

Statistical inference part 2

Practice session 2: Simple and multiple linear regression

Exercises on simple linear regression

6.2 (a) Show that $E(\hat{\beta}_1) = \beta_1$ as in (6.7).

(b) Show that $E(\hat{\beta}_0) = \beta_0$ as in (6.8).

$$\text{with } \hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad \text{and} \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\text{and } E[y_i] = \beta_0 + \beta_1 x_i$$

6.9 (a) Obtain a test for $H_0: \beta_0 = a$ versus $H_1: \beta_0 \neq a$.

(b) Obtain a confidence interval for β_0 .

$$\text{Use } \hat{\beta}_0 \sim \mathcal{N}\left(\beta_0, \sigma^2 \left[1/n + \bar{x}^2 / \sum_{i=1}^n (x_i - \bar{x})^2\right]\right), \quad (n-2) \frac{s^2}{\sigma^2} \sim \chi_{(n-2)}^2$$

and $\hat{\beta}_0$ and s^2 independent. Thus deduce the distribution of $\frac{\hat{\beta}_0 - \beta_0}{\sigma \sqrt{c}}$...

6.14 Table 6.1 (Weisberg 1985, p. 231) gives the data on daytime eruptions of Old Faithful Geyser in Yellowstone National Park during August 1–4, 1978. The variables are x = duration of an eruption and y = interval to the next eruption. Can x be used to successfully predict y using a simple linear model $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$? see file geyser.csv

(a) Find $\hat{\beta}_0$ and $\hat{\beta}_1$.

(b) Test $H_0: \beta_1 = 0$ using (6.14).

(c) Find a confidence interval for β_1 .

(d) Find r^2 using (6.16).

6.2

(a)

 $\hat{\beta}_1$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^m (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^m (x_i - \bar{x})^2}, \quad \text{we have } y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^m (x_i - \bar{x})(\beta_1(x_i - \bar{x}) + \epsilon_i - \bar{\epsilon})}{\sum_{i=1}^m (x_i - \bar{x})^2}$$

$$= \beta_1 + \frac{\sum_{i=1}^m (x_i - \bar{x})(\epsilon_i - \bar{\epsilon})}{\sum_{i=1}^m (x_i - \bar{x})^2}$$

$$E[\hat{\beta}_1] = \beta_1 + \frac{\sum_{i=1}^m (x_i - \bar{x})(\overbrace{E(\epsilon_i)}^0 - \overbrace{\bar{E}(\bar{\epsilon})}^0)}{\sum_{i=1}^m (x_i - \bar{x})^2} \quad \text{by linearity of the expectation}$$

(b)

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$E[\hat{\beta}_0] = E[\bar{y}] - E[\hat{\beta}_1] \cdot \bar{x} = \beta_0 + \beta_1 \bar{x} - \beta_1 \bar{x} = \beta_0$$

thus $\hat{\beta}_1$ and $\hat{\beta}_0$ are unbiased

6.3

(a)

$$\frac{\hat{\beta}_0 - \beta_0}{s\sqrt{c}} \sim N(0,1) \quad \text{and} \quad (n-2) \frac{s^2}{\sigma^2} \sim \chi^2_{(n-2)}$$

$$\text{and } \hat{\beta}_0 \perp s^2 \quad \text{thus} \quad \frac{\hat{\beta}_0 - \beta_0}{s\sqrt{c}} \sim t_{(n-2)} \quad \text{with } c = \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\text{Under } H_0: \beta_0 = a, \quad t := \frac{\hat{\beta}_0 - a}{s\sqrt{c}} \sim t_{(n-2)}$$

Thus we reject H_0 at risk α if $|t| > t_{\alpha/2, n-2}$

(b)

$$\frac{\hat{\beta}_0 - \beta_0}{s\sqrt{c}} \sim t_{(n-2)} \quad \text{thus}$$

$$P\left(-t_{\alpha/2, n-2} \leq \frac{\hat{\beta}_0 - \beta_0}{s\sqrt{c}} \leq t_{\alpha/2, n-2}\right) = 1 - \alpha$$

$$\Leftrightarrow P\left(\hat{\beta}_0 - t_{\alpha/2, n-2} s\sqrt{c} \leq \beta_0 \leq \hat{\beta}_0 + t_{\alpha/2, n-2} s\sqrt{c}\right) = 1 - \alpha$$

Thus an $1 - \alpha$ confidence interval for β_0 is

$$\boxed{\hat{\beta}_0 \pm t_{\alpha/2, n-2} s \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}}$$

6.14 : See Rmd file.

