

ISyE 6739 – Linear Regression (Chapters 11 & 12)

Instructor: Kamran Paynabar
H. Milton Stewart School of
Industrial and Systems Engineering
Georgia Tech

Kamran.paynabar@isye.gatech.edu
Office: Groseclose 436

Scatter Diagram

- Many problems in engineering and science involve exploring the relationships between two or more variables.
- **Regression analysis** is a statistical technique that is very useful for these types of problems.

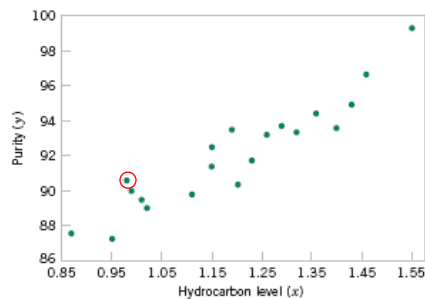


Table 11-1 Oxygen and Hydrocarbon Levels

Observation Number	Hydrocarbon Level x_i (%)	Purity y_i (%)
1	0.99	90.01
2	1.02	89.05
3	1.15	91.43
4	1.29	93.74
5	1.46	96.73
6	1.36	94.45
7	0.87	87.59
8	1.23	91.77
9	1.55	99.42
10	1.40	93.65
11	1.19	93.54
12	1.15	92.52
13	0.98	90.56
14	1.01	89.54
15	1.11	89.85
16	1.20	90.39
17	1.26	93.25
18	1.32	93.41
19	1.43	94.98
20	0.95	87.33

$$\hat{\rho} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \times \sum_{i=1}^n (y_i - \bar{y})^2}} \quad -1 \leq \hat{\rho} \leq 1$$

Hypothesis Test on Correlation

$H_0 : \rho = 0$ (population correlation ρ)

$H_0 : \rho \neq 0$

$$T_0 = \frac{\hat{\rho}\sqrt{n-2}}{\sqrt{1-\hat{\rho}^2}} \sim t_{n-2}; \quad (\text{sample correlation } r)$$

Cannot reject H_0 if $\left| \frac{\hat{\rho}\sqrt{n-2}}{\sqrt{1-\hat{\rho}^2}} \right| < t_{\alpha/2, n-2}$

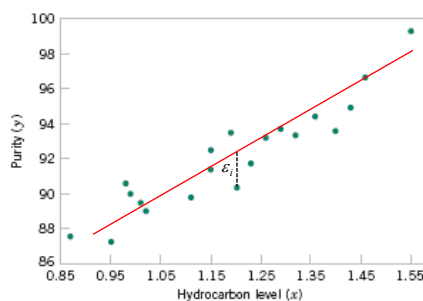
$$\left| \frac{\hat{\rho}\sqrt{n-2}}{\sqrt{1-\hat{\rho}^2}} \right| < t_{\alpha/2, n-2} \xrightarrow{\text{if } r \text{ is small}} \left| \hat{\rho}\sqrt{n-2} \right| < t_{\alpha/2, n-2} \xrightarrow{\text{if } n \text{ is large}} \left| \hat{\rho}\sqrt{n} \right| < z_{\alpha/2}$$

Approximate critical region for large n

$$\frac{2}{\sqrt{\text{number of observations, } n}} \quad \text{For } \alpha \text{ about } 0.05$$

Simple Linear Regression

Based on the scatter diagram, it is probably reasonable to assume that the mean of the random variable Y is related to X by the following **simple linear regression model**:



$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i \quad i = 1, 2, \dots, n$$

Response Regressor or Predictor

Intercept Slope Random error

$$\varepsilon_i \sim NID(0, \sigma^2)$$

where the slope and intercept of the line are called **regression coefficients**.

- The case of simple linear regression considers a single regressor or predictor x and a dependent or response variable Y.

Simple Linear Regression

The **method of least squares** is used to estimate the parameters, β_0 and β_1 by minimizing the sum of the squares of the vertical deviations in Figure 11-3.

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i \quad i = 1, 2, \dots, n$$

sum of the squares of the error

$$L = \sum_{i=1}^n \epsilon_i^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

Minimize

$$\left. \frac{\partial L}{\partial \beta_0} \right|_{\hat{\beta}_0, \hat{\beta}_1} = -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0$$

$$\left. \frac{\partial L}{\partial \beta_1} \right|_{\hat{\beta}_0, \hat{\beta}_1} = -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i = 0$$

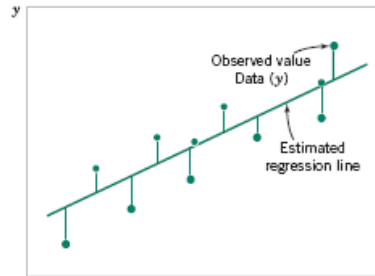


Figure 11-3 Deviations of the data from the estimated regression model.

