

Let  $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 x^*$ , where  $x^*$  is some fixed value of  $x$ . Then  
1. The mean value of  $\hat{Y}$  is  
$$E(\hat{Y}) = E(\hat{\beta}_0 + \hat{\beta}_1 x^*) = \mu_{\hat{\beta}_0 + \hat{\beta}_1 x^*} = \beta_0 + \beta_1 x^*$$
  
Thus  $\hat{\beta}_0 + \hat{\beta}_1 x^*$  is an unbiased estimator for  $\beta_0 + \beta_1 x^*$  (i.e., for  $\mu_{Y, x^*}$ ).  
2. The variance of  $\hat{Y}$  is  
$$V(\hat{Y}) = \sigma^2 = \sigma^2 \left[ \frac{1}{n} + \frac{(x^* - \bar{x})^2}{\sum x_i^2 - (\sum x_i)^2/n} \right] = \sigma^2 \left[ \frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}} \right]$$
  
and the standard deviation  $\sigma_{\hat{Y}}$  is the square root of this expression. The estimated standard deviation of  $\hat{\beta}_0 + \hat{\beta}_1 x^*$ , denoted by  $s_{\hat{Y}}$  or  $s_{\hat{\beta}_0 + \hat{\beta}_1 x^*}$ , results from replacing  $\sigma$  by its estimate  $s$ .

$$s_{\hat{Y}} = s_{\hat{\beta}_0 + \hat{\beta}_1 x^*} = s \sqrt{\frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}}}$$

3.  $\hat{Y}$  has a normal distribution.

Proposition

The variable

$$T = \frac{\hat{\beta}_0 + \hat{\beta}_1 x^* - (\beta_0 + \beta_1 x^*)}{S_{\hat{\beta}_0 + \hat{\beta}_1 x^*}} = \frac{\hat{Y} - (\beta_0 + \beta_1 x^*)}{s_{\hat{Y}}} \quad (12.5)$$

has a  $t$  distribution with  $n - 2$  df.

A  $100(1 - \alpha)\%$  CI for  $\mu_{Y, x^*}$ , the expected value of  $Y$  when  $x = x^*$ , is

$$\hat{\beta}_0 + \hat{\beta}_1 x^* \pm t_{\alpha/2, n-2} \cdot s_{\hat{\beta}_0 + \hat{\beta}_1 x^*} = \hat{Y} \pm t_{\alpha/2, n-2} \cdot s_{\hat{Y}} \quad (12.6)$$

A  $100(1 - \alpha)\%$  PI for a future  $Y$  observation to be made when  $x = x^*$  is

$$\begin{aligned} \hat{\beta}_0 + \hat{\beta}_1 x^* \pm t_{\alpha/2, n-2} \cdot s \sqrt{1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}}} \\ = \hat{\beta}_0 + \hat{\beta}_1 x^* \pm t_{\alpha/2, n-2} \cdot \sqrt{s^2 + s_{\hat{\beta}_0 + \hat{\beta}_1 x^*}^2} \\ = \hat{Y} \pm t_{\alpha/2, n-2} \cdot \sqrt{s^2 + s_{\hat{Y}}^2} \end{aligned} \quad (12.7)$$

PI: Prediction Interval for future  $Y$

Inference U\_yx\* and prediction of Y

$$s^2 = S_{yy} / (n-1)$$
$$S_{xy} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \sum_{i=1}^n x_i y_i - \frac{\left(\sum_{i=1}^n x_i\right)\left(\sum_{i=1}^n y_i\right)}{n}$$

The sample correlation coefficient for the  $n$  pairs  $(x_1, y_1), \dots, (x_n, y_n)$  is

$$r = \frac{S_{xy}}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}} = \frac{S_{xy}}{\sqrt{S_{xx}} \sqrt{S_{yy}}} \quad (12.8)$$

### Properties of $r$

The most important properties of  $r$  are as follows:

- The value of  $r$  does not depend on which of the two variables under study is labeled  $x$  and which is labeled  $y$ .
- The value of  $r$  is independent of the units in which  $x$  and  $y$  are measured.
- $-1 \leq r \leq 1$
- $r = 1$  if and only if (iff) all  $(x, y)$  pairs lie on a straight line with positive slope, and  $r = -1$  iff all  $(x, y)$  pairs lie on a straight line with negative slope.
- The square of the sample correlation coefficient gives the value of the coefficient of determination that would result from fitting the simple linear regression model—in symbols,  $(r)^2 = r^2$ .