Exercises on multiple linear regression

Prove (i) of theorem 7.6b

Theorem 7.6b. Suppose that \mathbf{y} is $N_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I})$, where \mathbf{X} is $n \times (k+1)$ of rank $k+1 < n$ and $\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_k)'$. Then the maximum likelihood estimators $\hat{\boldsymbol{\beta}}$ and $\hat{\sigma}^2$ given in Theorem 7.6a have the following distributional properties:

(i) $\hat{\boldsymbol{\beta}}$ is $N_{k+1}[\boldsymbol{\beta}, \sigma^2(\mathbf{X}'\mathbf{X})^{-1}]$.
use that $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y}$ | Thus give $\text{cov}(\hat{\boldsymbol{\beta}})$

7.2 Show that (7.10) follows from (7.9). Why is $\mathbf{X}'\mathbf{X}$ positive definite, as noted below (7.10)?

7.29 (a) Show that R^2 in (7.55) can be written in the form $R^2 = 1 - \text{SSE} / \sum_i (y_i - \bar{y})^2$.

(b) Replace SSE and $\sum_i (y_i - \bar{y})^2$ in part (a) by variance estimators $\text{SSE}/(n-k-1)$ and $\sum_i (y_i - \bar{y})^2/(n-1)$ and show that the result is the same as R_a^2 in (7.56).

7.53 When gasoline is pumped into the tank of a car, vapors are vented into the atmosphere. An experiment was conducted to determine whether y , the amount of vapor, can be predicted using the following four variables based on initial conditions of the tank and the dispensed gasoline:

x_1 = tank temperature ($^{\circ}\text{F}$)

x_2 = gasoline temperature ($^{\circ}\text{F}$)

x_3 = vapor pressure in tank (psi)

x_4 = vapor pressure of gasoline (psi)

data gas.csv

The data are given in Table 7.3 (Weisberg 1985, p. 138).

(a) Find $\hat{\boldsymbol{\beta}}$ and s^2 .

(b) Find an estimate of $\text{cov}(\hat{\boldsymbol{\beta}})$.

(c) Find $\hat{\beta}_1$ and $\hat{\beta}_0$ using \mathbf{S}_{xx} and \mathbf{s}_{xy} as in (7.46) and (7.47).

(d) Find R^2 and R_a^2 .

Proof of Thm 7.6

$$y \sim N_n(X\beta, \sigma^2 I)$$

$$\hat{\beta} = (X'X)^{-1} X'y = Ay \text{ with } A = (X'X)^{-1} X'$$

Thus $\hat{\beta} \sim N_{k+1}(\underbrace{A X \beta}_{(X'X)^{-1} X' X \beta}, \underbrace{A \sigma^2 I A'}_{\sigma^2 (X'X)^{-1}})$ due to the properties of linear combinations of Gaussian vectors

$$\Rightarrow \hat{\beta} \sim N_{k+1}(\beta, \sigma^2 (X'X)^{-1})$$

$\hat{\beta}$ follows a multivariate normal distribution centered in β and with variance-covariance matrix $\sigma^2 (X'X)^{-1}$

unfortunately σ^2 is unknown ..., but it can be estimated by

$$s^2 = \frac{1}{n-k-1} (y - X\hat{\beta})'(y - X\hat{\beta})$$

And $(n-k-1) \frac{s^2}{\sigma^2} \sim \chi^2_{(n-k-1)}$ and $\hat{\beta}$ and s^2

independent. This will be useful to build tests and confidence intervals on β .

$$(7.2) \quad (y - Xb)'(y - Xb) = (y - X\hat{\beta})' + (X\hat{\beta} - Xb)'(y - X\hat{\beta}) + (X\hat{\beta} - Xb)'(y - X\hat{\beta}) + (X\hat{\beta} - Xb)'(X\hat{\beta} - Xb) \quad (7.9)$$

$$= (y - X\hat{\beta})'(y - X\hat{\beta}) + (\hat{\beta} - b)'X'X(\hat{\beta} - b) + 2(\hat{\beta} - b)'(X'y - X'X\hat{\beta}). \quad (7.10)$$

from (7.9) we get

$$(y - X\hat{\beta})'(y - X\hat{\beta}) + (y - X\hat{\beta})'(\underbrace{X\hat{\beta} - Xb}_{X(\hat{\beta} - b)}) + (X\hat{\beta} - Xb)'(y - X\hat{\beta}) + (X\hat{\beta} - Xb)'(X\hat{\beta} - Xb)$$

$$= (y - X\hat{\beta})'(y - X\hat{\beta}) + 2(\hat{\beta} - b)'(X'y - X'X\hat{\beta}) + (\hat{\beta} - b)'X'X(\hat{\beta} - b)$$

