

Chapter 13

The Simple Linear Regression Model: Theory

13.1 The model

13.1.1 The data

observations	response variable	explanatory variable
1	y_1	x_1
2	y_2	x_2
\vdots	\vdots	\vdots
n	y_n	x_n

Plotting the data.

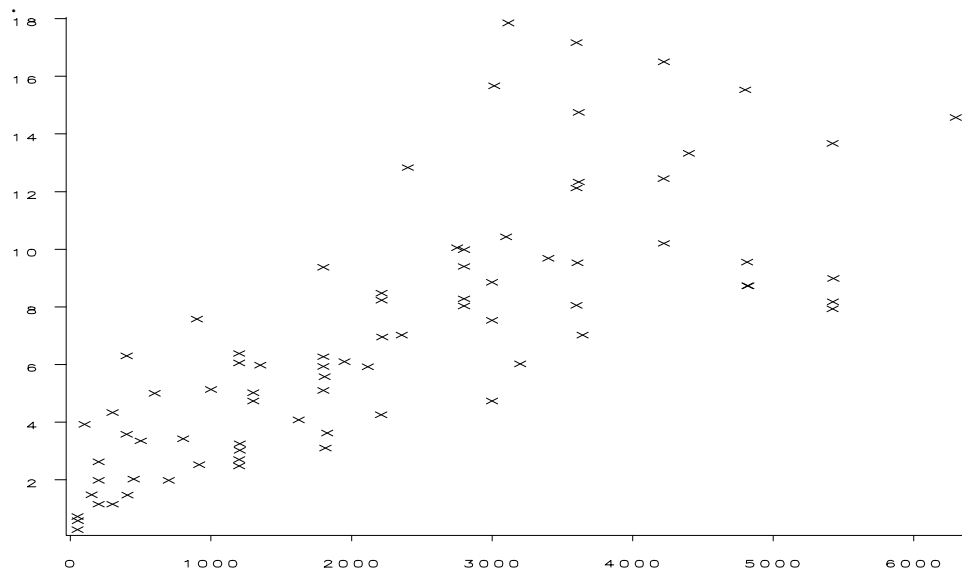


Figure 13.1: Displaying the cable data considered by Cohen at al (1993). There are 79 observations of the number of hours y needed to splice x pairs of wires for a particular type of telephone cable

If the plot is not linear try a simple transformation to linearity. i.e. log, square root, square.

13.1.2 Assumptions for the model

i) The assumption about the linearity of the model

$$Y_i = \alpha + \beta x_i + \varepsilon_i \quad \text{for} \quad i = 1, 2, \dots, n$$

ii) The assumption about the error distribution for ε_i

- a) Full distributional assumption for error term ε_i .

$$\varepsilon_i \sim N(0, \sigma^2) \text{ and } \varepsilon_i \text{ and } \varepsilon_j \text{ for } i \neq j \text{ are independent.}$$

Estimation in this case of the parameters α , β and σ^2 is achieved by Maximum Likelihood.

- b) Assumption about the first and second moments of the distribution for ε_i .

$$E(\varepsilon_i) = 0$$

$$Var(\varepsilon_i) = \sigma^2$$

$$Cov(\varepsilon_i, \varepsilon_j) = 0$$

Estimation in this case can be achieved by Least Squares.

iii) The assumption about the x-variable.

The x-variable is not a random variable and it is fixed at the observed values

13.2 Least squares estimation of parameters

$$\text{Let } S(\alpha, \beta) = \sum_{i=1}^n (y_i - \alpha - \beta x_i)^2$$

where the y_i are observed values for the random variable Y_i

In order to find the least square estimators for α and β we need to minimise $S(\alpha, \beta)$ (for fixed y 's and x 's) with respect to the parameters α and β .