Not dependent on unit

## The Sample Correlation Coefficient $r$

$$\rho = \rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \cdot \sigma_Y}$$

$$\text{Cov}(X, Y) = \begin{cases} \sum_x \sum_y (x - \mu_x)(y - \mu_y) p(x, y) & (X, Y) \text{ discrete} \\ \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x - \mu_x)(y - \mu_y) f(x, y) dx dy & (X, Y) \text{ continuous} \end{cases}$$

$$\hat{\rho} = R = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum (X_i - \bar{X})^2} \sqrt{\sum (Y_i - \bar{Y})^2}}$$

The joint probability distribution of  $(X, Y)$  is specified by

$$f(x, y) = \frac{1}{2\pi \cdot \sigma_1 \sigma_2 \sqrt{1 - \rho^2}} e^{-\frac{1}{2(1-\rho^2)} \{ [(x-\mu_1)/\sigma_1]^2 - 2\rho[(x-\mu_1)/\sigma_1][(y-\mu_2)/\sigma_2] + [(y-\mu_2)/\sigma_2]^2 \}} \quad (12.9)$$

where  $\mu_1$  and  $\sigma_1$  are the mean and standard deviation of  $X$ , and  $\mu_2$  and  $\sigma_2$  are the mean and standard deviation of  $Y$ ;  $f(x, y)$  is called the bivariate normal probability distribution.

### Testing for the Absence of Correlation

When  $H_0: \rho = 0$  is true, the test statistic

$$T = \frac{R\sqrt{n-2}}{\sqrt{1-R^2}}$$

has a  $t$  distribution with  $n - 2$  df.

Alternative Hypothesis

Rejection Region for Level  $\alpha$  Test

$$H_a: \rho > 0$$

$$t \geq t_{\alpha, n-2}$$

$$H_a: \rho < 0$$

$$t \leq -t_{\alpha, n-2}$$

$$H_a: \rho \neq 0$$

$$\text{either } t \geq t_{\alpha/2, n-2} \text{ or } t \leq -t_{\alpha/2, n-2}$$

A  $P$ -value based on  $n - 2$  df can be calculated as described previously.

When  $(X_1, Y_1), \dots, (X_n, Y_n)$  is a sample from a bivariate normal distribution, the rv

$$V = \frac{1}{2} \ln \left( \frac{1+R}{1-R} \right) \quad (12.10)$$

has approximately a normal distribution with mean and variance

$$\mu_V = \frac{1}{2} \ln \left( \frac{1+\rho}{1-\rho} \right) \quad \sigma_V^2 = \frac{1}{n-3}$$

The test statistic for testing  $H_0: \rho = \rho_0$  is

$$Z = \frac{V - \frac{1}{2} \ln[(1 + \rho_0)/(1 - \rho_0)]}{1/\sqrt{n-3}}$$

Alternative Hypothesis

Rejection Region for Level  $\alpha$  Test

$$H_a: \rho > \rho_0$$

$$z \geq z_\alpha$$

$$H_a: \rho < \rho_0$$

$$z \leq -z_\alpha$$

$$H_a: \rho \neq \rho_0$$

$$\text{either } z \geq z_{\alpha/2} \text{ or } z \leq -z_{\alpha/2}$$

A  $P$ -value can be calculated in the same manner as for previous  $z$  tests.

A  $100(1 - \alpha)\%$  confidence interval for  $\rho$  is

$$\left( \frac{e^{2c_1} - 1}{e^{2c_1} + 1}, \frac{e^{2c_2} - 1}{e^{2c_2} + 1} \right)$$

where  $c_1$  and  $c_2$  are the left and right endpoints, respectively, of the interval (12.11).

Other inference concerning rho

W12

The simple regression model

Relationship of 2 variables

$x$ : independent, predictor, explanatory variable

$y$ : dependent, response variable

$$Y = \beta_0 + \beta_1 x + \epsilon$$

$\epsilon$

random deviation, random error term

True regression line

$$y = \beta_0 + \beta_1 x$$

$\mu_{Y, x^*}$  = the expected (or mean) value of  $Y$  when  $x$  has value  $x^*$

$\sigma_{Y, x^*}^2$  = the variance of  $Y$  when  $x$  has value  $x^*$

$$\mu_{Y, x^*} = E(\beta_0 + \beta_1 x^* + \epsilon) = \beta_0 + \beta_1 x^* + E(\epsilon) = \beta_0 + \beta_1 x^*$$

$$\sigma_{Y, x^*}^2 = V(\beta_0 + \beta_1 x^* + \epsilon) = V(\beta_0 + \beta_1 x^*) + V(\epsilon) = 0 + \sigma^2 = \sigma^2$$

### Assumption

Linear relationship between  $x$  and  $y$

Error is normally distributed

Mean error is 0

Constant variance

Observation are independent

### Least squares estimates

$$f(h_0, h_1) = \sum_{i=1}^n [y_i - (h_0 + h_1 x_i)]^2$$

$$y = \hat{\beta}_0 + \hat{\beta}_1 x.$$

$$b_1 = \hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{S_{xy}}{S_{xx}}$$

$$b_0 = \hat{\beta}_0 = \frac{\sum y_i - \hat{\beta}_1 \sum x_i}{n} = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$S_{xy} = \sum x_i y_i - (\sum x_i)(\sum y_i)/n$$

$$S_{xx} = \sum x_i^2 - (\sum x_i)^2/n$$

Only within data range !