Let $u \in \mathbb{R}^{k+1}$ $u'X'Xu = v'v$ with $v = Xu$
 $= \|v\|^2$

if $u \neq 0$ then $v \neq 0$ since X is of rank $k+1$

(if there exist $u \neq 0$ such $Xu = 0$ this would mean that a column of X would be a combination of the other which would contradict the assumption that X is of rank $k+1$)

thus if $u \neq 0$ $u'X'Xu = \|v\|^2 > 0$ with $v = Xu$

thus $X'X$ is a positive definite matrix.

(7.29)

$$R^2 = \frac{SSR}{SST}$$

with $SSR = \sum_{i=1}^m (\hat{y}_i - \bar{y})^2$ and $SST = \sum_{i=1}^m (y_i - \bar{y})^2$

$$SSE = \sum_{i=1}^m (y_i - \hat{y}_i)^2.$$

Thus we need to show that $SST = SSR + SSE$ to deduce that

$$R^2 = 1 - \frac{SSE}{SST}$$

see solution p. 111

7.7 R^2 IN FIXED- x REGRESSION

In (7.39), we have $SSE = \sum_{i=1}^n (y_i - \bar{y})^2 - \hat{\beta}'_1 \mathbf{X}'_c \mathbf{y}$. Thus the corrected total sum of squares $SST = \sum_i (y_i - \bar{y})^2$ can be partitioned as

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \hat{\beta}'_1 \mathbf{X}'_c \mathbf{y} + SSE, \quad (7.53)$$

$$SST = SSR + SSE,$$

where $SSR = \hat{\beta}'_1 \mathbf{X}'_c \mathbf{y}$ is the *regression sum of squares*. From (7.37), we obtain $\mathbf{X}'_c \mathbf{y} = \mathbf{X}'_c \mathbf{X}_c \hat{\beta}_1$, and multiplying this by $\hat{\beta}'_1$ gives $\hat{\beta}'_1 \mathbf{X}'_c \mathbf{y} = \hat{\beta}'_1 \mathbf{X}'_c \mathbf{X}_c \hat{\beta}_1$. Then $SSR = \hat{\beta}'_1 \mathbf{X}'_c \mathbf{y}$ can be written as

$$SSR = \hat{\beta}'_1 \mathbf{X}'_c \mathbf{X}_c \hat{\beta}_1 = (\mathbf{X}_c \hat{\beta}_1)' (\mathbf{X}_c \hat{\beta}_1). \quad (7.54)$$

In this form, it is clear that SSR is due to $\beta_1 = (\beta_1, \beta_2, \dots, \beta_k)'$.

The proportion of the total sum of squares due to regression is

$$R^2 = \frac{\hat{\beta}'_1 \mathbf{X}'_c \mathbf{X}_c \hat{\beta}_1}{\sum_{i=1}^n (y_i - \bar{y})^2} = \frac{SSR}{SST}, \quad (7.55)$$

which is known as the *coefficient of determination* or the *squared multiple correlation*. The ratio in (7.55) is a measure of model fit and provides an indication of how well the x 's predict y .

The partitioning in (7.53) can be rewritten as the identity

$$\begin{aligned} \sum_{i=1}^n (y_i - \bar{y})^2 &= \mathbf{y}'\mathbf{y} - n\bar{y}^2 = (\hat{\beta}'\mathbf{X}'\mathbf{y} - n\bar{y}^2) + (\mathbf{y}'\mathbf{y} - \hat{\beta}'\mathbf{X}'\mathbf{y}) \\ &= SSR + SSE, \end{aligned}$$

which leads to an alternative expression for R^2 :

$$R^2 = \frac{\hat{\beta}'\mathbf{X}'\mathbf{y} - n\bar{y}^2}{\mathbf{y}'\mathbf{y} - n\bar{y}^2}. \quad (7.56)$$

7.53 see the Rmd file .