Least Square Normal Equations

Least Square Estimates

The least squares estimates of the intercept and slope in the simple linear regression model are

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad (11-7)$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n y_i x_i - \frac{\left(\sum_{i=1}^n y_i\right)\left(\sum_{i=1}^n x_i\right)}{n}}{\sum_{i=1}^n x_i^2 - \frac{\left(\sum_{i=1}^n x_i\right)^2}{n}} \quad (11-8)$$

where $\bar{y} = (1/n) \sum_{i=1}^n y_i$ and $\bar{x} = (1/n) \sum_{i=1}^n x_i$.

Alternative Notation

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - \frac{\left(\sum_{i=1}^n x_i\right)^2}{n} \quad S_{xy} = \sum_{i=1}^n y_i (x_i - \bar{x}) = \sum_{i=1}^n x_i y_i - \frac{\left(\sum_{i=1}^n x_i\right)\left(\sum_{i=1}^n y_i\right)}{n}$$

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

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$$

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

Fitted (estimated) regression model

Example

Find the least square estimates of the simple linear regression describing the relationship between Purity (y) and Hydrocarbon Levels (x).

Also, calculate the predicted purity when hydrocarbon level is 1.01. Find the prediction error.

$$n = 20 \quad \sum_{i=1}^{20} x_i = 23.92 \quad \sum_{i=1}^{20} y_i = 1,843.21 \quad \bar{x} = 1.1960 \quad \bar{y} = 92.1605$$

$$\sum_{i=1}^{20} y_i^2 = 170,044.5321 \quad \sum_{i=1}^{20} x_i^2 = 29.2892 \quad \sum_{i=1}^{20} x_i y_i = 2,214.6566$$

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - \frac{\left(\sum_{i=1}^n x_i\right)^2}{n}$$

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

Table 11-1 Oxygen and Hydrocarbon Levels

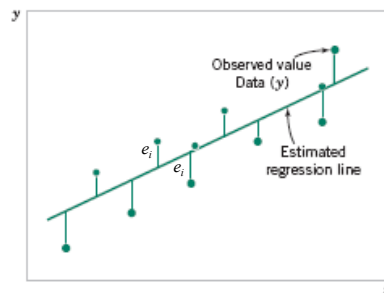
Observation Number	Hydrocarbon Level x(%)	Purity y(%)
1	0.99	90.01
2	1.02	89.05
3	1.15	91.43
4	1.29	93.74
5	1.46	96.73
6	1.36	94.45
7	0.87	87.59
8	1.23	91.77
9	1.55	99.42
10	1.40	93.65
11	1.19	93.54
12	1.15	92.52
13	0.98	90.56
14	1.01	89.54
15	1.11	89.85
16	1.20	90.39
17	1.26	93.25
18	1.32	93.41
19	1.43	94.98
20	0.95	87.33

Estimation of Variance (σ^2)

The error sum of squares is

$$SS_E = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$SS_E = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$$



An unbiased estimator of σ^2 is

$$\hat{\sigma}^2 = \frac{SS_E}{n - 2} \quad (11-13)$$

where SS_E can be easily computed using (easier formula)

$$SS_E = SS_T - \hat{\beta}_1 S_{xy} \quad SS_T = \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n y_i^2 - \frac{\left(\sum_{i=1}^n y_i\right)^2}{n}$$

Confidence Intervals for Coefficients

Slope Properties

$$E(\hat{\beta}_1) = \beta_1 \quad V(\hat{\beta}_1) = \frac{\sigma^2}{S_{xx}}$$

Intercept Properties

$$E(\hat{\beta}_0) = \beta_0 \quad \text{and} \quad V(\hat{\beta}_0) = \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right]$$

Under the assumption that the observations are normally and independently distributed, a $100(1 - \alpha)\%$ **confidence interval on the slope** β_1 in simple linear regression is

$$\hat{\beta}_1 - t_{\alpha/2, n-2} \sqrt{\frac{\hat{\sigma}^2}{S_{xx}}} \leq \beta_1 \leq \hat{\beta}_1 + t_{\alpha/2, n-2} \sqrt{\frac{\hat{\sigma}^2}{S_{xx}}} \quad (11-29)$$

Similarly, a $100(1 - \alpha)\%$ **confidence interval on the intercept** β_0 is

$$\begin{aligned} \hat{\beta}_0 - t_{\alpha/2, n-2} \sqrt{\hat{\sigma}^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right]} \\ \leq \beta_0 \leq \hat{\beta}_0 + t_{\alpha/2, n-2} \sqrt{\hat{\sigma}^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right]} \end{aligned} \quad (11-30)$$

