta}_0 - \beta_0)^2 & (\hat{\beta}_0 - \beta_0)(\hat{\beta}_1 - \beta_1) & \dots & (\hat{\beta}_0 - \beta_0)(\hat{\beta}_k - \beta_k) \\ (\hat{\beta}_1 - \beta_1)(\hat{\beta}_0 - \beta_0) & (\hat{\beta}_1 - \beta_1)^2 & \dots & (\hat{\beta}_1 - \beta_1)(\hat{\beta}_k - \beta_k) \\ \vdots & \vdots & \dots & \vdots \\ (\hat{\beta}_k - \beta_k)(\hat{\beta}_0 - \beta_0) & (\hat{\beta}_k - \beta_k)(\hat{\beta}_1 - \beta_1) & \dots & (\hat{\beta}_k - \beta_k)^2 \end{bmatrix} \\
&= \begin{bmatrix} Var(\hat{\beta}_0) & Cov(\hat{\beta}_0, \hat{\beta}_1) & \dots & Cov(\hat{\beta}_0, \hat{\beta}_k) \\ Cov(\hat{\beta}_1, \hat{\beta}_0) & Var(\hat{\beta}_1) & \dots & Cov(\hat{\beta}_1, \hat{\beta}_k) \\ \vdots & \vdots & \dots & \vdots \\ Cov(\hat{\beta}_k, \hat{\beta}_0) & Cov(\hat{\beta}_k, \hat{\beta}_1) & \dots & Var(\hat{\beta}_k) \end{bmatrix}.
\end{aligned}$$

We get :

$$\begin{aligned}
E \left[(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^t \right] &= E \left[(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \boldsymbol{\varepsilon} \boldsymbol{\varepsilon}' \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \right] \\
&= (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' E(\boldsymbol{\varepsilon} \boldsymbol{\varepsilon}') \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \sigma^2 \mathbf{I} \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} = \sigma^2 (\mathbf{X}'\mathbf{X})^{-1}.
\end{aligned}$$

An example. Simple linear regression.

The model can be written as:

$$\begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} 1 & x_{11} \\ 1 & x_{21} \\ \vdots & \vdots \\ 1 & x_{n1} \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

We get:

$$(\mathbf{X}'\mathbf{X}) = \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_{11} & x_{21} & \dots & x_{n1} \end{bmatrix} \begin{bmatrix} 1 & x_{11} \\ 1 & x_{21} \\ \vdots & \vdots \\ 1 & x_{n1} \end{bmatrix} = \begin{bmatrix} n & \sum_{i=1}^n x_{i1} \\ \sum_{i=1}^n x_{i1} & \sum_{i=1}^n x_{i1}^2 \end{bmatrix}$$

$$\mathbf{X}'\mathbf{Y} = \begin{bmatrix} \sum_{i=1}^n Y_i \\ \sum_{i=1}^n x_{i1} Y_i \end{bmatrix}$$

$$\text{and } (\mathbf{X}'\mathbf{X})^{-1} = \frac{1}{n \sum_{i=1}^n x_{i1}^2 - \left(\sum_{i=1}^n x_{i1} \right)^2} \begin{bmatrix} \sum_{i=1}^n x_{i1}^2 & -\sum_{i=1}^n x_{i1} \\ -\sum_{i=1}^n x_{i1} & n \end{bmatrix}$$

such that

$$\hat{\beta} = \frac{1}{\sum_{i=1}^n x_{i1}^2 - n \bar{x}_1^2} \left(-\bar{x}_1 \sum_{i=1}^n Y_i + \sum_{i=1}^n x_{i1} Y_i \right) = \frac{\sum_{i=1}^n Y_i (x_{i1} - \bar{x}_1)}{\sum_{i=1}^n x_{i1}^2 - n \bar{x}_1^2}$$

and

$$\hat{\alpha} = \frac{1}{\sum_{i=1}^n x_{i1}^2 - n \bar{x}_1^2} \left(\bar{Y} \left(\sum_{i=1}^n x_{i1}^2 - n \bar{x}_1^2 \right) + \bar{x}_1 \left(n \bar{Y} \bar{x}_1 - \sum_{i=1}^n x_{i1} Y_i \right) \right) = \bar{Y} - \hat{\beta} \bar{x}_1.$$

$$\text{Cov} \begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix} = \sigma^2 \begin{bmatrix} \frac{\sum_{i=1}^n x_{i1}^2}{n \left(\sum_{i=1}^n x_{i1}^2 - n \bar{x}_1^2 \right)} & -\frac{\bar{x}_1}{\sum_{i=1}^n x_{i1}^2 - n \bar{x}_1^2} \\ -\frac{\bar{x}_1}{\sum_{i=1}^n x_{i1}^2 - n \bar{x}_1^2} & \frac{1}{\sum_{i=1}^n x_{i1}^2 - n \bar{x}_1^2} \end{bmatrix}.$$

$\sum_{i=1}^n x_{i1}^2 - n \bar{x}_1^2$ can be written as $\sum_{i=1}^n (x_{i1} - \bar{x}_1)^2$. This follows because

$$\sum_{i=1}^n (x_{i1} - \bar{x}_1)^2 = \sum_{i=1}^n x_{i1}^2 - 2 \bar{x}_1 \sum_{i=1}^n x_{i1} + n \bar{x}_1^2 = \sum_{i=1}^n x_{i1}^2 - n \bar{x}_1^2.$$

Partitioning of variation

A measure of the total variance in the response is given by $\sum_{i=1}^n (y_i - \bar{y})^2$. We can split up this as follows:

First write

$$y_i - \bar{y} = y_i - \hat{y}_i + \hat{y}_i - \bar{y}$$

such that

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + 2 \sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2.$$

Now

$$\sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) = \sum_{i=1}^n (y_i - \hat{y}_i) \hat{y}_i - \sum_{i=1}^n (y_i - \hat{y}_i) \bar{y} = 0.$$

This follows because the two last expressions are zero, since the normal equations have to be fulfilled.

Thereby we have

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2.$$

This can also be written as:

$$SS_T = SS_E + SS_R.$$

Or the total sum of squares equals the error sum of squares + the regression sum of squares.

This partitioning of the total variation in the response is only true if the normal equations are fulfilled. In models with no constant term this partitioning is normally not valid.

Some results about idempotent matrices and distribution of quadratic forms and their application to multiple linear regression.

Most of the necessary theory behind procedures in regression analysis and in linear models in general relies upon properties of idempotent matrices and the spectral decomposition of symmetric matrices. In fact, in this course you will find these two mathematical topics involved in most of the theory we are building up.

Idempotent matrices

First define $\mathbf{J}_{n \times n}$ to be a $n \times n$ matrix with all entries equal to 1. Let $\mathbf{Y} = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix}$ be a vector of

random variables. We have:

$$\bar{\mathbf{Y}} = \begin{bmatrix} \bar{Y} \\ \bar{Y} \\ \vdots \\ \bar{Y} \end{bmatrix} = \frac{1}{n} \begin{bmatrix} 1 & 1 & \cdots & 1 \\ 1 & 1 & \cdots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \cdots & 1 \end{bmatrix} \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} = \frac{1}{n} \mathbf{J} \mathbf{Y}.$$

Now $\left(\mathbf{I}_{n \times n} - \frac{1}{n} \mathbf{J} \right)$ is obviously symmetric. Therefore $\mathbf{Y} - \bar{\mathbf{Y}} = \left(\mathbf{I}_{n \times n} - \frac{1}{n} \mathbf{J} \right) \mathbf{Y}$ and

$$(\mathbf{y} - \bar{\mathbf{y}})' (\mathbf{y} - \bar{\mathbf{y}}) = \mathbf{y}' \left(\mathbf{I} - \frac{1}{n} \mathbf{J} \right) \left(\mathbf{I} - \frac{1}{n} \mathbf{J} \right) \mathbf{y} = SS_T.$$

Further $\hat{\mathbf{Y}} = \mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{X} \left(\mathbf{X}' \mathbf{X} \right)^{-1} \mathbf{X}' \mathbf{Y} = \mathbf{H} \mathbf{Y}$ where $\mathbf{H} = \mathbf{X} \left(\mathbf{X}' \mathbf{X} \right)^{-1} \mathbf{X}'$ is the often referred to hat matrix in linear regression theory. $(\mathbf{I}_{n \times n} - \mathbf{H}_{n \times n})$ is symmetric and we obtain

$$\mathbf{Y} - \hat{\mathbf{Y}} = (\mathbf{I}_{n \times n} - \mathbf{H}_{n \times n}) \mathbf{Y} \text{ and } (\mathbf{y} - \hat{\mathbf{y}})' (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{y}' (\mathbf{I} - \mathbf{H}) (\mathbf{I} - \mathbf{H}) \mathbf{y} = SS_E$$

Also $\left(\mathbf{H} - \frac{1}{n} \mathbf{J} \right)$ is symmetric and $\hat{\mathbf{Y}} - \bar{\mathbf{Y}} = \left(\mathbf{H} - \frac{1}{n} \mathbf{J} \right) \mathbf{Y}$. Hence

$$(\hat{\mathbf{y}} - \bar{\mathbf{y}})' (\hat{\mathbf{y}} - \bar{\mathbf{y}}) = \mathbf{y}' \left(\mathbf{H} - \frac{1}{n} \mathbf{J} \right) \left(\mathbf{H} - \frac{1}{n} \mathbf{J} \right) \mathbf{y} = SS_R.$$

For the matrices $\mathbf{J}_{n \times n}$ and $\mathbf{H}_{n \times n}$ we have:

$$\frac{1}{n} \mathbf{J} \frac{1}{n} \mathbf{J} = \frac{1}{n} \mathbf{J} \text{ and } \mathbf{H} \mathbf{H} = \mathbf{X} \left(\mathbf{X}' \mathbf{X} \right)^{-1} \mathbf{X}' \mathbf{X} \left(\mathbf{X}' \mathbf{X} \right)^{-1} \mathbf{X}' = \mathbf{X} \left(\mathbf{X}' \mathbf{X} \right)^{-1} \mathbf{X}' = \mathbf{H}$$

Matrices that remain unchanged regardless of how many times we multiply them together are called idempotent.

We further get

$$\left(\mathbf{I} - \frac{1}{n} \mathbf{J} \right) \left(\mathbf{I} - \frac{1}{n} \mathbf{J} \right) = \mathbf{I} - \frac{1}{n} \mathbf{J} - \frac{1}{n} \mathbf{J} + \frac{1}{n} \mathbf{J} \frac{1}{n} \mathbf{J} = \mathbf{I} - \frac{1}{n} \mathbf{J} - \frac{1}{n} \mathbf{J} + \frac{1}{n} \mathbf{J} = \mathbf{I} - \frac{1}{n} \mathbf{J} \text{ and}$$

$$(\mathbf{I} - \mathbf{H}) (\mathbf{I} - \mathbf{H}) = \mathbf{I} - \mathbf{H} - \mathbf{H} + \mathbf{H}^2 = \mathbf{I} - \mathbf{H} - \mathbf{H} + \mathbf{H} = \mathbf{I} - \mathbf{H}.$$

The following property of the hat matrix is useful:

R1. $\mathbf{H} \mathbf{X} = \mathbf{X} \left(\mathbf{X}' \mathbf{X} \right)^{-1} \mathbf{X}' \mathbf{X} = \mathbf{X}.$

This means that

$$\mathbf{H} \begin{bmatrix} 1 & \mathbf{x}_1 & \cdots & \mathbf{x}_k \\ 1 \\ \vdots \\ 1 \end{bmatrix} = \begin{bmatrix} 1 & \mathbf{x}_1 & \cdots & \mathbf{x}_k \\ 1 \\ \vdots \\ 1 \end{bmatrix} \text{ where } \mathbf{x}_i, i=1,2,\dots,k, \text{ is the column of observed}$$

values for the i 'th regression variable. In particular we get $\mathbf{H}\mathbf{1}=\mathbf{1}$ or $\sum_{j=1}^n h_{ij} = \sum_{i=1}^n h_{ij} = 1$ since the hat matrix is symmetric. In other words, the elements in any row and columns in the hat matrix sum to 1.

Then we obtain

$$\left(\mathbf{H} - \frac{1}{n}\mathbf{J}\right)\left(\mathbf{H} - \frac{1}{n}\mathbf{J}\right) = \mathbf{H}^2 - \frac{1}{n}\mathbf{J}\mathbf{H} - \frac{1}{n}\mathbf{H}\mathbf{J} + \frac{1}{n}\mathbf{J}\frac{1}{n}\mathbf{J} = \mathbf{H} - \frac{1}{n}\mathbf{J} - \frac{1}{n}\mathbf{J} + \frac{1}{n}\mathbf{J} = \mathbf{H} - \frac{1}{n}\mathbf{J} \text{ since } \mathbf{JH} = \mathbf{HJ} = \mathbf{J}.$$

A matrix \mathbf{A} that satisfies $\mathbf{A}\mathbf{A} = \mathbf{A}$ is called idempotent. If \mathbf{A} also is symmetric it is also called a projection matrix.

R2. \mathbf{H} , $\frac{1}{n}\mathbf{J}$, $\mathbf{I} - \mathbf{H}$, $\mathbf{I} - \frac{1}{n}\mathbf{J}$ and $\mathbf{H} - \frac{1}{n}\mathbf{J}$ are all idempotent matrices and also projection matrices.

Let \mathbf{A} and \mathbf{B} be two matrices such that their product \mathbf{AB} is defined. Then

$$\text{rank}(\mathbf{AB}) \leq \min(\text{rank}(\mathbf{A}), \text{rank}(\mathbf{B})).$$

This follows since the columns of \mathbf{AB} are linear combinations of the columns of \mathbf{A} and the rows in \mathbf{AB} are linear combinations of the rows of \mathbf{B} .