That is we find $\frac{\partial S}{\partial \alpha}$ and $\frac{\partial S}{\partial \beta}$ and we set them equal to zero.

$$\frac{\partial S}{\partial \alpha} = \sum_{i=1}^n -2(y_i - \hat{\alpha} - \hat{\beta} x_i) = -2 \left(\sum_{i=1}^n y_i - n\hat{\alpha} - \hat{\beta} \sum_{i=1}^n x_i \right) = 0$$

$$\frac{\partial S}{\partial \beta} = \sum_{i=1}^n -2x_i(y_i - \hat{\alpha} - \hat{\beta} x_i) = -2 \left(\sum_{i=1}^n x_i y_i - \hat{\alpha} \sum_{i=1}^n x_i - \hat{\beta} \sum_{i=1}^n x_i^2 \right) = 0$$

with solutions

$$\hat{\alpha} = \bar{y} - \hat{\beta} \bar{x}$$

$$\hat{\beta} = \frac{\sum_{i=1}^n x_i y_i - n\bar{x}\bar{y}}{\sum_{i=1}^n x_i^2 - n\bar{x}^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

The quantities $\hat{y}_i = \hat{\alpha} + \hat{\beta} x_i$ are called the **fitted values**.

The quantities $\hat{\varepsilon}_i = y_i - \hat{y}_i$ are called the **residuals**.

13.3 Properties of the least square estimators

Note that both $\hat{\alpha}$ and $\hat{\beta}$ are linear functions of the y 's. For example for $\hat{\beta}$ we have

$$\hat{\beta} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})y_i}{S_{xx}} = \sum_{i=1}^n C_i y_i$$

where $S_{xx} = \sum (x_i - \bar{x})^2$ and $C_i = \frac{(x_i - \bar{x})}{S_{xx}}$.

(Prove the above statement for $\hat{\alpha}$).

13.3.1 Expected values for $\hat{\alpha}$ and $\hat{\beta}$

i) $E(\hat{\beta}) = \beta$: $\hat{\beta}$ is an unbiased estimator of β .

Proof

$$\begin{aligned}
 E(\hat{\beta}) &= \frac{\sum_{i=1}^n (x_i - \bar{x}) E(Y_i)}{S_{xx}} = \frac{\sum_{i=1}^n (x_i - \bar{x}) (a + \beta x_i)}{S_{xx}} \quad (1) \\
 &= \frac{\beta \sum_{i=1}^n (x_i - \bar{x}) x_i}{S_{xx}} \quad (2) \\
 &= \beta \quad (3)
 \end{aligned}$$

(1) since if $y = \sum c_i z_i \Rightarrow E(y) = \sum c_i E(z_i)$

(2) since $\sum (x_i - \bar{x}) a = 0$ (prove it)

(3) since $S_{xx} = \sum (x_i - \bar{x})^2 = \sum (x_i - \bar{x}) x_i$ (prove it)

ii) $E(\hat{\alpha}) = \alpha$: $\hat{\alpha}$ is an unbiased estimator of α .

Proof:
$$n\hat{\alpha} = \sum_{i=1}^n Y_i - \hat{\beta} \sum_{i=1}^n x_i$$

$$\begin{aligned}
E(n\hat{\alpha}) &= n E(\hat{\alpha}) = \sum_{i=1}^n E(y_i) - E(\hat{\beta}) \sum_{i=1}^n x_i \\
&= \sum_{i=1}^n (\alpha + \beta x_i) - \beta \sum_{i=1}^n x_i \\
&= n\alpha + \beta \sum_{i=1}^n x_i - \beta \sum_{i=1}^n x_i \\
&= n\alpha
\end{aligned}$$

or

$$E(\tilde{\alpha}) = \alpha$$

13.3.2 The Variances of $\hat{\alpha}$ and $\hat{\beta}$

$$\text{i) } \text{Var}(\hat{\beta}) = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sigma^2}{S_{xx}}$$

Proof.

