

In order to be able to assess the "goodness" of our estimation:
 $\hat{\alpha} = \alpha$, $\hat{\beta} = b$, we need to make some assumptions concerning the distribution of the random variables we are dealing with and about the true nature of the regression line.

Inferences based on the L.S.E
The L.S.E (least square estimates) are used when the relationship between x and the mean of y (EY) is linear or close enough.

Assumptions:

(12)

- Regression is linear

- y_1, y_2, \dots, y_n are independent

$$\cdot y_i \sim N(\alpha + \beta x_i, \sigma^2), i=1, 2, \dots, n$$



Mean is
Variance constant

The deterministic,
i.e. no random

Part of the regression

$$y_i = \alpha + \beta x_i + \epsilon_i \quad (i=1, \dots, n)$$

Model for
Straight
line regression

Want $\epsilon_i \sim N(0, \sigma^2)$

Show that if $\epsilon_i \sim N(0, \sigma^2)$ then

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

Inference on LSE.

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i - (\hat{a} + \hat{b}x_i))^2$$

residual sum of squares

$\hat{\sigma}$ - standard error of estimate

$$\text{OR} \quad \hat{\sigma}^2 = \frac{S_{yy} - (S_{xy})^2/S_{xx}}{n-2} = \text{MSE}$$

Another test statistic for hypothesis

$$t = \frac{\hat{a} - a}{\hat{\sigma}} \cdot \sqrt{\frac{n S_{xx}}{S_{xx} + n(\bar{x})^2}} \sim t(n-2)$$

$$t = \frac{\hat{\beta}_1 - \beta}{\hat{\sigma}} \cdot \sqrt{S_{xx}} \sim t(n-2)$$

There are two statistics we will use.

to do inference on a and β . That

i. Confidence Intervals and Hypothesis tests. Point Est. have been already derived.

Confidence Intervals:

(14)

For a given confidence level α ,
 $\alpha (1-\alpha) \cdot 100\%$ confidence interval,
 $(C.I)$ for the regression parameters
 α and β is given by:

$$\hat{\alpha} \in (\hat{\alpha} + t_{\alpha/2} \cdot S_{\hat{\alpha}}, \hat{\alpha} - t_{\alpha/2} \cdot S_{\hat{\alpha}})$$

$$\hat{\beta} \in (\hat{\beta} + t_{\alpha/2} \cdot \frac{S_{\hat{\beta}}}{\sqrt{n}}, \hat{\beta} - t_{\alpha/2} \cdot \frac{S_{\hat{\beta}}}{\sqrt{n}})$$

Remember:

$$\hat{\alpha} = \bar{y} - b \cdot \bar{x}$$

$$\hat{\beta} = b = \frac{\sum xy}{\sum x^2}$$

$$S_{\hat{\alpha}} = \sqrt{\frac{\sum (y_i - \hat{y}_i)^2}{n-2}} \cdot \frac{1}{\sqrt{\sum x_i^2}}$$

$$S_{\hat{\beta}} = \sqrt{\frac{\sum (y_i - \hat{y}_i)^2}{n-2}} \cdot \frac{1}{\sqrt{n-2}}$$

Confidence Interval for the α (IS)

Example 3 (continue)

95% C.I for the regression coefficient α .

We know from Example 3 that

$$n = 10, \bar{X} = 200, S_{xx} = 132000, S_{yy} = 2.13745 \\ S_{xy} = 505.40.$$

Then:

$$\hat{\sigma}^2 = \frac{2.13745 - \frac{(505.40)^2}{132000}}{8} = \frac{0.20238}{8} = 0.0253$$

$$\Rightarrow \hat{\delta} = \sqrt{0.0253} = 0.159$$

standard error of
estimate.

$$t_{0.025, 8} = 2.306, \text{ a } 95\% \text{ C.I for } \alpha$$

is given by

$$\alpha \in \left(\hat{\alpha} \pm t_{\alpha/2, n-2} \cdot \hat{\sigma} \cdot \sqrt{\frac{1}{n} + \frac{(\bar{X})^2}{S_{xx}}} \right) = \left(0.069 \pm 2.306 \cdot 0.159 \cdot \sqrt{\frac{1}{10} + \frac{200^2}{132,000}} \right) = (-0.164, 0.302) \text{ through origin!}$$