For a symmetric matrix $\mathbf{A}_{n \times n}$ we have $\mathbf{A} = \mathbf{P}\mathbf{\Lambda}\mathbf{P}^t$ where $\mathbf{P} = [\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n]$, the matrix of

$$\text{eigenvectors and } \mathbf{\Lambda} = \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_n \end{bmatrix} \text{ is a diagonal matrix of eigenvalues. This is known}$$

as the spectral decomposition of the matrix \mathbf{A} .

R3. $\mathbf{A}_{n \times n}$ symmetric implies $\text{rank}(\mathbf{A}) = \text{the number of nonzero eigenvalues}$

Proof.

We have

$$\text{rank}(A) = \text{rank}(P \Lambda P^t) \leq \min(\text{rank}(P^t), \text{rank}(\Lambda P^t)) \leq \text{rank}(\Lambda P^t)$$

$$\text{and } \text{rank}(\Lambda P^t) = \text{rank}(P^t P \Lambda P^t) \leq \min(\text{rank}(P^t), \text{rank}(P \Lambda P^t)) \leq \text{rank}(P \Lambda P^t) = \text{rank}(A). \text{ Also } \text{rank}(\Lambda P^t) \leq \min(\text{rank}(\Lambda), \text{rank}(P^t)) \leq \text{rank}(\Lambda) = \text{rank}(\Lambda P^t P) \leq \text{rank}(\Lambda P^t).$$

$$\text{Hence } \text{rank}(A) = \text{rank}(\Lambda P^t) = \text{rank}(\Lambda).$$

R4. $A_{n \times n}$ symmetric and idempotent with rank r implies r eigenvalues are 1 and $n-r$ are zero.

Proof.

$Ax = \lambda x$ and $\lambda x^t x = x^t Ax = x^t A Ax = \lambda^2 x^t x$. Hence $x^t x \lambda (\lambda - 1) = 0$ which implies $\lambda = 1$ or $\lambda = 0$. Since the number of nonzero eigenvalues are r , r eigenvalues are 1 and $n-r$ are zero.

R5. $A_{n \times n}$ symmetric and idempotent implies $\text{tr}(A) = \text{rank}(A)$.

Proof .

$$\text{tr}(A) = \text{tr}(P \Lambda P^t) = \text{tr}(P^t P \Lambda) = \text{tr}(\Lambda) = \text{rank}(A).$$

Distribution of Quadratic forms

R6. Let $A_{n \times n}$ be a symmetric and idempotent matrix of rank r and let Y_1, Y_2, \dots, Y_n be independent normally distributed with expectation $\mu_1, \mu_2, \dots, \mu_n$ respectively and equal

variance σ^2 . Define $Y = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix}$ and $\underline{\mu} = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_n \end{bmatrix}$. Then $\frac{(Y - \underline{\mu})^t (A) (Y - \underline{\mu})}{\sigma^2}$ is $\chi^2(r)$.

Proof . $A = P \Lambda P^t$. Let $U = \begin{bmatrix} U_1 \\ \vdots \\ U_n \end{bmatrix}$ where $U_i = \frac{Y_i - \mu_i}{\sigma}$, $i=1, 2, \dots, n$. U_1, U_2, \dots, U_n are

independent random variables. Also $E(U_i) = 0$ and $SD(U_i) = 1$, $i=1, 2, \dots, n$. We get:

$$\frac{(Y - \underline{\mu})^t A (Y - \underline{\mu})}{\sigma^2} = U^t A U = U^t P \Lambda P^t U = Z^t \Lambda Z = \sum_{i=1}^r Z_i^2 \text{ where } Z = P^t U.$$

Further $\text{Cov}(\mathbf{Z}) = \text{Cov}(\mathbf{P}'\mathbf{U}) = E(\mathbf{P}'\mathbf{U}\mathbf{U}'\mathbf{P}) = \mathbf{P}' E(\mathbf{U}\mathbf{U}')\mathbf{P} = \mathbf{P}'\mathbf{I}\mathbf{P} = \mathbf{I}$.

Obviously $E(\mathbf{P}'\mathbf{U}) = \mathbf{P}' E(\mathbf{U}) = \mathbf{0}$. Therefore, $Z_i \sim N(0,1), i=1,2,\dots,n$ and independent

which implies that $\sum_{i=1}^r Z_i^2$ is $\chi^2(r)$.

R7. Let \mathbf{A} and \mathbf{B} be symmetric and idempotent matrices such that $\mathbf{AB} = \mathbf{0}$. Assume further that Y_1, Y_2, \dots, Y_n are independent normally distributed with expectation $\mu_1, \mu_2, \dots, \mu_n$

respectively and equal variance σ^2 . Let $\mathbf{Y} = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix}$ and $\underline{\mu} = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_n \end{bmatrix}$. Then $(\mathbf{Y} - \underline{\mu})' \mathbf{A} (\mathbf{Y} - \underline{\mu})$

and $(\mathbf{Y} - \underline{\mu})' \mathbf{B} (\mathbf{Y} - \underline{\mu})$ are independent.

Proof.

Let $\mathbf{U} = \mathbf{Y} - \underline{\mu}$. Then $E(\mathbf{U}) = \mathbf{0}$ and $\text{Cov}(\mathbf{U}) = \sigma^2 \mathbf{I}$. Further $\mathbf{U}'\mathbf{A}\mathbf{U} = \mathbf{U}'\mathbf{A}\mathbf{A}\mathbf{U} = \mathbf{Z}'\mathbf{Z}$ where $\mathbf{Z} = \mathbf{A}'\mathbf{U}$, and $\mathbf{U}'\mathbf{B}\mathbf{U} = \mathbf{U}'\mathbf{B}\mathbf{B}\mathbf{U} = \mathbf{V}'\mathbf{V}$ where $\mathbf{V} = \mathbf{B}'\mathbf{U}$. \mathbf{Z} and \mathbf{V} are independent if $E(\mathbf{Z}\mathbf{V}') = \mathbf{0}$.

$$E(\mathbf{Z}\mathbf{V}') = E(\mathbf{A}'\mathbf{U}\mathbf{U}'\mathbf{B}) = \mathbf{A} E(\mathbf{U}\mathbf{U}')\mathbf{B} = \mathbf{A}\sigma^2\mathbf{I}\mathbf{B} = \sigma^2\mathbf{AB} = \mathbf{0}.$$

Some specific results for multippel linear regression models.

Theorem 12.1

In the linear regression model $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} = \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix}$ with $E(\varepsilon_i \varepsilon_j) = \begin{cases} \sigma^2, & i=j \\ 0, & \text{elles} \end{cases}$ an estimator for

σ^2 is given by $\hat{\sigma}^2 = \sum_{i=1}^n \frac{(Y_i - \hat{Y}_i)^2}{n - (k+1)}$ where k is the number of regression variables.

Proof.

$$\sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = (\mathbf{Y} - \hat{\mathbf{Y}})'(\mathbf{Y} - \hat{\mathbf{Y}}) = \boldsymbol{\varepsilon}'\boldsymbol{\varepsilon}.$$

$$\boldsymbol{\varepsilon} = \mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{Y} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y} = (\mathbf{I} - \mathbf{H})\mathbf{Y} = (\mathbf{I} - \mathbf{H})(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})$$

since $(\mathbf{I} - \mathbf{H})\mathbf{X}\boldsymbol{\beta} = \mathbf{X}\boldsymbol{\beta} - \mathbf{X}\boldsymbol{\beta} = \mathbf{0}$.

Hence $\hat{\boldsymbol{\varepsilon}}^t \hat{\boldsymbol{\varepsilon}} = \boldsymbol{\varepsilon}^t (\mathbf{I} - \mathbf{H})(\mathbf{I} - \mathbf{H}) \boldsymbol{\varepsilon}$,

$$\begin{aligned} \text{and } E(\hat{\boldsymbol{\varepsilon}}^t \hat{\boldsymbol{\varepsilon}}) &= E(\boldsymbol{\varepsilon}^t (\mathbf{I} - \mathbf{H})(\mathbf{I} - \mathbf{H}) \boldsymbol{\varepsilon}) = E\left(\sum_{j=1}^n \sum_{i=1}^n \varepsilon_i (\mathbf{I} - \mathbf{H})_{ij} \varepsilon_j\right) = \sum_{i=1}^n (\mathbf{I} - \mathbf{H})_{ii} E(\varepsilon_i^2) \\ &= \sigma^2 \sum_{i=1}^n (\mathbf{I} - \mathbf{H})_{ii} = \sigma^2 \text{tr}(\mathbf{I} - \mathbf{H}) = \sigma^2 (\text{tr}(\mathbf{I}) - \text{tr}(\mathbf{H})) = \sigma^2 n - \sigma^2 \text{tr}\left(\mathbf{X}(\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t\right) \\ &= \sigma^2 n - \sigma^2 \text{tr}\left(\mathbf{X}^t \mathbf{X} (\mathbf{X}^t \mathbf{X})^{-1}\right) = \sigma^2 \left(n - \text{tr}(\mathbf{I}_{(k+1) \times (k+1)})\right) = (n - k - 1) \sigma^2. \end{aligned}$$

R8. In the linear regression model $\mathbf{Y} = \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\varepsilon} = \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix}$, with $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$,

$$\frac{SS_E}{\sigma^2} = \frac{(\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}})^t (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}})}{\sigma^2} \text{ is } \chi^2(n - k - 1).$$

Proof

$$(\mathbf{I} - \mathbf{H})(\mathbf{Y} - \mathbf{X} \boldsymbol{\beta}) = \mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}} - \mathbf{X} \boldsymbol{\beta} + \mathbf{X} \boldsymbol{\beta} = \mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}}.$$

Hence

$$(\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}})^t (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}}) = (\mathbf{Y} - \mathbf{X} \boldsymbol{\beta})^t (\mathbf{I} - \mathbf{H})^t (\mathbf{I} - \mathbf{H})(\mathbf{Y} - \mathbf{X} \boldsymbol{\beta}) = (\mathbf{Y} - \mathbf{X} \boldsymbol{\beta})^t (\mathbf{I} - \mathbf{H})(\mathbf{Y} - \mathbf{X} \boldsymbol{\beta})$$

$$\text{Rank}(\mathbf{I} - \mathbf{H}) = \text{tr}(\mathbf{I} - \mathbf{H}) = n - (k + 1).$$

$$\text{Hence } \frac{(\mathbf{Y} - \mathbf{X} \boldsymbol{\beta})^t (\mathbf{I} - \mathbf{H})(\mathbf{Y} - \mathbf{X} \boldsymbol{\beta})}{\sigma^2} \text{ is } \chi^2(n - (k + 1)).$$

Now assume $\beta_1 = \beta_2 = \dots = \beta_k = 0$ i.e. $Y_i = \beta_0 + \varepsilon_i$, $i = 1, 2, \dots, n$. Let $\mathbf{1}_{n \times 1} = [1, 1, \dots, 1]^t$

$$\text{Then } \left(\mathbf{I} - \frac{1}{n} \mathbf{J}\right)(\mathbf{Y} - \mathbf{X} \boldsymbol{\beta}) = \mathbf{Y} - \mathbf{X} \boldsymbol{\beta} - \bar{Y} \mathbf{1} + \frac{1}{n} \mathbf{J} \mathbf{X} \boldsymbol{\beta} = \mathbf{Y} - \bar{Y} \mathbf{1} - \mathbf{1} \beta_0 + \mathbf{1} \beta_0 = \mathbf{Y} - \bar{Y} \mathbf{1}.$$

$$\text{and } \left(\mathbf{H} - \frac{1}{n} \mathbf{J}\right)(\mathbf{Y} - \mathbf{X} \boldsymbol{\beta}) = \mathbf{X} \hat{\boldsymbol{\beta}} - \bar{Y} \mathbf{1} - \mathbf{X} \boldsymbol{\beta} + \frac{1}{n} \mathbf{J} \mathbf{X} \boldsymbol{\beta} = \hat{Y} \mathbf{1} - \bar{Y} \mathbf{1}.$$

$$\text{Therefore } \frac{SS_T}{\sigma^2} = \frac{\sum_{i=1}^n (Y_i - \bar{Y})^2}{\sigma^2} = \frac{(\mathbf{Y} - \bar{Y} \mathbf{1})^t (\mathbf{Y} - \bar{Y} \mathbf{1})}{\sigma^2} = \frac{(\mathbf{Y} - \mathbf{X} \boldsymbol{\beta})^t \left(\mathbf{I} - \frac{1}{n} \mathbf{J}\right)(\mathbf{Y} - \mathbf{X} \boldsymbol{\beta})}{\sigma^2} \text{ is } \chi^2(n - 1)$$

since $\text{tr}\left(\mathbf{I} - \frac{1}{n}\mathbf{J}\right) = n - \text{tr}\left(\frac{1}{n}\mathbf{J}\right) = n - 1$ and

$$\frac{SS_R}{\sigma^2} = \left(\frac{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}{\sigma^2} \right) = \frac{(\hat{\mathbf{Y}} - \bar{\mathbf{Y}})'(\hat{\mathbf{Y}} - \bar{\mathbf{Y}})}{\sigma^2} = \frac{(\mathbf{Y} - \mathbf{X}\beta)' \left(\mathbf{H} - \frac{1}{n}\mathbf{J} \right) (\mathbf{Y} - \mathbf{X}\beta)}{\sigma^2} \text{ is } \chi^2(k)$$

since $\text{rank}\left(\mathbf{H} - \frac{1}{n}\mathbf{J}\right) = \text{tr}(\mathbf{H}) - \text{tr}\left(\frac{1}{n}\mathbf{J}\right) = k + 1 - 1 = k$.

We also have $(\mathbf{I} - \mathbf{H})\left(\mathbf{H} - \frac{1}{n}\mathbf{J}\right) = \mathbf{0}$ from which it follows that SS_E and SS_R are independent.