$$\begin{aligned}
\text{Var}(\hat{\beta}) &= \left\{ \frac{\sum_{i=1}^n (x_i - \bar{x})}{S_{xx}} \right\}^2 \text{var}(Y_i) \\
&= \frac{\sigma^2}{S_{xx}}
\end{aligned}$$

so

$$\text{Var}(\hat{\beta}) = \frac{\hat{\sigma}^2}{S_{xx}}$$

$$\text{ii) } \text{var}(\hat{\alpha}) = \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right]$$

Proof.

$$\text{var}(\hat{\alpha}) = \text{var}(\bar{y}) + \bar{x}^2 \text{var}(\hat{\beta}) - 2\bar{x} \text{cov}(\bar{y}, \hat{\beta})$$

But $\text{cov}(\bar{y}, \hat{\beta}) = 0$ (see Exercise 13.2), so we have

$$\begin{aligned}\text{var}(\hat{\alpha}) &= \frac{\sigma^2}{n} + \bar{x}^2 \frac{\sigma^2}{S_{xx}} \\ &= \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right]\end{aligned}$$

hence

$$\hat{\text{var}}(\hat{\alpha}) = \hat{\sigma}^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right]$$

13.3.3 The Gauss-Markoff theorem

The least-squares estimators $\hat{\alpha}$ and $\hat{\beta}$ have minimum variances among all the linear unbiased estimators.

13.3.4 The Normality assumption of $\hat{\alpha}$ and $\hat{\beta}$

Note that if Y is a linear function of normally distributed variables U_i i.e.

$$Y = c_1 U_1 + c_2 U_2$$

Y will be Normally distributed i.e.

$$Y \sim N(\mu, \sigma^2).$$

The L.S. estimators $\hat{\alpha}$ and $\hat{\beta}$ are linear functions of Y_i which is

$$Y_i \sim N(\alpha + \beta x_i, \sigma^2)$$

so $\hat{\alpha}$ and $\hat{\beta}$ will be Normally distributed as

$$\hat{\alpha} \sim N\left(\alpha, \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right]\right)$$

$$\hat{\beta} \sim N\left(\beta, \frac{\sigma^2}{S_{xx}}\right)$$

13.4 Hypothesis testing

13.4.1 Estimation of σ^2

$$\hat{\sigma}^2 = s^2 = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n-2} = \frac{\sum_{i=1}^n \hat{\varepsilon}_i^2}{n-2}$$

where $\hat{y}_i = \hat{\alpha} + \hat{\beta} x_i$ are the fitted values $\hat{\varepsilon}_i = y_i - \hat{y}_i$ the residuals and $n - 2$ are the residual degrees of freedom (df).

13.4.2 *t*-test for $\hat{\beta}$ and $\hat{\alpha}$

$$\text{Note that if } \left. \begin{array}{l} z \sim N(0,1) \\ \omega \sim \chi^2(\eta) \end{array} \right\} \text{ then } t = \frac{z}{\sqrt{\frac{\omega}{\eta}}} \sim t(\eta)$$

and z and ω are independent we have that

$$\frac{\hat{\beta} - \beta}{\frac{\sigma}{\sqrt{S_{xx}}}} \sim N(0,1)$$

and that

$$\frac{\sum (y_i - \hat{y}_i)^2}{\sigma^2} \sim \chi^2(n-2)$$

Also $\hat{\beta}$ and $\sum (y_i - \hat{y})^2$ are independent (not proven). So

$$t = \frac{\frac{\hat{\beta} - \beta}{\frac{\sigma}{\sqrt{S_{xx}}}}}{\sqrt{\frac{\sum (y_i - \hat{y})^2}{\sigma^2} \cdot \frac{1}{n-2}}} = \frac{\hat{\beta} - \beta}{\frac{\hat{\sigma}}{\sqrt{S_{xx}}}} = \frac{\hat{\beta} - \beta}{se(\hat{\beta})} \sim t(n-2)$$

where $se(\hat{\beta}) = \frac{s}{\sqrt{S_{xx}}} = \frac{\hat{\sigma}}{\sqrt{S_{xx}}}$ is the standard error of $\hat{\beta}$.