β - slope of regression = the change
in the mean of Y corresponding to a
unit increase in x .

$$\beta = 0$$

Line horizontal! EY does not
depend on x

$$\beta > 0$$

Positive association
large (small) values of Y
are associated with
large (small) values of x

$$\beta < 0$$

Negative association
large (small) values of Y
are associated with
small (large) values of x .

Hence it is of interest to test

whether $\beta = 0$.

Hypothesis test on β .

(17)

$$H_0 : \beta = \beta_0 \quad t_0 = \frac{\hat{\beta} - \beta}{\sigma} \cdot \sqrt{S_{xx}} \sim t_{(n-2)}$$

$$H_\alpha : \beta \neq \beta_0$$

Reject H_0 if $t_0 < t_{\alpha/2, n-2}$ or $t_0 > t_{\alpha/2, n-2}$.

Example 3 continue II

$$H_0 : \beta = 0 \quad \alpha = 0.05 \Rightarrow t_{0.025, 8} = 2.306$$
$$H_\alpha : \beta \neq 0$$

$$t_0 = \frac{0.00383 - 0}{0.159} \cdot \sqrt{132000} = 8.75$$

Since $t_0 = 8.75 > 2.306 = t_{0.025, 8}$ we

reject the null hypothesis and conclude:

that there is a relationship between air velocity and the average evaporation coefficient. Scatter plot tells me is a linear relationship.

(18)

Mean Response

We are interested in estimating the regression $\alpha + \beta x$, i.e. the mean of the distribution of Y for a given value of x , say x_0 . It is reasonable to do so by using

$$\hat{Y} = \hat{\alpha} + \hat{\beta} \cdot x_0 = \alpha + b x_0, \quad [a, b - L.S.E]$$

A. (1-a).100% C.I. for the mean response -

$$Y_e \left(\hat{Y} \pm t_{\alpha/2} \cdot \hat{\sigma} \cdot \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_x^2}} \right)$$

Example 3.

Construct a 95% C.I. for

the mean evaporation coefficient when the air velocity is 190 cm/s.

$$Y_e (0.80 + 2.306 * 0.159) \sqrt{\frac{1}{10} + \frac{(190 - 200)^2}{132,000}}$$

$$= (1.068, 0.92)$$

Hence a 95% C.I for the mean evaporation coefficient of $\alpha + \beta x_0$ is $(0.68, 0.92)$

Prediction of new value

Even more important than estimating $\alpha + \beta x_0$ is the prediction of a future value of Y when $x = x_0$ and x_0 is within the range of experimentation.

Method: We want to predict a future obs Y when $x = x_0$ with probability α . If α and β were known, then since

$$Y \sim N(\alpha + \beta x_0, \sigma^2) \Leftrightarrow Y - \alpha - \beta x_0 \sim N(0, \sigma^2)$$

BUT α, β are NOT known and estimated by $\hat{\alpha} = a, \hat{\beta} = b$. At $(1-\alpha)\cdot 100\%$ C.I. $a + b x_0$

$$Y_{\text{pred}} \in \left(\hat{Y}_{\text{pred}} \pm t_{\alpha/2, n-2} \cdot \hat{\sigma} \cdot \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}} \right)$$

Example 3 A 95% "prediction" interval (20)

for an observation of the evaporation coefficient when air velocity = 190 cm/s.

$$\text{Known } (0.80 \pm 2.306 * 0.159 \cdot \sqrt{1 + \frac{1}{10} + \frac{(190-200)^2}{132000}}) = \\ = (0.41, 1.19)$$

Interval of prediction is wider than the C.I for the mean response.

More vigorous "proofs" of the facts presented here follow closely the arguments we have already presented in Ch. 5 and C.T. If you want to have a look they are repeated at Ch. 7 of your book!!!