R9. Assume $\varepsilon \sim N(0, \sigma^2 \mathbf{I})$. Then $\hat{\beta}$ is independent of $\hat{\sigma}^2$.

$$\begin{aligned} \text{Proof. Cov}(\hat{\beta}, \mathbf{Y} - \mathbf{X}\beta) &= E\left((\hat{\beta} - \beta), (\mathbf{Y} - \mathbf{H}\mathbf{Y})'\right) = E\left((\hat{\beta} - \beta), \varepsilon'\right) \\ &= E\left((\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'(\mathbf{Y} - \mathbf{X}\beta) \varepsilon' (\mathbf{I} - \mathbf{H})\right) = E\left((\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \varepsilon \varepsilon' (\mathbf{I} - \mathbf{H})\right) \\ &= \sigma^2 (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \mathbf{I} (\mathbf{I} - \mathbf{H}) = \sigma^2 \left((\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \right) - (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' = \mathbf{0}. \end{aligned}$$

The Fisher (F) distribution and the analysis of variance table for linear regression.

R10. Assume the random variable X is chi-squared distributed with ν_1 degrees of freedom, writes $\chi^2(\nu_1)$, and $Y \sim \chi^2(\nu_2)$ and that X and Y are independent. Then $\frac{X/\nu_1}{Y/\nu_2}$ is Fisher distributed with ν_1 and ν_2 degrees of freedom, writes F_{ν_1, ν_2} .

Note

$$F \sim F_{\nu_1, \nu_2} \Rightarrow \frac{1}{F} \sim F_{\nu_2, \nu_1}$$

Now define α by

$$P(F_{\nu_1, \nu_2} \leq f_{\alpha, \nu_1, \nu_2}) = \alpha$$

and since

$$P\left(f_{1-\frac{\alpha}{2}, \nu_1, \nu_2} \leq F \leq f_{\frac{\alpha}{2}, \nu_1, \nu_2}\right) = 1 - \alpha \Leftrightarrow P\left(\frac{1}{f_{\frac{\alpha}{2}, \nu_1, \nu_2}} \leq \frac{1}{F} \leq \frac{1}{f_{1-\frac{\alpha}{2}, \nu_1, \nu_2}}\right) = 1 - \alpha$$

we have $\frac{1}{f_{\frac{\alpha}{2}, \nu_1, \nu_2}} = f_{1-\frac{\alpha}{2}, \nu_2, \nu_1}$ and $f_{1-\frac{\alpha}{2}, \nu_1, \nu_2} = \frac{1}{f_{\frac{\alpha}{2}, \nu_2, \nu_1}}$

The density function in the F distribution is given by:

$$f(x) = k_{\nu_1, \nu_2} \cdot x^{\frac{\nu_1}{2}-1} \left(1 + \frac{\nu_1}{\nu_2} x\right)^{-\frac{\nu_1 + \nu_2}{2}}, x > 0$$

where $k_{\nu_1, \nu_2} = \frac{\Gamma\left(\frac{\nu_1 + \nu_2}{2}\right)}{\Gamma\left(\frac{\nu_1}{2}\right) \cdot \Gamma\left(\frac{\nu_2}{2}\right)} \cdot \left(\frac{\nu_1}{\nu_2}\right)^{\frac{\nu_1}{2}}$

$$E(F) = \frac{\nu_1}{\nu_2 - 2} \text{ and } Var(F_{\nu_1, \nu_2}) = \frac{2\nu_2^2(\nu_1 + \nu_2 - 2)}{\nu_1(\nu_2 - 2)^2(\nu_2 - 4)} \text{ provided } \nu_2 > 4.$$

Some examples of how the density function in the Fisher distribution looks like is given in Figure 6.

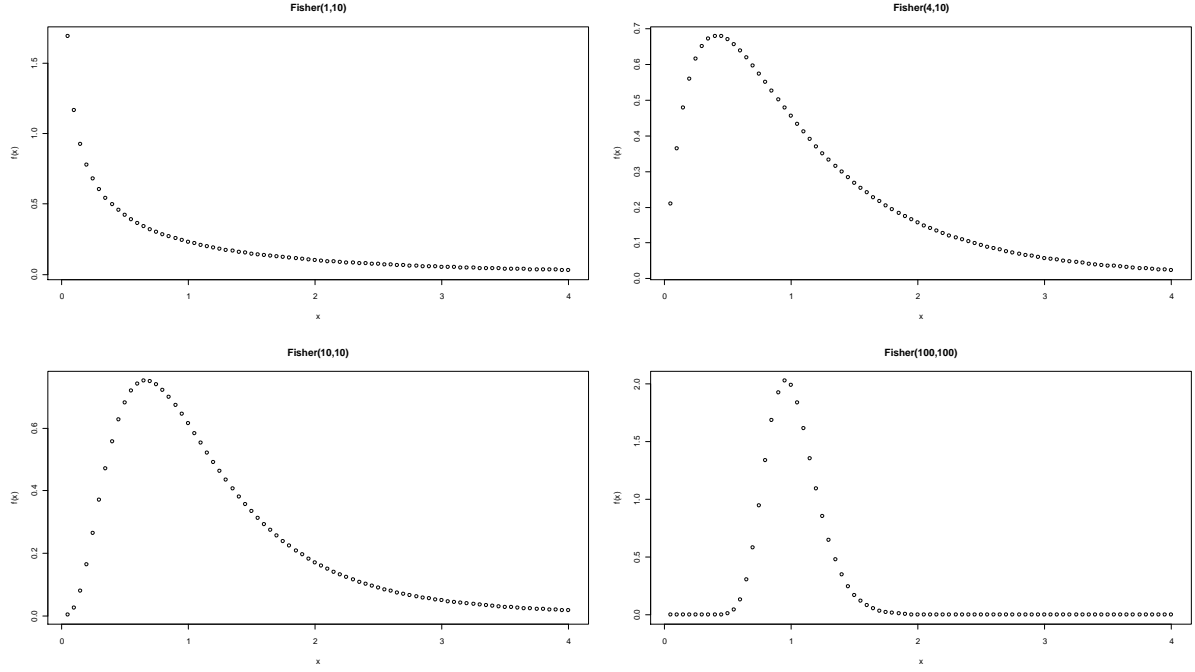


Figure 6.

Inference in multipel linear regression

A test on whether the regression is significant (or if the regression equation differs from a constant) is given by:

$H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$ against $H_1: \text{at least one } \beta_i, i=1,2,\dots,k, \text{ is different from } 0.$

With $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$, we know that $\frac{SS_E}{\sigma^2}$ is $\chi^2(n-k-1)$. Under H_0 we also have that

$\frac{SS_R}{\sigma^2}$ is $\chi^2(k)$ and since they are independent, we get $\frac{SS_R}{k} \bigg/ \frac{SS_E}{n-k-1} \sim F_{k,n-k-1}$. Since the denominator is independent of H_0 and SS_R under H_1 is greater than under H_0 a natural test for H_0 can be built on $F_{k,n-k-1}$.

We reject H_0 if $f \geq f_{\alpha,k,n-(k+1)}$ where f is the observed value of $F_{k,n-k-1}$.

It is normal to collect the sum of squares values from a regression analysis in a variance of analysis table as follows:

Table. The analysis of variance table for regression analysis

Source	Sum of squares	Df	Mean sum of squares	F
Regression	$SS_R = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$	k	$MSS_R = \frac{SS_R}{k}$	$f = \frac{MSS_R}{MSS_E}$
Error	$SS_E = \sum_{i=1}^n (y_i - \hat{y}_i)^2$	$n - (k+1)$	$MSS_E = \frac{SS_E}{n - (k+1)}$	
Total	$SS_T = \sum_{i=1}^n (y_i - \bar{y})^2$	$n - 1$		

The partial F-test.

The partial F-test is useful if we want to test the significance of a group of regression variables. Assume we want to test the following hypothesis.

$$H_0: \beta_{r+1}, \beta_{r+2}, \dots, \beta_k = 0 \quad H_1: \text{at least one of } \beta_{r+1}, \dots, \beta_k \text{ are } \neq 0.$$

Let $X = [\mathbf{1}, x_1, x_2, \dots, x_r, \dots, x_k]$ and $X_1 = [\mathbf{1}, x_1, x_2, \dots, x_r]$ and define H and H_1 by

$$H = X(X^t X)^{-1} X^t \text{ and } H_1 = X_1(X_1^t X_1)^{-1} X_1^t. \text{ We get for all } y:$$

$$HH_1 y = HX_1 b_1 = HX \begin{bmatrix} b_1 \\ 0 \end{bmatrix} = X \begin{bmatrix} b_1 \\ 0 \end{bmatrix} = X_1 b_1 = H_1 y. \text{ Thereby } HH_1 = H_1 \text{ and by transposing}$$

it we get $H_1 H = H_1$. Therefore $(H - H_1)^2 = H - HH_1 - H_1 H + H_1^2 = H - H_1$ i.e. $(H - H_1)$ is idempotent.

$$\text{Let } \beta = \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_k \end{bmatrix} \text{ and } \beta_1 = \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_r \end{bmatrix}.$$

$$SS_r(\beta) - SS_r(\beta_1) = SS_T - Y^t(I - H)Y - SS_T + Y^t(I - H_1)Y = Y^t(H - H_1)Y = (Y - X_1 \beta_1)^t (H - H_1)(Y - X_1 \beta_1)$$

since $(H - H_1)X_1 \beta_1 = X_1 \beta_1 - X_1 \beta_1$.

$$\text{Rank}(H - H_1) = \text{tr}(H) - \text{tr}(H_1) = k + 1 - (r + 1) = k - r. \text{ Therefore under } H_0$$

$$\frac{SS_r(\beta) - SS_r(\beta_1)}{\sigma^2} = \frac{(Y - X_1 \beta_1)^t (H - H_1)(Y - X_1 \beta_1)}{\sigma^2} \text{ is } \chi^2(k - r).$$

Also $(\mathbf{H} - \mathbf{H}_1)(\mathbf{I} - \mathbf{H}) = \mathbf{H} - \mathbf{H}_1 - \mathbf{H} + \mathbf{H}_1 = \mathbf{0}$ which implies $\mathbf{Y}'(\mathbf{H} - \mathbf{H}_1)\mathbf{Y}$ and $\hat{\sigma}^2$ are independent.

$$\text{The test statistic for } H_0 \text{ is } F = \frac{\frac{SS_r(\boldsymbol{\beta}) - SS_r(\boldsymbol{\beta}_1)}{k-r}}{\frac{SS_E(\boldsymbol{\beta})}{n-(k+1)}} \sim F_{(k-r), (n-k-1)},$$

and we reject H_0 if $f_{obs} \geq f_{\alpha, (k-r), (n-k-1)}$. Note that a partial F-test can be used to test on single variables given that the rest of the variables are in the model, on groups of variables given that the rest of the variables are in the model and if the regression is significant.

Testing a general linear hypothesis

Suppose $E[Y] = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$ and we want to test the hypothesis

$$H_0: \begin{bmatrix} \beta_1 - \beta_2 \\ \beta_2 - \beta_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \text{ or equivalently } H_0: \beta_1 = \beta_2 = \beta_3 = \beta.$$

The hypothesis can be formulated as: $H_0: \mathbf{C}\boldsymbol{\beta} = \mathbf{0}$ vs $H_1: \mathbf{C}\boldsymbol{\beta} \neq \mathbf{0}$, where

$$\mathbf{C} = \begin{bmatrix} 0 & 1 & -1 & 0 \\ 0 & 0 & 1 & -1 \end{bmatrix} \text{ and } \boldsymbol{\beta} = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \end{bmatrix}. \text{ We have thus } r \text{ (in this case 2) constraints on } \boldsymbol{\beta}.$$

With the assumption of normally distributed errors,

$$\hat{\boldsymbol{\beta}} \sim N(\boldsymbol{\beta}, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}) \text{ and } \mathbf{C}\hat{\boldsymbol{\beta}} \sim N(\mathbf{C}\boldsymbol{\beta}, \sigma^2 \mathbf{C}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{C}^T).$$

Hence, under H_0 , $(\mathbf{C}\hat{\boldsymbol{\beta}})^T (\sigma^2 \mathbf{C}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{C}^T)^{-1} (\mathbf{C}\hat{\boldsymbol{\beta}}) \sim \chi^2(r)$.

Since $\hat{\boldsymbol{\beta}}$ and SS_E (full model) are independent and $\frac{SS_E}{\sigma^2} \sim \chi^2(n-k-1)$, we get

$$F = \frac{(\mathbf{C}\hat{\boldsymbol{\beta}})^T (\mathbf{C}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{C}^T)^{-1} (\mathbf{C}\hat{\boldsymbol{\beta}}) / \sigma^2 r}{SS_E / \sigma^2 (n-k-1)} = \frac{(\mathbf{C}\hat{\boldsymbol{\beta}})^T (\mathbf{C}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{C}^T)^{-1} (\mathbf{C}\hat{\boldsymbol{\beta}}) / r}{SS_E / (n-k-1)} \sim F_{r, n-k-1}.$$

We reject H_0 if F is large.

If we want to test $H_0: \mathbf{C}\boldsymbol{\beta} = \mathbf{d}$ vs $H_1: \mathbf{C}\boldsymbol{\beta} \neq \mathbf{d}$, we just substitute $\mathbf{C}\hat{\boldsymbol{\beta}}$ with $\mathbf{C}\hat{\boldsymbol{\beta}} - \mathbf{d}$ in F .

Test on single coefficients, intervals and regions.

We will now be mainly concerned with how regression analysis can be performed in practice. Assume that $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$ is a sequence of independent, normally distributed random variables with expectation equal to zero and constant variance σ^2 .