9

Example

EXAMPLE 11-4 Oxygen Purity Confidence Interval on the Slope

We will find a 95% confidence interval on the slope of the regression line using the data in Example 11-1. Recall that $\hat{\beta}_1 = 14.947$, $S_{xx} = 0.68088$, and $\hat{\sigma}^2 = 1.18$ (see Table 11-2). Then, from Equation 11-29 we find

This simplifies to

$$12.181 \leq \beta_1 \leq 17.713$$

$$\hat{\beta}_1 - t_{0.025, 18} \sqrt{\frac{\hat{\sigma}^2}{S_{xx}}} \leq \beta_1 \leq \hat{\beta}_1 + t_{0.025, 18} \sqrt{\frac{\hat{\sigma}^2}{S_{xx}}}$$

or

$$\begin{aligned} 14.947 - 2.101 \sqrt{\frac{1.18}{0.68088}} \leq \beta_1 \leq 14.947 \\ + 2.101 \sqrt{\frac{1.18}{0.68088}} \end{aligned}$$

Practical Interpretation: This CI does not include zero, so there is strong evidence (at $\alpha = 0.05$) that the slope is not zero. The CI is reasonably narrow (± 2.766) because the error variance is fairly small.

10

Confidence Intervals for Coefficients

Slope Properties

$$E(\hat{\beta}_1) = \beta_1 \quad V(\hat{\beta}_1) = \frac{\sigma^2}{S_{xx}}$$

Intercept Properties

$$E(\hat{\beta}_0) = \beta_0 \quad \text{and} \quad V(\hat{\beta}_0) = \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right]$$

Under the assumption that the observations are normally and independently distributed, a $100(1 - \alpha)\%$ **confidence interval on the slope** β_1 in simple linear regression is

$$\hat{\beta}_1 - t_{\alpha/2, n-2} \sqrt{\frac{\hat{\sigma}^2}{S_{xx}}} \leq \beta_1 \leq \hat{\beta}_1 + t_{\alpha/2, n-2} \sqrt{\frac{\hat{\sigma}^2}{S_{xx}}} \quad (11-29)$$

Similarly, a $100(1 - \alpha)\%$ **confidence interval on the intercept** β_0 is

$$\begin{aligned} \hat{\beta}_0 - t_{\alpha/2, n-2} \sqrt{\hat{\sigma}^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right]} \\ \leq \beta_0 \leq \hat{\beta}_0 + t_{\alpha/2, n-2} \sqrt{\hat{\sigma}^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right]} \end{aligned} \quad (11-30)$$

11

Hypothesis Tests in Simple Linear Regression

Suppose we wish to test $H_0: \beta_1 = \beta_{1,0}$

$$H_1: \beta_1 \neq \beta_{1,0}$$

An appropriate test statistic would be

$$T_0 = \frac{\hat{\beta}_1 - \beta_{1,0}}{se(\hat{\beta}_1)}$$

We would reject the null hypothesis if

$$|t_0| > t_{\alpha/2, n-2}$$

12

Hypothesis Tests in Simple Linear Regression

Suppose we wish to test

$$H_0: \beta_0 = \beta_{0,0}$$

$$H_1: \beta_0 \neq \beta_{0,0}$$

An appropriate test statistic would be

$$T_0 = \frac{\hat{\beta}_0 - \beta_{0,0}}{\sqrt{\hat{\sigma}^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right]}} = \frac{\hat{\beta}_0 - \beta_{0,0}}{se(\hat{\beta}_0)}$$

We would reject the null hypothesis if

$$|t_0| > t_{\alpha/2, n-2}$$

13

Example

EXAMPLE 11-2 Oxygen Purity Tests of Coefficients

We will test for significance of regression using the model for the oxygen purity data from Example 11-1. The hypotheses are

$$H_0: \beta_1 = 0$$

$$H_1: \beta_1 \neq 0$$

and we will use $\alpha = 0.01$. From Example 11-1 and Table 11-2 we have

$$\hat{\beta}_1 = 14.947 \quad n = 20, \quad S_{xx} = 0.68088, \quad \hat{\sigma}^2 = 1.18$$

so the t -statistic in Equation 10-20 becomes

$$t_0 = \frac{\hat{\beta}_1}{\sqrt{\hat{\sigma}^2/S_{xx}}} = \frac{\hat{\beta}_1}{se(\hat{\beta}_1)} = \frac{14.947}{\sqrt{1.18/0.68088}} = 11.35$$

Practical Interpretation: Since the reference value of t is $t_{0.005, 18} = 2.88$, the value of the test statistic is very far into the critical region, implying that $H_0: \beta_1 = 0$ should be rejected. There is strong evidence to support this claim. The P -value for this test is $P \approx 1.23 \times 10^{-9}$. This was obtained manually with a calculator.