We can test hypothesis for β using the statistics t .

For example to test

$$H_0 : \beta = \beta_o \qquad H_1 : \beta \neq \beta_o$$

calculate

$$t = \frac{\hat{\beta} - \beta_o}{se(\hat{\beta})}$$

Now if

$$|t| > t_{n-2, \frac{a}{2}}$$

reject the null hypothesis H_0 and accept the alternative H_1 , otherwise accept H_0 .

Note that a is the *significant level* of the test and not the constant parameter α of the linear model.

To test hypothesis about α i.e.

$$H_0 : \alpha = \alpha_o \qquad H_1 : \alpha \neq \alpha_o$$

use the test statistic

$$t = \frac{\hat{\alpha} - \alpha_o}{se(\hat{\alpha})}.$$

13.4.3 C.I. for α and β

A $(1-a)100\%$ C. I. for β is given by

$$\hat{\beta} \pm t_{n-2, \frac{a}{2}} \times se(\hat{\beta})$$

and for α is given by

$$\hat{\alpha} \pm t_{n-2, \frac{a}{2}} \times se(\hat{\alpha})$$

13.5 Prediction and Confidence Intervals

13.5.1 Confidence Intervals for $\mu_o = a + bx_o$

Note that the expected value for y_o the value of the y-variable when the explanatory variable is at x_o is

$$E(y_o) = \mu_o = \alpha + \beta x_o \quad :$$

The fitted value at the point x_o is defined as

$$\hat{y}_o = \hat{\mu}_o = \hat{\alpha} + \hat{\beta} x_o = \bar{y} - \hat{\beta} \bar{x} + \hat{\beta} x_o = \bar{y} + \hat{\beta} (x_o - \bar{x})$$

with expected values

$$E(\hat{y}_o) = \alpha + \beta x_o = \mu_o \text{ as } E(\hat{\alpha}) = \alpha \text{ and } E(\hat{\beta}) = \beta$$

So the fitted value \hat{y}_o is unbiased for μ_o . The variance for \hat{y}_o is

$$\begin{aligned}\text{var}(\hat{y}_o) &= \text{var}(\bar{y}) + \text{var}(\hat{\beta})(x_o - \bar{x})^2 + 2(x_o - \bar{x})\text{cov}(\bar{y}, \hat{\beta}) \\ &= \frac{\sigma^2}{n} + \frac{\sigma^2}{S_{xx}}(x_o - \bar{x})^2 = \sigma^2 \left[\frac{1}{n} + \frac{(x_o - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right]\end{aligned}$$

so an estimate for the variance is given by.

$$\text{var}(\hat{y}_o) = s^2 \left[\frac{1}{n} + \frac{(x_o - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right]$$

where $s^2 = \sum (y_i - \hat{y})^2 / n - 2$.

Since \hat{y} is a linear combination of Normally distributed variables, it is Normally distributed; i.e.

$$\begin{aligned}\hat{y}_o &\sim N \left(\mu_o, \sigma^2 \left[\frac{1}{n} + \frac{(x_o - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right] \right) \\ \Rightarrow z_o &= \frac{\hat{y}_o - \mu_o}{\sqrt{\sigma^2 \left[\frac{1}{n} + \frac{(x_o - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right]}} \sim N(0,1)\end{aligned}$$

or

$$t = \frac{\hat{y}_o - \mu_o}{S^2 \left[\frac{1}{n} + \frac{(x_o - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right]} \sim t_{n-2}$$

so a C.I for μ_o is given by

$$\hat{y}_o \pm t_{n-2, \frac{\alpha}{2}} \times \text{se}(\hat{y}_o)$$

where

$$se(\hat{y}_o) = S \left[\frac{1}{n} + \frac{(x_o - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right]^{\frac{1}{2}}$$

13.52 Prediction Interval for y_o^* / x_o , a future observation for y_o .

Let y_o^* / x_o denote a future observation of the y-variable at the x-variable value x_o . Then

$$E(y_o^* / x_o) = \alpha + \beta x_o = \mu_o$$

Since $\hat{y}_o = \hat{\alpha} + \hat{\beta} x_o$ is an unbiased estimator for μ_o it can be used to predict the mean of a future observation y_o^* / x_o .