In multiple linear regression we have:

$$E(Y) = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k.$$

Thus β_i , $i = 1, 2, \dots, k$ is the change in $E(Y)$ if x_i is changed with one unit and all the rest of the variables are kept unchanged.

We want a test for: Does a variable have a significant impact on the response given that the other variables are in the model.

Such a test for x_j , $j = 1, 2, \dots, k$ is:

$$H_0: \beta_j = 0 \text{ against } H_1: \beta_j \neq 0$$

We know that $\hat{\beta}$ is independent of $S^2 = \frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{n - (k+1)}$.

Now let c_{jj} be the $(j+1)$ th diagonal element in $(\mathbf{X}^t \mathbf{X})^{-1}$. Then

$$T = \frac{\hat{\beta}_j}{S \sqrt{c_{jj}}} = \frac{\frac{\hat{\beta}_j}{\sigma \sqrt{c_{jj}}}}{\frac{S}{\sigma}} \text{ and hence the numerator is } N(0,1) \text{ and since } \frac{S^2 (n-k-1)}{\sigma^2} \text{ is}$$

$\chi^2 (n-k-1)$ the denominator can be written as a variable that is: $\sqrt{\frac{\chi^2 (n-k-1)}{n-k-1}}$. Therefore T is t-distributed with $(n-k-1)$ degrees of freedom and we reject H_0 if $|t_{obs}| \geq t_{\frac{\alpha}{2}, n-(k+1)}$. For the test

$$H_0: \beta_j = \beta_{j0} \text{ against } H_1: \beta_j = \beta_{j0}, \text{ we use } T = \frac{\hat{\beta}_j - \beta_{j0}}{S \sqrt{c_{jj}}} \text{ the same way except that } \frac{\alpha}{2} \text{ is}$$

substituted by α if the test is one-sided.

Individual confidence intervals for each β_j , $j = 1, 2, \dots, k$, can be deduced using T .

We get $P\left(\hat{\beta}_j - t_{\frac{\alpha}{2}, n-k-1} S\sqrt{c_{jj}} \leq \beta_j \leq \hat{\beta}_j + t_{\frac{\alpha}{2}, n-k-1} S\sqrt{c_{jj}}\right) = 1 - \alpha$.

The confidence ellipsoid

A confidence ellipsoid may be useful if several parameters are to be studied jointly.

Exploiting that $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}) \Rightarrow \hat{\boldsymbol{\beta}} \sim N(\boldsymbol{\beta}, \sigma^2 (\mathbf{X}'\mathbf{X})^{-1})$ we get

$(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})' \frac{\mathbf{X}'\mathbf{X}}{\sigma^2} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \sim \chi^2(k+1)$. We also know that $\frac{SS_E}{\sigma^2} \sim \chi^2(n-k-1)$ and

$\frac{SS_E}{n-k-1} = S^2$. Therefore, $\frac{S^2(n-k-1)}{\sigma^2} \sim \chi^2(n-k-1)$. Thereby

$$\frac{(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})' \mathbf{X}'\mathbf{X} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})}{\frac{(n-k-1)S^2}{\sigma^2(n-k-1)}} \text{ and a } 100(1-\alpha)\% \text{ confidence region for } \boldsymbol{\beta} \text{ is given by those } \boldsymbol{\beta}$$

that satisfy $(\mathbf{b} - \boldsymbol{\beta})' \frac{\mathbf{X}'\mathbf{X}}{\sigma^2} (\mathbf{b} - \boldsymbol{\beta}) \leq (k+1)s^2 f_{\alpha, k+1, n-k-1}$. Here \mathbf{b} and s^2 are estimates.

Confidence intervals for expected response and predictions.

We will now construct a confidence interval for the expected response given values $x_{10}, x_{20}, \dots, x_{k0}$ for x_1, x_2, \dots, x_k .

Since $\hat{\boldsymbol{\beta}}$ is independent of S^2 we have that $\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$ also must be independent of S^2 .

Further

$$E(Y|x_{10}, x_{20}, \dots, x_{k0}) = \beta_0 + \beta_1 x_{10} + \dots + \beta_k x_{k0} = [1, x_{10}, x_{20}, \dots, x_{k0}] \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix} = \mathbf{x}_0' \boldsymbol{\beta} = \mu_{Y|x_{10}, \dots, x_{k0}}. \text{ Also}$$

$$\hat{\mu}_{Y|x_{10}, \mathbf{K} x_{k0}} = \mathbf{x}_0' \hat{\boldsymbol{\beta}} = \hat{Y}_0.$$

$$\text{Var}(\hat{Y}_0) = E\left[\mathbf{x}_0' (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})' \mathbf{x}_0\right] = \mathbf{x}_0' E\left[(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})'\right] \mathbf{x}_0 = \sigma^2 \mathbf{x}_0' (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_0$$

Since $\hat{\boldsymbol{\beta}}$ and S^2 are independent we get:

$$T = \frac{\hat{Y}_0 - \mu_{Y|x_{10}, x_{20}, \dots, x_{k0}}}{S \sqrt{\mathbf{x}_0^t (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{x}_0}} \text{ is t-distributed with } (n-k-1) \text{ degrees of freedom.}$$

Hence a $100(1-\alpha)\%$ confidence interval for expected response when $\mathbf{x} = \mathbf{x}_0$ is:

$$\hat{y}_0 \pm t_{\frac{\alpha}{2}, n-k-1} S \sqrt{\mathbf{x}_0^t (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{x}_0}$$

Prediction interval for a new value of the response when $\mathbf{x} = \mathbf{x}_0$.

Let Y_0 be the new value. Then

$$E[Y_0 - \hat{Y}_0] = E[\mathbf{x}_0^t \boldsymbol{\beta} + \varepsilon] - E[\mathbf{x}_0^t \boldsymbol{\beta}] = \mathbf{x}_0^t \boldsymbol{\beta} - \mathbf{x}_0^t \boldsymbol{\beta} = 0.$$

$$\text{Var}[Y_0 - \hat{Y}_0] = \text{Var}[Y_0] + \text{Var}[\hat{Y}_0] = \sigma^2 + \sigma^2 \mathbf{x}_0^t (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{x}_0$$

such that

$$T = \frac{\hat{Y}_0 - Y_0}{S \sqrt{1 + \mathbf{x}_0^t (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{x}_0}}$$

is t-distributed with $(n-k-1)$ degrees of freedom. Hence a $100(1-\alpha)\%$ prediction interval

for Y given $\mathbf{x} = \mathbf{x}_0$ is:

$$\hat{y}_0 \pm t_{\frac{\alpha}{2}, n-k-1} S \sqrt{\left(1 + \mathbf{x}_0^t (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{x}_0\right)}.$$

Choice of a fitted model

Let Y be the response and x_1, x_2, \dots, x_k be the explanatory variable. Explanatory variables with small or no influence on the response can give the model a bad prediction ability.

Criteria for evaluating the adequacy of a model

We have

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \text{ or } SS_T = SS_E + SS_R.$$

In a good model $y_i - \hat{y}_i$ should be small and $\hat{y}_i - \bar{y} \approx y_i - \bar{y}, i = 1, 2, \dots, n$. A measure of how much variation in the data that can be explained by the model is the coefficient of multiple determination:

$$R^2 = \frac{SS_R}{SS_T} = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}.$$

We have $0 \leq R^2 \leq 1$. $R^2 = 0.84$ tells us that 84% of the variation in the data can be explained by the model. Note that $R^2 = \frac{SS_R}{SS_T} = \frac{SS_T - SS_E}{SS_T} = 1 - \frac{SS_E}{SS_T}$.

A problem with R^2 is that it will always increase when we increase the number of explanatory variables. Therefore, $R^2_{adjusted}$ has been introduced defined as

$$R^2_{adjusted} = 1 - \frac{\frac{SS_E}{n - (k + 1)}}{\frac{SS_T}{n - 1}} = 1 - \frac{(n - 1)s^2}{SS_T}.$$

Therefore maximizing $R^2_{adjusted}$ is equivalent to minimizing s^2 .

An example. Sour precipitation

In order to study the influence of sour precipitation in Norwegian lakes, SFT in 1986 carried through an investigation where data from 1005 lakes were collected. In this example 26 random lakes, 16 from Telemark and 10 from Trøndelag is chosen out.

y = Measured pH-value

x_1 = Content of SO_4

x_2 = Content of NO_3

x_3 = Content of Ca

x_4 = Content of latent aluminum

x_5 = Content of organic material

x_6 = Area of water

x_7 = Location (0 – Telemark, 1-Trøndelag)

The observed data is given below:

	y	x1	x2	x3	x4	x5	x6	x7
1	5.38	4.9	39	1.54	78	2.02	0.30	0
2	5.68	4.1	75	1.55	17	2.98	1.85	0

3	5.04	3.5	80	0.83	157	3.40	0.25	0
4	4.81	3.8	75	0.53	163	3.42	0.30	0
5	4.92	3.8	90	0.82	105	3.41	0.25	0
6	5.34	2.6	49	0.62	114	1.90	0.65	0
7	5.74	2.7	79	1.08	15	2.53	1.15	0
8	5.17	2.6	90	0.67	90	1.89	0.91	0
9	5.02	2.4	64	0.41	107	0.97	0.60	0
10	5.88	2.8	27	1.15	12	3.04	0.37	0
11	5.36	3.4	13	0.89	93	2.95	0.58	0
12	5.26	2.7	14	0.74	72	2.75	0.15	0
13	5.69	3.2	13	1.03	99	3.34	0.53	0
14	5.51	2.5	79	0.67	80	0.78	0.56	0
15	5.25	1.5	77	0.33	72	0.11	1.08	0
16	6.06	3.7	15	1.94	4	5.53	0.72	0
17	6.08	1.9	16	1.05	0	5.63	0.48	1
18	6.08	1.3	2	0.81	0	3.84	0.25	1
19	6.20	2.2	32	1.40	7	3.77	1.87	1
20	5.64	2.2	21	0.75	13	5.62	0.53	1
21	5.75	1.6	7	0.79	2	3.95	0.35	1
22	5.43	2.0	27	0.47	16	2.39	1.03	1
23	5.82	1.8	17	0.74	3	3.08	0.32	1
24	5.50	2.0	6	0.49	10	2.82	0.28	1
25	5.62	1.5	3	0.36	7	1.71	0.10	1
26	5.41	1.7	7	0.54	9	6.48	0.21	1

In R we have the possibility of plotting pairs of columns like this:

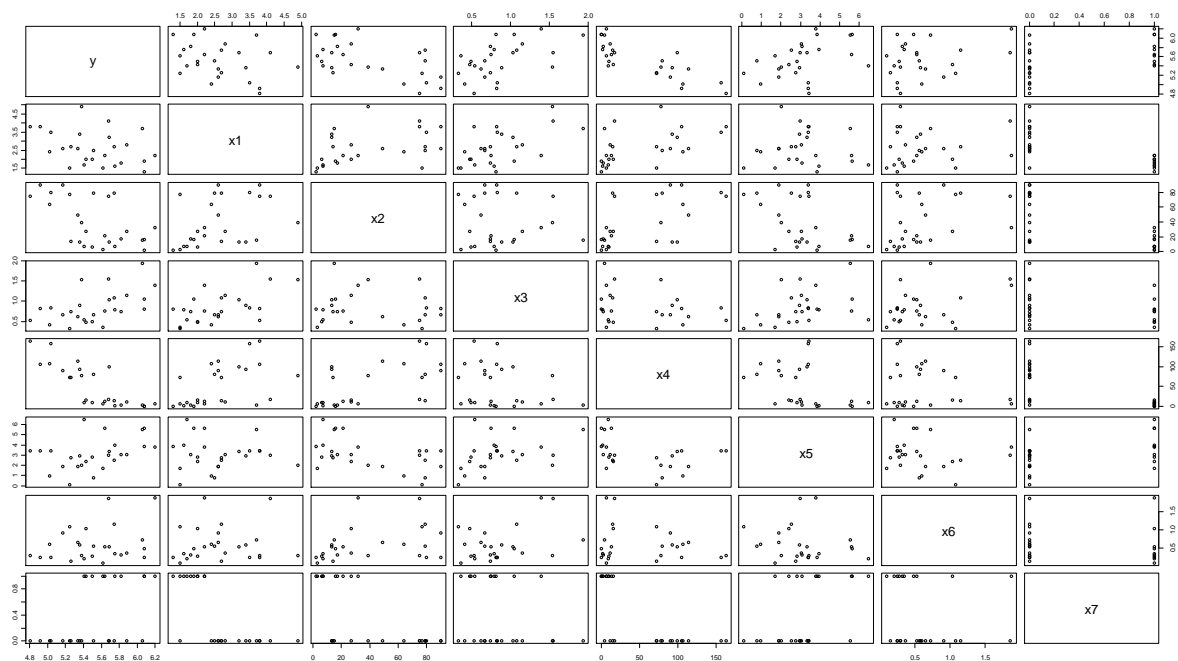


Figure7. In this figure the response is plotted against each regression variable and each regression variable is plotted against other regression variables

From the figure it seems like there might be a relationship between y and x_2, x_3, x_4 and x_5 .

We can also calculate the correlation between all pairs of two variables. The correlation between

Y and x_1 is calculated as: $r_{Y, x_1} = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_{1i} - \bar{x}_1)}{\sqrt{(y_i - \bar{y})^2 (x_{1i} - \bar{x})^2}}$, for instance

	y	x1	x2	x3	x4	x5	x6	x7
y	1.0000000	-0.33530173	-0.5703950	0.5263695	-0.8128681	0.40547721	0.2764476	0.4930767
x1	-0.3353017	1.0000000	0.4090205	0.5830377	0.5144954	-0.02402104	0.0891068	-0.7029792
x2	-0.5703950	0.40902048	1.0000000	-0.0585331	0.5973926	-0.47822779	0.3481284	-0.6470683
x3	0.5263695	0.58303773	-0.0585331	1.0000000	-0.2453771	0.36579209	0.3744328	-0.2292290
x4	-0.8128681	0.51449542	0.5973926	-0.2453771	1.0000000	-0.39443507	-0.1857614	-0.6966124
x5	0.4054772	-0.02402104	-0.4782278	0.3657921	-0.3944351	1.0000000	-0.1519310	0.4444162
x6	0.2764476	0.08910679	0.3481284	0.3744328	-0.1857614	-0.15193098	1.0000000	-0.1045081
x7	0.4930767	-0.70297917	-0.6470683	-0.2292290	-0.6966124	0.44441622	-0.1045081	1.0000000

It is clear that the variable that is the most correlated with the response is x_4 , but there is also some correlation between Y and x_2, x_3, x_5 and x_7 .