$$Y_1 - \hat{Y}_1, Y_2 - \hat{Y}_2, \dots, Y_n - \hat{Y}_n$$

Residuals

Estimate sigma

Error Sum of Square: variation unexplained by the model

$$SSE = \sum (y_i - \hat{y}_i)^2 = \sum [y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i)]^2$$

$$\hat{\sigma}^2 = s^2 = \frac{SSE}{n-2} = \frac{\sum (y_i - \hat{y}_i)^2}{n-2}$$

$$SSE = \sum y_i^2 - \hat{\beta}_0 \sum y_i - \hat{\beta}_1 \sum x_i y_i$$

$$\text{Total Sum of Square} \quad SST = S_{yy} = \sum (y_i - \bar{y})^2 = \sum y_i^2 - (\sum y_i)^2/n$$

The coefficient of determination,  $r^2$ : What proportion of variation could be explained by the model

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

$$SST = SSR + SSE$$

$$r^2 = 1 - SSE/SST = (SST - SSE)/SST = SSR/SST$$

SSR: Regression sum of square. variation explained by the model

Regression effect

### Estimator of Beta\_1

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(Y_i - \bar{Y})}{\sum (x_i - \bar{x})^2}$$

$$\hat{\beta}_0 = \frac{\sum Y_i - \hat{\beta}_1 \sum x_i}{n}$$

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})Y_i}{S_{xx}} = \sum c_i Y_i \quad \text{where } c_i = (x_i - \bar{x})/S_{xx}$$

$$\hat{\sigma}^2 = s^2 = \frac{\sum Y_i^2 - \hat{\beta}_0 \sum Y_i - \hat{\beta}_1 \sum x_i Y_i}{n-2}$$

$$V(\hat{\beta}_1) = \sigma_{\hat{\beta}_1}^2 = \frac{\sigma^2}{S_{xx}} \quad \sigma_{\hat{\beta}_1} = \frac{\sigma}{\sqrt{S_{xx}}}$$

$$S_{xx} = \sum (x_i - \bar{x})^2 = \sum x_i^2 - (\sum x_i)^2/n$$

$$s_{\hat{\beta}_1} = \frac{s}{\sqrt{S_{xx}}}$$

$$T = \frac{\hat{\beta}_1 - \beta_1}{s_{\hat{\beta}_1}} = \frac{\hat{\beta}_1 - \beta_1}{S_{\hat{\beta}_1}} \quad df = n-2$$

$$P(-t_{\alpha/2, n-2} < \frac{\hat{\beta}_1 - \beta_1}{s_{\hat{\beta}_1}} < t_{\alpha/2, n-2}) = 1 - \alpha$$

A  $100(1 - \alpha)\%$  CI for the slope  $\beta_1$  of the true regression line is

$$\hat{\beta}_1 \pm t_{\alpha/2, n-2} \cdot s_{\hat{\beta}_1}$$

CI for beta\_1

Inference about slop, beta\_1

Null hypothesis:  $H_0: \beta_1 = \beta_{10}$   
Test statistic value:  $t = \frac{\hat{\beta}_1 - \beta_{10}}{s_{\hat{\beta}_1}}$   
Alternative Hypothesis  
 $H_a: \beta_1 > \beta_{10}$   
 $H_a: \beta_1 < \beta_{10}$   
 $H_a: \beta_1 \neq \beta_{10}$   
Rejection Region for Level  $\alpha$  Test  
 $t \geq t_{\alpha, n-2}$   
 $t \leq -t_{\alpha, n-2}$   
either  $t \geq t_{\alpha/2, n-2}$  or  $t \leq -t_{\alpha/2, n-2}$   
A  $P$ -value based on  $n - 2$  df can be calculated just as was done previously for  $t$  tests in Chapters 8 and 9.  
The model utility test is the test of  $H_0: \beta_1 = 0$  versus  $H_a: \beta_1 \neq 0$ , in which case the test statistic value is the  $t$  ratio  $t = \hat{\beta}_1/s_{\hat{\beta}_1}$ .

Test procedure

$$t^2 = f \text{ and } t_{\alpha/2, n-2}^2 = F_{\alpha, 1, n-2}$$

Regression and ANOVA

Table 12.2 ANOVA Table for Simple Linear Regression

Source of Variation	df	Sum of Squares	Mean Square	$f$
Regression	1	SSR	SSR	SSR/SSE/(n-2)
Error	n-2	SSE	$s^2 = \frac{SSE}{n-2}$	
Total	n-1	SST		