In general in order to evaluate how good our predictor \hat{y} is for predicting a further observation y^* we have to know the mean square error for prediction or PSE.

Definition: $PSE(y^*) = E(y^* - \hat{y})^2$

Theorem: Let \hat{y} be an estimate of μ and let y^* be a new observation such that $E(y^*) = \mu$. Then $PSE(y^*) = Var(y^*) + MSE(\hat{y})$ where $MSE(\hat{y}) = E(\hat{y} - \mu)^2$.

Proof:

$$\begin{aligned} PSE(y^*) &= E(y^* - \hat{y})^2 \\ &= E[(y^* - \mu) - (\hat{y} - \mu)]^2 \\ &= E[(y^* - \mu)^2 - 2(y^* - \mu)(\hat{y} - \mu) + (\hat{y} - \mu)^2] \end{aligned}$$

y^* is independent of \hat{y} , as y^* is a new observation. Hence

$$E[(y^* - \mu)(\hat{y} - \mu)] = E(y^* - \mu)E(\hat{y} - \mu) = 0 \text{ as } E(y^*) = \mu$$

Hence

$$\begin{aligned} PSE(y^*) &= E[(y^* - \mu)^2 + E(\hat{y} - \mu)^2] \\ &= Var(y^*) + MSE(\hat{y}) \\ &= Var(y^*) + Var(\hat{y}) + (bias)^2 \end{aligned}$$

In the simple linear regression example we have

$$\begin{aligned}
 E(y_o^* - \hat{y}_o)^2 &= \text{Var}(y_o^*) + \text{Var}(\hat{y}_o) \quad \text{since } \hat{y}_o \text{ is unbiased} \\
 &= \sigma^2 + \sigma^2 \left[\frac{1}{n} + \frac{(x_o - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right] \\
 &= \sigma^2 \left[1 + \frac{1}{n} + \frac{(x_o - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right]
 \end{aligned}$$

A $100(1-a)\%$ prediction interval for y_o^* / x_o is given by

$$\hat{y}_o \pm t_{n-2, \frac{a}{2}} s \left[1 + \frac{1}{n} + \frac{(x_o - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right]^{\frac{1}{2}}$$

13.6 Maximum likelihood estimation of the parameters α , β and σ^2 in the simple linear regression.

The **likelihood function** is the probability of observing the sample seeing as a function of the parameter rather than a function of the random variables.

For independent random variables $x_1, x_2 \dots x_n$ the likelihood will be

$$L(\theta) = f(x_1; \theta) f(x_2; \theta) \dots f(x_n; \theta)$$

In the simple regression model we have

$$Y_i = \alpha + \beta x_i + \varepsilon_i$$

where

$$\varepsilon_i \stackrel{\text{ind}}{\sim} N(0, \sigma^2) \Rightarrow y_i \stackrel{\text{ind}}{\sim} N(\alpha + \beta x_i, \sigma^2)$$

i.e.

$$E(Y_i) = \alpha + \beta x_i$$

$$\text{var}(Y_i) = \sigma^2$$

and $\theta = (\alpha, \beta, \sigma^2)$.

The likelihood for one observation is

$$L(\alpha, \beta, \sigma^2; y_i) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(y_i - \alpha - \beta x_i)^2\right\}$$

for n independent observations the likelihood will be

$$\begin{aligned} L(\alpha, \beta, \sigma^2 / y_1 \dots y_n) &= \prod_{i=1}^n \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) \exp\left\{-\frac{1}{2\sigma^2}(y_i - \alpha - \beta x_i)^2\right\} \\ &= \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right)^n \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \alpha - \beta x_i)^2\right\}. \end{aligned}$$

Note that $S = \sum_{i=1}^n (y_i - \alpha - \beta x_i)^2$ is the function that we minimised in the least square estimation approach.