Using x_4 as the single regression variable we get the following output from R:

```

              Estimate Std. Error  t value Pr(>|t|)
(Intercept)   5.8259182   0.0619470   94.047   < 2e-16 ***
x4            -0.0058244   0.0008519   -6.837   4.52e-07 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

```

Residual standard error: 0.222 on 24 degrees of freedom
Multiple R-squared:  0.6608,    Adjusted R-squared:  0.6466
F-statistic: 46.75 on 1 and 24 DF,  p-value: 4.521e-07

```

With the following analysis of variance table

	x4	Residuals
Sum of Squares	2.303553	1.182693
Deg. of Freedom	1	24

Residual standard error: 0.2219885

According to the t-test, the variable x_4 and the constant term is highly significant. Further we get that

$s = 0.222$ or 0.2219885 and $24s^2 = 1.182693 = SS_E$. $R^2 = \frac{SS_R}{SS_T} = 0.6608$ and $R^2_{adj} = 0.6466$. $SS_R = 2.303553$ and $SS_T = 2.303553 + 1.182693 = 3.486245$

If we try to estimate a model with x_2, x_3, x_4, x_5 and x_7 as the regression variables we get the following results:

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	5.444654	0.205522	26.492	< 2e-16 ***
x2	-0.002140	0.001620	-1.321	0.20139
x3	0.468389	0.139062	3.368	0.00306 **
x4	-0.003517	0.001256	-2.800	0.01107 *
x5	-0.038055	0.033321	-1.142	0.26691
x7	0.164354	0.155003	1.060	0.30163

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

Residual standard error: 0.1815 on 20 degrees of freedom
Multiple R-squared: 0.8111, Adjusted R-squared: 0.7639
F-statistic: 17.18 on 5 and 20 DF, p-value: 1.237e-06

We observe that with these variables in the model, x_3 and x_4 seem to be significant.

Finally, with all the variables in the model x_1 and x_3 seem to be the significant ones while x_4 is not even close to being significant.

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	5.6764334	0.1389162	40.862	< 2e-16 ***
x1	-0.3150444	0.0587512	-5.362	4.27e-05 ***
x2	-0.0018533	0.0012587	-1.472	0.158
x3	0.9751745	0.1449075	6.730	2.62e-06 ***
x4	-0.0002268	0.0010038	-0.226	0.824
x5	-0.0334242	0.0225009	-1.485	0.155
x6	-0.0039399	0.0724339	-0.054	0.957
x7	0.0888722	0.1025724	0.866	0.398

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

Residual standard error: 0.1165 on 18 degrees of freedom
Multiple R-squared: 0.93, Adjusted R-squared: 0.9027
F-statistic: 34.15 on 7 and 18 DF, p-value: 3.904e-09

This example shows that it is not an easy task to find the model that best explains the variation in the response. Fortunately, there are other methods than the ones above that are helpful when doing model selection. First, we will show some procedures based on significant tests that are useful to select variables to be in the model.

Variable selection methods

The most common variable selection methods are forward selection, backward elimination, and stepwise regression. Normally these algorithms are performed as partial F-tests, adding or removing one variable at the time. We will now give an algorithmic presentation of these methods.

Forward selection

1. Start with only β_0 in the model.
2. Find $\max_j R(\beta_j) = \max_j SS_R(\beta_j | \beta_0) = \max_j \{SS_R(\beta_0, \beta_j) - SS_R(\beta_0)\}$.
3. If $\max_j \frac{R(\beta_j)}{\frac{SS_E}{n-2}} = \frac{R(\beta_m)}{\frac{SS_E}{n-2}} < f_{\alpha, 1, n-2}$ stop, no variable is entered into the model.
4. If $\frac{R(\beta_m)}{\frac{SS_E}{n-2}} \geq f_{\alpha, 1, n-2}$ add x_m to the model.

Find $\max_{j \neq m} R(\beta_j | \beta_m) = \max_{j \neq m} \{SS_R(\beta_0, \beta_m, \beta_j) - SS_R(\beta_0, \beta_m)\}$ and go to step 3. If

$\max_{j \neq m} \frac{R(\beta_j | \beta_m)}{\frac{SS_E}{n-3}} < f_{\alpha, 1, n-3}$, no more variables are entered into the model. Otherwise proceed in

the same way. Note that the degrees of freedom in the partial F-test is reduced by one for each variable that is entered into the model.

Backward elimination

Define $\beta \setminus \beta_j = (\beta_0, \dots, \beta_{j-1}, \beta_{j+1}, \dots, \beta_k)$.

1. Start with all the variables in the model.
2. Find $\min_j R(\beta_j | \beta \setminus \beta_j) = \min_j \{SS_R(\beta_0, \dots, \beta_k) - SS_R(\beta_0, \dots, \beta_{j-1}, \beta_{j+1}, \dots, \beta_k)\}$.
3. If $\min_j \frac{R(\beta_j | \beta \setminus \beta_j)}{\frac{SS_E}{n-k-1}} = \frac{R(\beta_m | \beta \setminus \beta_m)}{\frac{SS_E}{n-k-1}} > f_{\alpha, 1, n-k-1}$ stop, no variable is removed from the model.
4. If $\frac{R(\beta_m | \beta \setminus \beta_m)}{\frac{SSE}{n-k-1}} < f_{\alpha, 1, n-k-1}$, remove x_m .

Find $\min_{j \neq m} R(\beta_j | \beta \setminus \{\beta_m, \beta_j\}) = \min_{j \neq m} \{SS_R(\beta_0, \dots, \beta_{m-1}, \beta_{m+1}, \dots, \beta_k) - SS_R(\beta \setminus \{\beta_m, \beta_j\})\}$

and go to step 3.

If $\frac{\min_{j \neq m} R(\beta_j | \beta \setminus \{\beta_m, \beta_j\})}{\frac{SS_E}{n-k}} > f_{\alpha, 1, n-k}$ no more variables are taken out of the model.

Otherwise proceed in the same way. Note that the degrees of freedom in the partial F-test is increased by one in each step.

Stepwise regression

This method combines forward selection and backward elimination.

1. Start as with forward selection. Assume x_n and x_m are chosen as the first two variables to enter the model.
2. Find $\min_{j=n,m} R(\beta_j | \{\beta_n, \beta_m\} \setminus \beta_j)$ and check if one of the variables x_n or x_m can be taken out as with backward selection.
3. Continue as with forward selection but check in each step if variables that are chosen to be in the model in earlier steps can be taken out.

Since an F-test where the numerator has one degrees of freedom is equivalent to a t-test for evaluating the significance of a parameter, t-tests could be used instead of F-tests. In R these three selection methods are performed based on criteria and not on significance evaluation. These criteria penalize large models. Default in R is the Akaike's information criterion given as:

$$AIC = n \ln \left(\frac{SS_E}{n} \right) + 2(k+2)$$

and small values for AIC are to be preferred. Note that now also σ^2 is considered a parameter, therefore $k+2$. An alternative that penalizes large models more is the Bayesian information criterion:

$$BIC = n \ln \left(\frac{SS_E}{n} \right) + \ln(n)(k+2).$$

Other methods are best subset regression, where one seek through the set of regression variables and for each s , $s = 1, 2, \dots, k$ find the subset of s regression variables that gives the best fit to the response measured in terms of criteria such as R^2 , $R^2_{adjusted}$, AIC, BIC or Mallows's- C_p .

Mallows- C_p

The idea behind Mallows- C_p is that most (all) models are wrong and therefore

$E(Y_i) \neq E(\hat{Y}_i)$. Now writing

$$\hat{Y}_i - E(Y_i) = \hat{Y}_i - E(\hat{Y}_i) + E(\hat{Y}_i) - E(Y_i),$$

we get

$$E\left(\hat{Y}_i - E(Y_i)\right)^2 = E\left(\hat{Y}_i - E(\hat{Y}_i)\right)^2 + E\left(E(\hat{Y}_i) - E(Y_i)\right)^2 + 2E\left[\left(\hat{Y}_i - E(\hat{Y}_i)\right)\left(E(\hat{Y}_i) - E(Y_i)\right)\right].$$

The last term is zero and we obtain

$$E\left[\left(\hat{Y}_i - E(Y_i)\right)^2\right] = \text{Var}(\hat{Y}_i) + \left(\text{Bias}\hat{Y}_i\right)^2$$

where $\text{Bias}\hat{Y}_i = E(\hat{Y}_i) - E(Y_i)$.

We want to minimize

$$E\left[\sum_{i=1}^n \frac{\left(\hat{Y}_i - E(Y_i)\right)^2}{\sigma^2}\right] = \sum_{i=1}^n \frac{\text{Var}(\hat{Y}_i)}{\sigma^2} + \sum_{i=1}^n \frac{\left(\text{Bias}\hat{Y}_i\right)^2}{\sigma^2}.$$

Now $\sum_{i=1}^n \text{Var}(\hat{Y}_i) = \sigma^2 \sum_{i=1}^n h_{ii} = \sigma^2 (k+1)$. An estimate for Mallows- C_p is given by:

$$C_p = k+1 + \frac{(s^2 - \hat{\sigma}^2)(n-k-1)}{\hat{\sigma}^2}.$$

Correct model has $C_p = k+1 = p$.

s^2 is an estimate for $\text{Var}(Y_i)$ in the model with p -parameters and $\hat{\sigma}^2$ is an estimate for σ^2 . Often it is the estimate of σ^2 from the full model.

As an example, forward regression is applied to the data set about sour precipitation using the default algorithm in R.

```
>step(lm(y~1), y~1+x1+x2+x3+x4+x5+x6+x7, direction="forward")
Start:  AIC=-50.24
y ~ 1
```

	Df	Sum of Sq	RSS	AIC
+ x4	1	2.30355	1.1827	-76.348
+ x2	1	1.13425	2.3520	-58.474
+ x3	1	0.96592	2.5203	-56.676
+ x7	1	0.84759	2.6387	-55.484
+ x5	1	0.57318	2.9131	-52.911
+ x1	1	0.39195	3.0943	-51.342
+ x6	1	0.26643	3.2198	-50.308
<none>			3.4862	-50.241

```
Step:  AIC=-76.35
y ~ x4
```

	Df	Sum of Sq	RSS	AIC
--	----	-----------	-----	-----

```

+ x3      1    0.39645 0.78625 -84.963
<none>                1.18269 -76.348
+ x6      1    0.05682 1.12587 -75.628
+ x2      1    0.03898 1.14372 -75.219
+ x7      1    0.03627 1.14642 -75.158
+ x1      1    0.03260 1.15010 -75.075
+ x5      1    0.02973 1.15297 -75.010

```

Step: AIC=-84.96

y ~ x4 + x3

	Df	Sum of Sq	RSS	AIC
+ x1	1	0.46780	0.31844	-106.463
+ x2	1	0.07316	0.71309	-85.502
<none>			0.78625	-84.963
+ x7	1	0.04413	0.74212	-84.465
+ x6	1	0.00050	0.78575	-82.980
+ x5	1	0.00035	0.78590	-82.975

Step: AIC=-106.46

y ~ x4 + x3 + x1

	Df	Sum of Sq	RSS	AIC
+ x2	1	0.041571	0.27687	-108.10
<none>			0.31844	-106.46
+ x7	1	0.010093	0.30835	-105.30
+ x6	1	0.007048	0.31140	-105.04
+ x5	1	0.003697	0.31475	-104.77

Step: AIC=-108.1

y ~ x4 + x3 + x1 + x2

	Df	Sum of Sq	RSS	AIC
+ x5	1	0.0223623	0.25451	-108.29
<none>			0.27687	-108.10
+ x6	1	0.0025490	0.27432	-106.34
+ x7	1	0.0003829	0.27649	-106.14

Step: AIC=-108.29

y ~ x4 + x3 + x1 + x2 + x5

	Df	Sum of Sq	RSS	AIC
<none>			0.25451	-108.29
+ x7	1	0.0103492	0.24416	-107.37
+ x6	1	0.0002081	0.25430	-106.31

Call:

lm(formula = y ~ x4 + x3 + x1 + x2 + x5)

Coefficients:

(Intercept)	x4	x3	x1	x2	x5
5.7668186	-0.0006377	0.9292754	-0.3210916	-0.0021220	-0.0242642

In the first step x_4 is added as the regression variable that gives the smallest *AIC*. Then x_3 is the only variable that together with x_4 will reduce the *AIC* criterion, so x_3 is added. In the next step both x_1 and x_2 are candidates to enter the model, but x_1 is the one that reduces the *AIC* criterion the most. x_2 is entered in the next step and after that x_5 . The procedure stops with x_1, x_2, x_3, x_4 and x_5 in the model.

Performing backward elimination, we get the result below.

```
> step(mod)
Start:  AIC=-105.37
y ~ x1 + x2 + x3 + x4 + x5 + x6 + x7
```

	Df	Sum of Sq	RSS	AIC
- x6	1	0.00004	0.24416	-107.369
- x4	1	0.00069	0.24481	-107.299
- x7	1	0.01018	0.25430	-106.311
<none>			0.24412	-105.373
- x2	1	0.02940	0.27352	-104.416
- x5	1	0.02993	0.27405	-104.366
- x1	1	0.38998	0.63410	-82.555
- x3	1	0.61421	0.85833	-74.682

```
Step:  AIC=-107.37
y ~ x1 + x2 + x3 + x4 + x5 + x7
```

	Df	Sum of Sq	RSS	AIC
- x4	1	0.00068	0.24484	-109.297
- x7	1	0.01035	0.25451	-108.289
<none>			0.24416	-107.369
- x5	1	0.03233	0.27649	-106.136
- x2	1	0.04480	0.28896	-104.988
- x1	1	0.41435	0.65851	-83.572
- x3	1	0.78768	1.03184	-71.896

```
Step:  AIC=-109.3
y ~ x1 + x2 + x3 + x5 + x7
```

	Df	Sum of Sq	RSS	AIC
- x7	1	0.01674	0.26158	-109.577
<none>			0.24484	-109.297
- x5	1	0.03582	0.28066	-107.746

```
- x2      1    0.04690 0.29174 -106.740
- x1      1    0.67174 0.91657  -76.975
- x3      1    1.94498 2.18982  -54.331
```

```
Step:  AIC=-109.58
y ~ x1 + x2 + x3 + x5
```

	Df	Sum of Sq	RSS	AIC
<none>			0.26158	-109.577
- x5	1	0.02188	0.28346	-109.488
- x2	1	0.07908	0.34066	-104.709
- x1	1	1.23395	1.49553	-66.246
- x3	1	1.93493	2.19651	-56.252

```
Call:
lm(formula = y ~ x1 + x2 + x3 + x5, data = sourprec)
```

```
Coefficients:
(Intercept)      x1      x2      x3      x5
   5.77017   -0.35181  -0.00233   0.98990  -0.02400
```

In the first step, taking out x_6 reduces the AIC the most. In the next steps first x_4 and thereafter x_7 are taken out and we are left with x_1, x_2, x_3 and x_5 in the model.

Finally, if we perform stepwise regression we get:

```
> step(lm(y~1), y~1+x1+x2+x3+x4+x5+x6+x7, direction="both")
Start:  AIC=-50.24
```

```
y ~ 1
      Df Sum of Sq    RSS    AIC
+ x4   1   2.30355  1.1827 -76.348
+ x2   1   1.13425  2.3520 -58.474
+ x3   1   0.96592  2.5203 -56.676
+ x7   1   0.84759  2.6387 -55.484
+ x5   1   0.57318  2.9131 -52.911
+ x1   1   0.39195  3.0943 -51.342
+ x6   1   0.26643  3.2198 -50.308
<none>      3.4862 -50.241
```

```
Step:  AIC=-76.35
y ~ x4
```

	Df	Sum of Sq	RSS	AIC
+ x3	1	0.39645	0.7862	-84.963
<none>			1.1827	-76.348
+ x6	1	0.05682	1.1259	-75.628

```

+ x2      1    0.03898 1.1437 -75.219
+ x7      1    0.03627 1.1464 -75.158
+ x1      1    0.03260 1.1501 -75.075
+ x5      1    0.02973 1.1530 -75.010
- x4      1    2.30355 3.4862 -50.241

```

Step: AIC=-84.96

y ~ x4 + x3

	Df	Sum of Sq	RSS	AIC
+ x1	1	0.46780	0.31844	-106.463
+ x2	1	0.07316	0.71309	-85.502
<none>			0.78625	-84.963
+ x7	1	0.04413	0.74212	-84.465
+ x6	1	0.00050	0.78575	-82.980
+ x5	1	0.00035	0.78590	-82.975
- x3	1	0.39645	1.18269	-76.348
- x4	1	1.73408	2.52033	-56.676

Step: AIC=-106.46

y ~ x4 + x3 + x1

	Df	Sum of Sq	RSS	AIC
+ x2	1	0.04157	0.27687	-108.100
- x4	1	0.02365	0.34210	-106.600
<none>			0.31844	-106.463
+ x7	1	0.01009	0.30835	-105.300
+ x6	1	0.00705	0.31140	-105.044
+ x5	1	0.00370	0.31475	-104.766
- x1	1	0.46780	0.78625	-84.963
- x3	1	0.83165	1.15010	-75.075

Step: AIC=-108.1

y ~ x4 + x3 + x1 + x2

	Df	Sum of Sq	RSS	AIC
- x4	1	0.00659	0.28346	-109.488
+ x5	1	0.02236	0.25451	-108.289
<none>			0.27687	-108.100
- x2	1	0.04157	0.31844	-106.463
+ x6	1	0.00255	0.27432	-106.340
+ x7	1	0.00038	0.27649	-106.136
- x1	1	0.43622	0.71309	-85.502
- x3	1	0.82188	1.09875	-74.262

Step: AIC=-109.49

y ~ x3 + x1 + x2

	Df	Sum of Sq	RSS	AIC
+ x5	1	0.02188	0.26158	-109.577
<none>			0.28346	-109.488
+ x4	1	0.00659	0.27687	-108.100
+ x6	1	0.00391	0.27955	-107.849
+ x7	1	0.00280	0.28066	-107.746
- x2	1	0.05864	0.34210	-106.600
- x1	1	1.21835	1.50181	-68.137
- x3	1	2.02498	2.30844	-56.960

Step: AIC=-109.58
y ~ x3 + x1 + x2 + x5

	Df	Sum of Sq	RSS	AIC
<none>			0.26158	-109.577
- x5	1	0.02188	0.28346	-109.488
+ x7	1	0.01674	0.24484	-109.297
+ x4	1	0.00707	0.25451	-108.289
+ x6	1	0.00079	0.26079	-107.656
- x2	1	0.07908	0.34066	-104.709
- x1	1	1.23395	1.49553	-66.246
- x3	1	1.93493	2.19651	-56.252

Call:
lm(formula = y ~ x3 + x1 + x2 + x5)

Coefficients:
(Intercept) x3 x1 x2 x5
5.77017 0.98990 -0.35181 -0.00233 -0.02400

As expected, x_4 is entered in the model in the first step, but with x_1, x_2 and x_3 also in the model x_4 is taken out and we end up with the same model as we did when performing backward elimination.

As should be noted there is no way of removing variable once entered when doing forward selection. Also, in common with stepwise regression, it may in the beginning be difficult for variable to be judged significant when using the F-test since the variation in the response caused by the rest of the regression variables is included in the error variance. Therefore, one often runs the procedure several times with different choice of α . If possible, it may be an advantage to use backward elimination. The advantage of the forward selection and stepwise regression algorithms is that they are more general applicable. To my knowledge there are no algorithms implemented for using F-tests or (t-tests) in R. However, backward elimination can be performed the following way where we first start with the full model.

```
> mod=lm(y~x1+x2+x3+x4+x5+x6+x7,data=sourprec)
> summary(mod)
```

```
lm(formula = y ~ x1 + x2 + x3 + x4 + x5 + x6 + x7, data = sourprec)
Residuals:
      Min       1Q   Median       3Q      Max
-0.15266 -0.09355  0.01429  0.06801  0.17755
Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  5.6764334   0.1389162  40.862 < 2e-16 ***
x1          -0.3150444   0.0587512  -5.362 4.27e-05 ***
x2          -0.0018533   0.0012587  -1.472  0.158
x3           0.9751745   0.1449075   6.730 2.62e-06 ***
x4          -0.0002268   0.0010038  -0.226  0.824
x5          -0.0334242   0.0225009  -1.485  0.155
x6          -0.0039399   0.0724339  -0.054  0.957
x7           0.0888722   0.1025724   0.866  0.398
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.1165 on 18 degrees of freedom
Multiple R-squared:  0.93,    Adjusted R-squared: 0.9027
F-statistic: 34.15 on 7 and 18 DF,  p-value: 3.904e-09
```

The least significant one is x_6 . Therefore, in the next step we can use the command update:

```
> mod1=(update(mod,~.-x6))
Call:
lm(formula = y ~ x1 + x2 + x3 + x4 + x5 + x7, data = sourprec)

Residuals:
      Min       1Q   Median       3Q      Max
-0.15269 -0.09218  0.01337  0.06853  0.17636

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  5.6757654   0.1346927  42.139 < 2e-16 ***
x1          -0.3142383   0.0553395  -5.678 1.79e-05 ***
x2          -0.0018918   0.0010132  -1.867  0.0774 .
x3           0.9714256   0.1240784   7.829 2.31e-07 ***
x4          -0.0002239   0.0009757  -0.229  0.8210
x5          -0.0330468   0.0208352  -1.586  0.1292
x7           0.0877409   0.0977709   0.897  0.3807
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.1134 on 19 degrees of freedom
Multiple R-squared:  0.93,    Adjusted R-squared: 0.9078
F-statistic: 42.05 on 6 and 19 DF,  p-value: 5.683e-10
```

By continuing the updating, we end up with the model:

```

Model 1
Call:
lm(formula = y ~ x1 + x2 + x3, data = sourprec)

Residuals:
    Min       1Q   Median       3Q      Max
-0.19867 -0.09314  0.03027  0.07030  0.18855

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  5.6990584  0.0697816  81.670 < 2e-16 ***
x1          -0.3489065  0.0358803  -9.724 2.00e-09 ***
x2          -0.0018364  0.0008608  -2.133  0.0443 *
x3           0.9548356  0.0761644  12.537 1.71e-11 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.1135 on 22 degrees of freedom
Multiple R-squared:  0.9187,    Adjusted R-squared:  0.9076
F-statistic: 82.86 on 3 and 22 DF, p-value: 3.821e-12

```

From backward elimination and stepwise regression using the AIC-criterion we get:

```

Model 2
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  5.7701749  0.0871028  66.246 < 2e-16 ***
x1          -0.3518082  0.0353467  -9.953 2.10e-09 ***
x2          -0.0023299  0.0009247  -2.520  0.0199 *
x3           0.9899007  0.0794240  12.463 3.61e-11 ***
x5          -0.0239962  0.0181057  -1.325  0.1993
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.1116 on 21 degrees of freedom
Multiple R-squared:  0.925,    Adjusted R-squared:  0.9107
F-statistic: 64.72 on 4 and 21 DF, p-value: 1.659e-11

```

And from the forward selection method using AIC

```

Model 3
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  5.7668186  0.0881546  65.417 < 2e-16 ***
x1          -0.3210916  0.0545427  -5.887 9.29e-06 ***
x2          -0.0021220  0.0009753  -2.176  0.0417 *
x3           0.9292754  0.1142848   8.131 9.06e-08 ***
x4          -0.0006377  0.0008557  -0.745  0.4647

```