In order to find the MLE's for α , β and σ^2 we have to maximise $L(\alpha, \beta, \sigma^2)$ with respect to the parameters or equivalently maximise $\log L(\alpha, \beta, \sigma^2) = \ell(\alpha, \beta, \sigma^2)$

Now

$$\ell(\alpha, \beta, \sigma^2) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \alpha - \beta x_i)^2$$

so we differentiate with respect to α , β and σ^2

$$\frac{\partial \ell}{\partial \alpha} = \frac{1}{\sigma^2} \sum_{i=1}^n (y_i - \hat{\alpha} - \hat{\beta} x_i) = 0$$

$$\frac{\partial \ell}{\partial \beta} = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i (y_i - \hat{\alpha} - \hat{\beta} x_i)) = 0$$

$$\frac{\partial \ell}{\partial \sigma^2} = \frac{n}{2\hat{\sigma}^2} + \frac{1}{2\hat{\sigma}^4} \sum_{i=1}^n (y_i - \hat{\alpha} - \hat{\beta} x_i)^2 = 0$$

solving for $\hat{\alpha}, \hat{\beta}, \hat{\sigma}^2$ we have

$$\hat{\alpha} = \bar{y} - \hat{\beta} \bar{x}$$

$$\hat{\beta} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$\hat{\sigma}^2 = \frac{\sum (y_i - \hat{\alpha} - \hat{\beta} x_i)^2}{n}$$

Note:

- i) $\hat{\alpha}$ and $\hat{\beta}$ are also the least-square estimators. This is because the maximisation of the log-likelihood (for fixed σ^2) is the equivalent of the minimisation of the least-square quantity $S = \sum_{i=1}^n (y_i - \alpha - \beta x_i)^2$.

- ii) We generally prefer to use an unbiased estimator of σ^2 given by

$$s^2 = \frac{\sum_{i=1}^n (y_i - \hat{\alpha} - \hat{\beta} x_i)^2}{n-2} = \frac{D}{df} \quad \leftarrow \text{deviance}$$

↑ the residual degrees of freedom

Exercise 13.1: Simple linear regression theory

- a) Consider the simple linear regression model of the form

$$Y_i = a + bx_i + \varepsilon_i \quad \text{for } i = 1, 2, \dots, n$$

where

Y_i is the response variable,

x_i is the independent variable,

a and b are parameters to be estimated.

ε_i , for $i = 1, 2, \dots, n$ are independent Normally distributed variables with mean 0 and variance σ^2 .

- i) Find the likelihood function for a single observation and hence show that the log-likelihood function for all n observations from the above model is

$$l(a, b, \sigma^2) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - a - bx_i)^2$$

- ii) Give the Normal equations used to find the Maximum likelihood estimators for the parameters a , b and σ^2 , and state the resulting maximum likelihood estimators of a , b and σ^2 .

b)

- i) State the distribution of \hat{b} and $\frac{\sum_{i=1}^n \hat{e}_i^2}{\sigma^2}$ where $\hat{e}_i = y_i - \hat{a} - \hat{b}x_i$, i.e. the residual for the i th observation and the numerator in the expression is the Residual Sum of Squares (RSS). Note that the variance of \hat{b} is $\frac{\sigma^2}{S_{xx}}$

$$\text{where } S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2.$$

- ii) We know that \hat{b} and $\sum_{i=1}^n \hat{e}_i^2$ are independent.

We also know that if $z \sim N(0,1)$ and $w \sim \chi_v^2$, independently,

$$\text{then } \frac{z}{\sqrt{(w/v)}} \sim t_v.$$

Use this to construct a $100(1-\alpha)\%$ confidence interval for b .