```

x5          -0.0242642  0.0183040  -1.326    0.1999
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.1128 on 20 degrees of freedom
Multiple R-squared:  0.927,    Adjusted R-squared:  0.9087 
F-statistic: 50.79 on 5 and 20 DF,  p-value: 1.100e-10

```

Ranked according to $R^2_{adjusted}$ model 2 is better than model 3 which again is better than model 1. But model 3 has two non-significant terms and model 2 has one. The two candidate models would in practice be model 1 and 2.

Now regardless of if we include the amount of organic material or not, we have arrived at a model saying that adding SO_4 and NO_3 makes the lakes more sour and adding Ca makes it less sour which should be in agreement with common knowledge. But what about the content of latent aluminum that is so highly correlated with the ph value. The explanation is simply that in sour lakes aluminum will be released from the rock. Therefore, sour lakes cause a high amount of aluminum, not the other way around.

Diagnostic checks of regression models.

The model we arrive at when performing regression analysis depends, strongly on the quality of the data we have i. e. the observations for the response and the regression variables. To thrust that the model is useful, we should check the quality of these. We will now discuss some problems that may arise and some methods that can be used to detect if such problems are present in our data.

Multicollinearity

Multicollinearity arises when two or more columns in the X -matrix are strongly correlated (almost linear dependent). Multicollinearity can sometimes be discovered by estimating the correlation between regression variables. The correlation between regression variable x_i and x_j is given by:

$$r_{ij} = \frac{\sum_{k=1}^n (x_{ik} - \bar{x}_i)(x_{jk} - \bar{x}_j)}{\sqrt{\sum_{k=1}^n (x_{ik} - \bar{x}_i)^2 \sum_{k=1}^n (x_{jk} - \bar{x}_j)^2}}.$$

Multicollinearity often occurs when the variation interval of two or more variables are approximately equal.

An example is if $x_{ki} \in (0.95, 1.05)$, $i=1,2,\dots,n$, the column with 1's, x_k and x_k^2 are strongly correlated. A good measure of multicollinearity is given by the variance inflation factor, VIF.

$$\text{VIF} = \frac{1}{1 - R_j^2}$$

where R_j^2 is the multiple determination coefficient when regression variable x_j is regressed on the other regression variables. If you ever meet the word tolerance, it is given by $1 - R_j^2$. VIF should be less than 10. If multicollinearity occurs one should consider:

1. Remove columns.
2. Collect more data.
3. Respecify the model. For instance if multicollinearity occurs between the three regression variables x_1, x_2 and x_3 , a new variable $\frac{x_1 + x_2 + x_3}{3}$ or $x_1 x_2 x_3$ may solve the problem.
4. Center variables i.e. use $x_i - \bar{x}$ instead of x_i if one works with polynomial models.

There are also several statistical methods like Principal component regression (PCR), Partial least square (PLS) and Ridge regression for handling such problems. These are, however, beyond the scope of this course.

Observations with influence

Influential points are points that when removed from the dataset could cause a large change in the fit. The performance of variable selection methods may depend heavily on influential observation. Therefore, it is a good idea to check for leverage points and the Cook's distance before starting a variable selection procedure.

Leverage points

We have $\hat{\mathbf{y}} = \mathbf{H}\mathbf{y} \Rightarrow \hat{y}_i = \sum_{j=1}^n h_{ij} y_j = h_{ii} y_i + \sum_{j \neq i} h_{ij} y_j$. If h_{ii} is large y_i could have a large influence on \hat{y}_i .

Now $\text{rank}(\mathbf{H}) = k+1 = \text{tr}(\mathbf{H}) = \sum_{i=1}^n h_{ii}$.

Hence the average value of $h_{ii} = \frac{k+1}{n}$. Observations where $h_{ii} > \frac{2(k+1)}{n}$ are said to be

leverage points. Note that if \mathbf{x}_i' is the i -th row in the \mathbf{X} -matrix, then $h_{ii} = \mathbf{x}_i' (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_i$.

Hence for an observation with a high leverage, the uncertainty in the estimator for the expected value of the response will be relatively high, and a corresponding prediction interval will be relatively wide.

Leverage points are entirely determined by the values in the X -matrix and need not necessarily influence our estimated coefficients. A check for this is done by the Cook's distance, D_i .

First: $(\hat{y} - \hat{y}_{(-i)})^t (\hat{y} - \hat{y}_{(-i)}) = (Xb - Xb_{(-i)})^t (Xb - Xb_{(-i)}) = (b - b_{(-i)})^t X^t X (b - b_{(-i)})$. The $(-i)$ notation means that the i -th set of observations $(y_i, x_{i1}, \dots, x_{ik})$ is taken out.

The Cook's distances D_i is defined as $D_i = \frac{(\hat{\beta} - \hat{\beta}_{(-i)})^t X^t X (\hat{\beta} - \hat{\beta}_{(-i)})}{p \hat{\sigma}^2}$, $i = 1, 2, \dots, n$.

Usually, we consider an observation $(y_i, x_{i1}, \dots, x_{ik})$ to be influential if $D_i > 1$.

Study of residuals

We expect the residuals from a good model to be uncorrelated with mean zero and constant variance. Now, if we have the correct model we can argue as follows:

$$\hat{\epsilon} = Y - \hat{Y} = (I - H)Y = (I - H)\epsilon.$$

Thereby $E(\hat{\epsilon}) = (I - H)E(\epsilon) = 0,$

and $Cov(\hat{\epsilon}) = E[(I - H)\epsilon\epsilon^t(I - H)^t] = \sigma^2(I - H).$

which shows that $\hat{\epsilon}_1, \hat{\epsilon}_2, \dots, \hat{\epsilon}_n$ are not uncorrelated. Further we get

$$Var(\hat{\epsilon}_i) = \sigma^2(1 - h_{ii}).$$

Hence the variances of $\hat{\epsilon}_1, \hat{\epsilon}_2, \dots, \hat{\epsilon}_n$ need not be equal either.

Since $Var\left(\frac{\hat{\epsilon}_i}{\sigma\sqrt{1-h_{ii}}}\right) = 1$ we often use $r_i = \frac{e_i}{s\sqrt{1-h_{ii}}}$ in the study of the residuals . These residuals are often said to be the standardized residuals. Some books call them studentized.

The residuals $\delta_i = Y_i - \hat{Y}_{i,-i}$, $i = 1, 2, \dots, n$ where $\hat{Y}_{i,-i}$ is the estimator for the i -th fitted value when the observation $(y_i, x_{i1}, \dots, x_{ik})$ is taken out are called the Press residuals . These can be

shown to be equal to $\frac{\hat{\epsilon}_i}{1-h_{ii}}$ and $Var(\delta_i) = \frac{Var(\hat{\epsilon}_i)}{(1-h_{ii})^2} = \frac{\sigma^2(1-h_{ii})}{(1-h_{ii})^2} = \frac{\sigma^2}{1-h_{ii}}.$

Hence $\frac{\delta_i}{Sd(\delta_i)} = \frac{\hat{\epsilon}_i}{(1-h_{ii})\sigma} = \frac{\hat{\epsilon}_i}{\sigma\sqrt{1-h_{ii}}}$

Now assume we estimate σ^2 by $S_{(-i)}^2$ which is the estimator for the variance when $(y_i, x_{i1}, \dots, x_{ik})$ is taken out. Then $\delta_i = Y_i - \hat{Y}_{i,-i}$ and $S_{(-i)}^2$ are independent and the residuals

$$\frac{\hat{\varepsilon}_i}{S_{-i}\sqrt{1-h_{ii}}}$$

is t-distributed with $n-k-2$ degrees of freedom and can be used to investigate if the expected value of a residual is different from zero, i.e. if we have an outlier. The residuals $\frac{e_i}{s_{-i}\sqrt{1-h_{ii}}}$ are called R-studentized, sometimes external studentized.

To summarize. Check of residuals provide information about:

1. Outliers
2. Heterogeneity in variance
3. Misspecification of models (or if other variables should be included)
4. Normal distribution

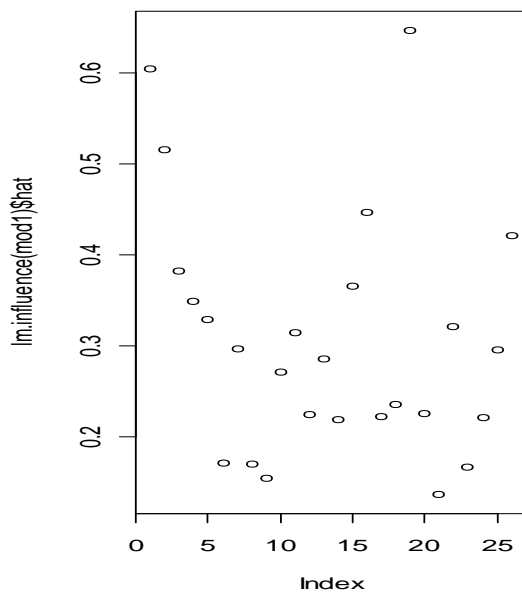
There are extremely many diagnostic tools developed for regression models and one can easily be confused about the importance of the information provided from each of them. In practice it is most common to check for multicollinearity, check the Cook's distance and perform residuals plots. The most common residual plots are plot of residuals against fitted values to check for heterogeneity in variance, against other regression variables to see if they should be included and normal plot to check for normally distributed data and outliers.

Example. Sour precipitation

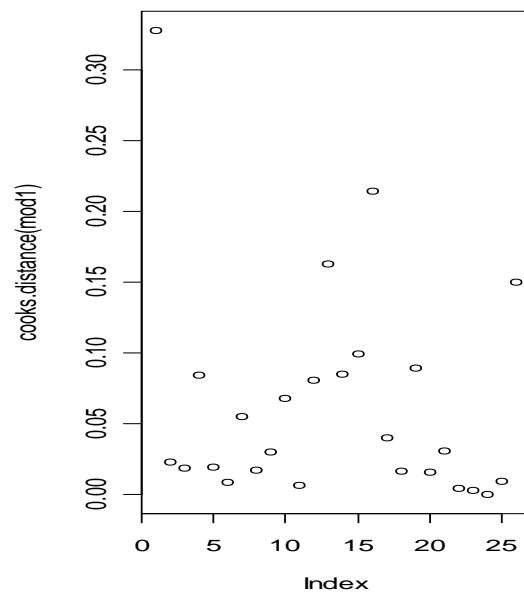
Let us return to our example about sour precipitation. The correlation matrix between pairs of variables indicated some high correlations between x_7 and other regression variables, but no clear sign of multicollinearity. The VIF-factor is in the R-library car and not in the basic R-package. This library car can be downloaded to any personal PC for free.

The critical value for the leverage is for this data set $\frac{2 \cdot 8}{26} = 0.62$. From the plot below there may be one observation with too high leverage. However, from the Cook's distance there

does not seem to be any observation that is really influential on our estimated coefficients.

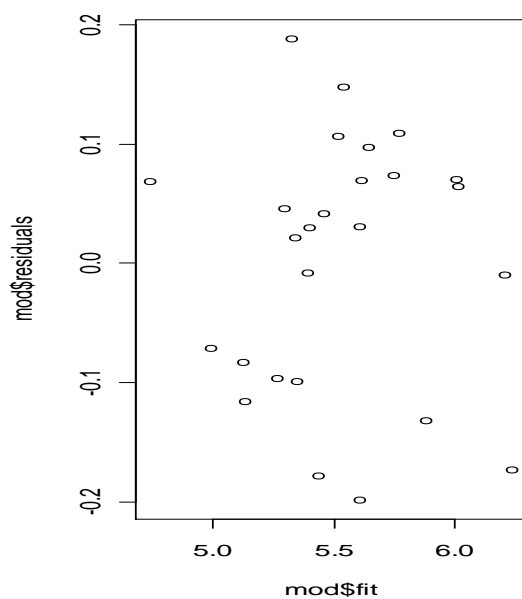


Leverages

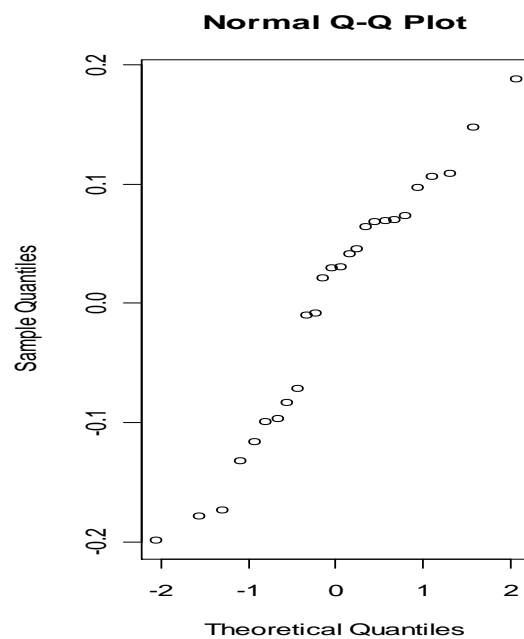


Cook's distances

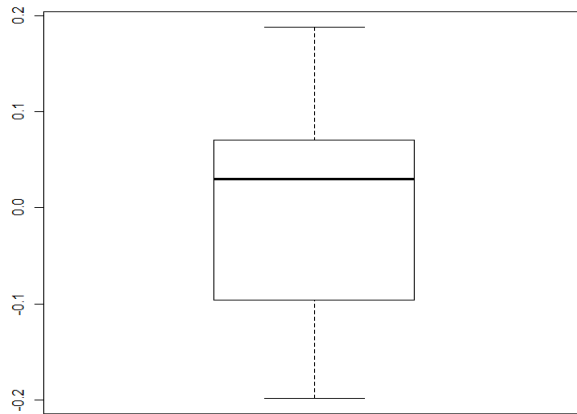
Next is shown a plot of residuals against fitted values and a normal plot. The residuals are obtained using x_1 , x_2 and x_3 as regression variables, since these were the only three significant ones. Examining the plots, the variance does not seem to depend on the value of the fitted model and the normal plot is close to a straight line. The box plot reveal that the median of the residuals is above zero, but part from that there is no observation that is close to being an outlier.



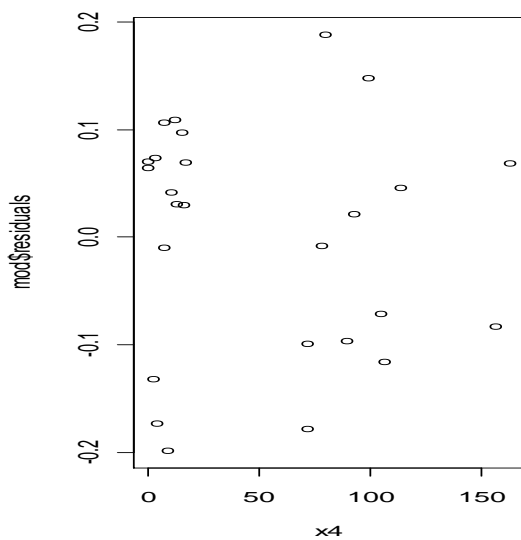
Plot of residuals against fitted values



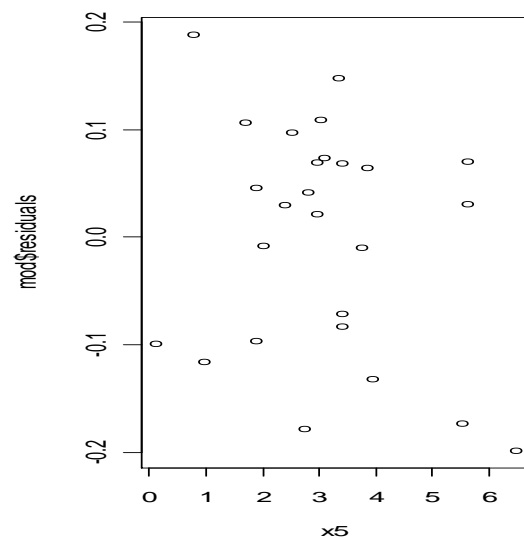
Normal plot of residuals



Boxplot of the residuals



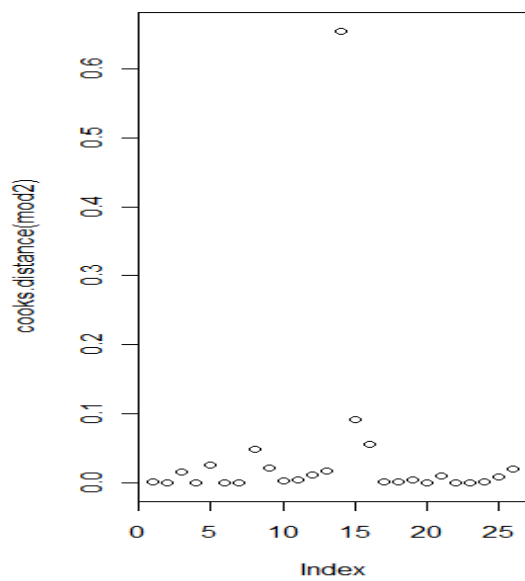
Plot of residuals against x_4



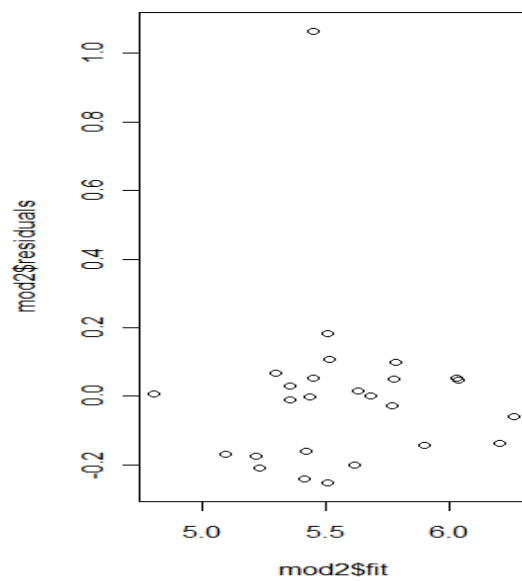
Plot of residuals against x_5

Finally plots of the residuals against x_4 and x_5 (content of aluminum and organic material) do not show any particular pattern and there is no obvious reason to include these in the model given that using x_1 , x_2 and x_3 already are there. Hence these diagnostic tools give no reason to not trust our model.

To get an impression of what might happen if we had observations of poor quality, we have added 1 to one of the response values, while all the rest is kept unchanged. We estimate the same model.



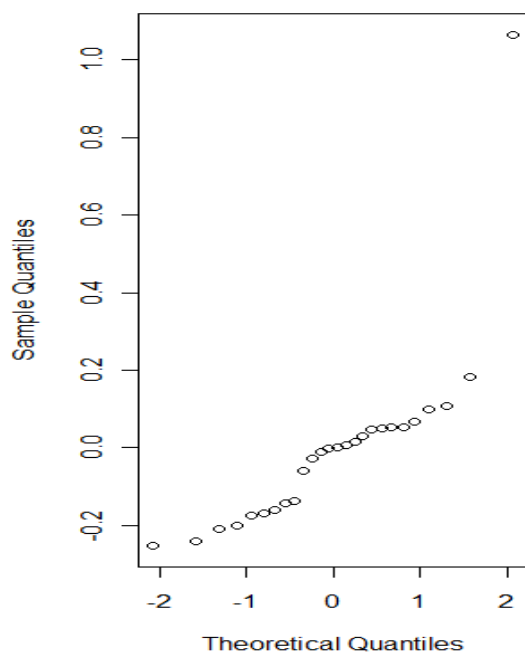
Cook's distances



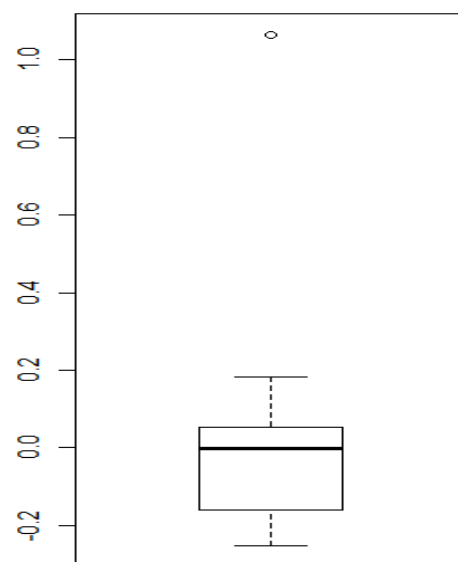
Plot of residuals against fitted values

We observe that the Cook's distance does not classify this observation as being strongly influential, while the normal plot finds one residual that does not fit into the assumption of having mean equal to zero and residuals being normally distributed. The boxplot of the residuals clearly signals an outlier. All the plots signal a very special observation.

Normal Q-Q Plot



Normal plot of residuals



Box plot of residuals

If we look at the estimated models we get:

Original data set:

```
lm(formula = y ~ x1 + x2 + x3, data = sour)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	5.6990584	0.0697816	81.670	< 2e-16	***
x1	-0.3489065	0.0358803	-9.724	2.00e-09	***
x2	-0.0018364	0.0008608	-2.133	0.0443	*
x3	0.9548356	0.0761644	12.537	1.71e-11	***

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

Residual standard error: 0.1135 on 22 degrees of freedom

Multiple R-squared: 0.9187, Adjusted R-squared: 0.9076

F-statistic: 82.86 on 3 and 22 DF, p-value: 3.821e-12

manipulated data set

```
lm(formula = y ~ x1 + x2 + x3, data = sour1)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	5.7446812	0.1624766	35.357	< 2e-16	***
x1	-0.3882316	0.0835422	-4.647	0.000124	***
x2	0.0002539	0.0020043	0.127	0.900338	
x3	0.9718510	0.1773381	5.480	1.66e-05	***

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

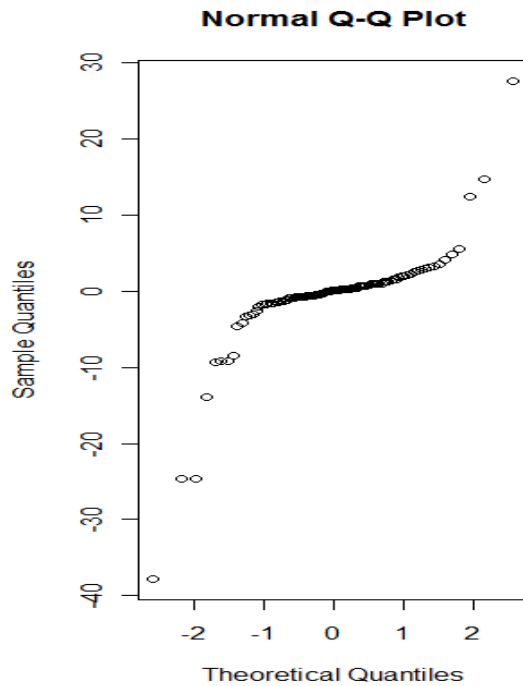
Residual standard error: 0.2643 on 22 degrees of freedom

Multiple R-squared: 0.6522, Adjusted R-squared: 0.6048

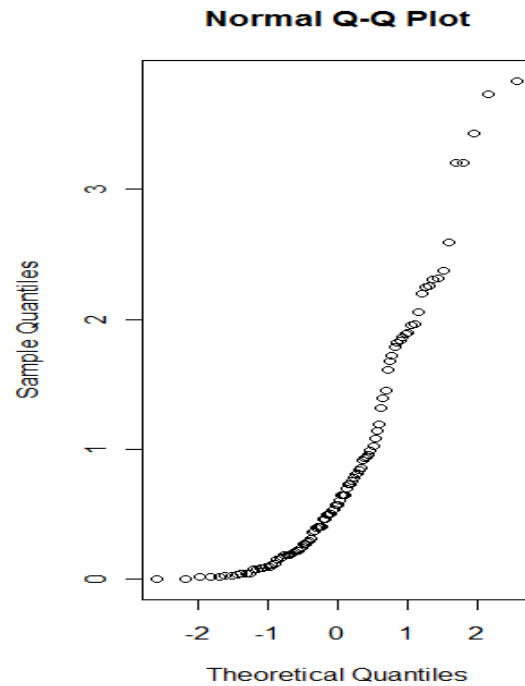
F-statistic: 13.75 on 3 and 22 DF, p-value: 2.874e-05

We observe that the estimated coefficients have not changed much and that is why the observation is not identified as an influential observation. The estimated standard deviations of the coefficients, however, have changed dramatically and x_2 is no longer significant. This illustrates that residual plots may provide information in addition to the leverages and the Cook's distances. A correct procedure would now be to try find out why observation number 14 is an outlier. Is it just a wrong recording? If a good reason is found the data are reanalyzed with the corresponding observation taken out or corrected.

To illustrate some other situations that may occur we show normal plots performed on data from a distribution with heavy tails and data from a skew distribution. If we have data from a distribution with lighter tails than the normal distribution and with expected value equal to 0, the form of the normal-plot would be an approximate mirror-image from the one from a heavy tailed distribution, mirrored around the line $y=x$.



Data from a distribution with heavy tails.



Data from a skew distribution

Other methods for choosing the best model

Calculating the PRESS-residuals (Predictive sum of squares) $\delta_i^* = y_i - y_{i,-i}$ is an example of doing leave one out cross validation. These residuals measure how good the i -th observation of the response can be estimated or predicted in terms of the other observations.

The model with the “best prediction ability” should be the one that minimizes $\sum_{i=1}^n |\delta_i^*|$ or

$$\sum_{i=1}^n (\delta_i^*)^2 = \text{PRESS}, \text{ eventually the one that maximizes } R_{pred}^2 = 1 - \frac{\sum_{i=1}^n (\delta_i^*)^2}{SS_T}$$

s^2 is an estimate for $\text{Var}(Y_i)$ in the model with p -parameters and $\hat{\sigma}^2$ is an estimate for σ^2 . Often it is the estimate of σ^2 from the full model.

Categorical or Indicator variables

Suppose

$$Y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \beta_3 Z_i + \varepsilon_i$$

Where the variable Z is a categorical variable with l categories. Such regression problems are normally best solved by introducing $l-1$ indicator variables. For instance if we have three different catalyst the regression model can be:

$$Y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \beta_3 Z_{1i} + \beta_4 Z_{2i} + \varepsilon_i \quad i = 1, 2, \dots, n$$

	Z_1	Z_2
	1	0
	1	0
	\vdots	\vdots
	1	0
	0	1
	0	1
	\vdots	\vdots
	0	1
	0	0
	0	0
	\vdots	\vdots
	0	0

where the columns Z_1 and Z_2 are

The $(1,0)$ combination corresponds to catalyst 1, the $(0,1)$ to catalyst 2 and the $(0,0)$ combination to catalyst 3.

The model is then:

$$Y_i = (\beta_0 + \beta_3) + \beta_1 x_{1i} + \beta_2 x_{2i} + \varepsilon_i \quad \text{for catalyst 1}$$

$$Y_i = (\beta_0 + \beta_4) + \beta_1 x_{1i} + \beta_2 x_{2i} + \varepsilon_i \quad \text{for catalyst 2}$$

and

$$Y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \varepsilon_i \quad \text{for catalyst 3 .}$$

If we believe the coefficients for the continuous variables depend on the level of the category variables, we introduce product terms. Suppose we construct the model

$$Y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \beta_3 Z_{1i} + \beta_4 Z_{2i} + \beta_5 x_{1i} z_{1i} + \beta_6 x_{2i} z_{2i} + \varepsilon_i \quad i = 1, 2, \dots, n$$

Then we have

$$Y_i = (\beta_0 + \beta_3) + (\beta_1 + \beta_5) x_{1i} + \beta_2 x_{2i} + \varepsilon_i \quad \text{for catalyst 1}$$

$$Y_i = (\beta_0 + \beta_4) + \beta_1 x_{1i} + (\beta_2 + \beta_6) x_{2i} + \varepsilon_i \quad \text{for catalyst 2}$$

$$Y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \varepsilon_i \quad \text{for catalyst 3}$$

Orthogonal columns in the Design matrix

Let the design matrix $X = [1, \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k]$.

If $\mathbf{x}_p^t \mathbf{x}_q = 0$ i.e. $\sum_{i=1}^n x_{pi} x_{qi} = 0$, $p \neq q$ we say that the columns \mathbf{x}_p and \mathbf{x}_q are orthogonal. If all the columns are orthogonal we have:

$$(\mathbf{X}^t \mathbf{X}) = \begin{bmatrix} n & 0 & \cdots & 0 \\ 0 & \mathbf{x}_1^t \mathbf{x}_1 & \cdots & 0 \\ 0 & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \mathbf{x}_k^t \mathbf{x}_k \end{bmatrix} \text{ and}$$

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \mathbf{Y} = \begin{bmatrix} 1/n & 0 & \cdots & 0 \\ 0 & (\mathbf{x}_1^t \mathbf{x}_1)^{-1} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & (\mathbf{x}_k^t \mathbf{x}_k)^{-1} \end{bmatrix} \begin{bmatrix} \sum_{i=1}^n Y_i \\ \mathbf{x}_1^t \mathbf{Y} \\ \vdots \\ \mathbf{x}_k^t \mathbf{Y} \end{bmatrix}.$$

We observe that the estimator for β_j only depends upon \mathbf{x}_j and \mathbf{Y} . In addition, we get

$$SS_R = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = \sum_{i=1}^n (b_0 + b_1 x_{1i} + \cdots + b_k x_{ki} - b_0)^2 = b_1^2 \sum_{i=1}^n x_{1i}^2 + b_2^2 \sum_{i=1}^n x_{2i}^2 + \cdots + b_k^2 \sum_{i=1}^n x_{ki}^2 \text{ or}$$

$$SS_R(\beta_0, \beta_1, \dots, \beta_k) = SS_R(\beta_1 | \beta_0) + SS_R(\beta_2 | \beta_0) + \cdots + SS_R(\beta_k | \beta_0) \text{ where}$$

$$SS_R(\beta_j | \beta_0) = b_j^2 \sum_{i=1}^n x_{ji}^2, j=1, 2, \dots, n.$$

Regression analysis can easily fill a whole course or more. In cases where you find that more knowledge is needed than what is covered here you may benefit from checking out one or more of these books:

Draper, N. R. & Smith, H. Applied Regression Analysis. Third edition.

Montgomery, D. C. Peck, E. A. & Vining, G. G. Introduction to Linear Regression Analysis, Third Edition.

Abraham, B & Ledolter, J. Introduction to Regression Modeling

If you want to to increase your theoretical skills these books may be good choices:

Rencher, A. C. & Schaalje, G. B. Linear Models in Statistics.

Seber, G. A. F. Linear Regression Analysis