Exercise 13.2: Simple linear regression theory

For a simple linear regression model, prove that $cov(\bar{y}, \hat{\beta}) = 0$

Note: an outlined method is as follows:

$$\begin{aligned}
 cov(\bar{y}, \hat{\beta}) &= E[(\bar{y} - \mu)(\hat{\beta} - \beta)] \\
 &= E\left[(\bar{y} - \mu) \left(\frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} - \beta \right)\right] \\
 &= E\left[(\bar{y} - \mu) \left(\frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} - \beta \right)\right] - E[(\bar{y} - \mu)\beta] \\
 &= E\left[(\bar{y} - \mu) \left(\frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} - \beta \right)\right] - 0 \quad (\text{Why?})
 \end{aligned}$$

Now

$$\begin{aligned}
 E[(y_i - \bar{y})(\bar{y} - \mu)] &= E\left[\left((n-1)y_i - y_1 - \dots - y_{i-1} - y_{i+1} - \dots - y_n\right)(y_1 + \dots + y_n - n\mu)/n^2\right] \\
 &= E\left[\left((n-1)\{y_i - \mu\} - \{y_1 - \mu\} - \dots - \{y_{i-1} - \mu\} - \{y_{i+1} - \mu\} - \dots - \{y_n - \mu\}\right)(\{y_1 - \mu\} + \dots + \{y_n - \mu\})/n^2\right] \\
 &= (n-1)\sigma^2 - \sigma^2 - \sigma^2 \dots - \sigma^2 \\
 &= 0
 \end{aligned}$$

as

$$E[(y_i - \mu)(y_j - \mu)] = 0 \quad (\text{Why?})$$

Hence deduce $cov(\bar{y}, \hat{\beta}) = 0$.

Exercise 13.3: Simple linear regression theory

For the regression $y_i = a + bx_i + \varepsilon_i$, $\varepsilon_i \sim N(0,1)$ show that the least squares estimate of b is given by

$$\hat{b} = \frac{S_{xy} - n\bar{x}\bar{y}}{S_{xx} - n\bar{x}^2}.$$

A test of $b=0$ can be based upon

$$\begin{aligned} T &= S_{xy} - n\bar{y}\bar{x} = \sum_{i=1}^n x_i y_i - \sum_{i=1}^n y_i \bar{x} \\ &= \sum_{i=1}^n y_i (x_i - \bar{x}) \end{aligned}$$

Show that, under $H_0 : b = 0$, we have $E(T) = 0$.

Also show that $V(T) = \sigma^2 \sum (x_i - \bar{x})^2$.

Deduce that $\frac{S_{xy} - n\bar{y}\bar{x}}{\sigma \sqrt{\sum (x_i - \bar{x})^2}} \sim N(0,1)$

And hence that $\frac{S_{xy} - n\bar{y}\bar{x}}{s \sqrt{\sum (x_i - \bar{x})^2}} \sim t_{n-1}$

Practical 13: Simple Linear Regression

The data set in the file SHARED (K):\SCTMS\SOM\MA2010\REGRESSION\FOOT GESTATION TIME.SAV comprises measurements of foetal foot length in mm (Y) and gestational age in weeks (X) for 450 fetuses.

1. Produce a scatter plot of Y against X using the procedure
> Graphs > Scatter > Simple Scatter > Define | Y Axis 'foot' | X Axis 'gest'
> OK.
Comment on this plot.
2. Fit a simple linear regression line using
> Analyse > Regression > Linear | Dependent: 'foot' | Independent(s): 'gest' to declare your y and x variable and fit the model.
You can use the PLOTS option to get the residual plots.
> Plots | Y: ZRESID | X: ZPRED | ✓ Histogram | ✓ Normal probability plot
> Continue > OK
 - i) State the model fitted and its parameter estimates. Interpret these estimates.
 - ii) Test whether there is a linear relationship between foot length and gestational age.
 - iii) State the assumptions necessary for your model to be valid.
 - iv) Do the residual plots show that any of the assumptions does not hold?