Table 11-2 presents the Minitab output for this problem. Notice that the t -statistic value for the slope is computed as 11.35 and that the reported P -value is $P = 0.000$. Minitab also reports the t -statistic for testing the hypothesis $H_0: \beta_0 = 0$. This statistic is computed from Equation 11-22, with $\beta_{0,0} = 0$, as $t_0 = 46.62$. Clearly, then, the hypothesis that the intercept is zero is rejected.

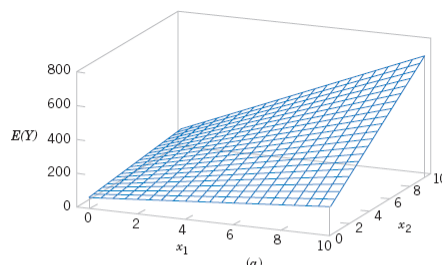
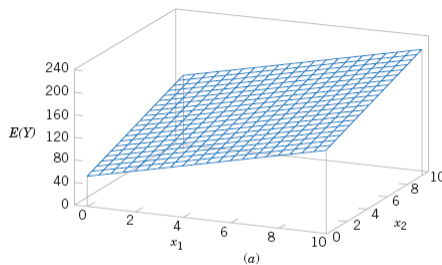
14

Multiple Linear Regression

A regression model that contains more than one regressor variable is called a **multiple regression model**.

For example, suppose that the effective life of a cutting tool depends on the cutting speed and the tool angle. A possible multiple regression model could be

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon \quad Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_{12} X_1 X_2 + \epsilon$$



ISyE 6739, Regression

15

Multiple Linear Regression – Least Square Estimates

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1k} \\ 1 & x_{21} & x_{22} & \dots & x_{2k} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nk} \end{bmatrix} \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix} \quad \text{and} \quad \boldsymbol{\epsilon} = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

We wish to find the vector of least squares estimators that minimizes:

$$L = \sum_{i=1}^n \epsilon_i^2 = \boldsymbol{\epsilon}'\boldsymbol{\epsilon} = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

The resulting least squares estimate is

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y} \quad (12-13)$$

ISyE 6739, Regression

16

Example

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$$

where y is the observed pull strength for a wire bond, x_1 is the wire length, and x_2 is the die height.

$X =$	$\begin{bmatrix} 1 & 2 & 50 \\ 1 & 8 & 110 \\ 1 & 11 & 120 \\ 1 & 10 & 550 \\ 1 & 8 & 295 \\ 1 & 4 & 200 \\ 1 & 2 & 375 \\ 1 & 2 & 52 \\ 1 & 9 & 100 \\ 1 & 8 & 300 \\ 1 & 4 & 412 \\ 1 & 11 & 400 \\ 1 & 12 & 500 \\ 1 & 2 & 360 \\ 1 & 4 & 205 \\ 1 & 4 & 400 \\ 1 & 20 & 600 \\ 1 & 1 & 585 \\ 1 & 10 & 540 \\ 1 & 15 & 250 \\ 1 & 15 & 290 \\ 1 & 16 & 510 \\ 1 & 17 & 590 \\ 1 & 6 & 100 \\ 1 & 5 & 400 \end{bmatrix}$	$y =$	$\begin{bmatrix} 9.95 \\ 24.45 \\ 31.75 \\ 35.00 \\ 25.02 \\ 16.86 \\ 14.38 \\ 9.60 \\ 24.35 \\ 27.50 \\ 17.08 \\ 37.00 \\ 41.95 \\ 11.66 \\ 21.65 \\ 17.89 \\ 69.00 \\ 10.30 \\ 34.93 \\ 46.59 \\ 44.88 \\ 54.12 \\ 56.63 \\ 22.13 \\ 21.15 \end{bmatrix}$
-------	---	-------	---

ISyE 6739, Regression

Example (continued)

The $X'X$ matrix is

$$X'X = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ 2 & 8 & \cdots & 5 \\ 50 & 110 & \cdots & 400 \end{bmatrix} \begin{bmatrix} 1 & 2 & 50 \\ 1 & 8 & 110 \\ \vdots & \vdots & \vdots \\ 1 & 5 & 400 \end{bmatrix}$$

$$= \begin{bmatrix} 25 & 206 & 8,294 \\ 206 & 2,396 & 77,177 \\ 8,294 & 77,177 & 3,531,848 \end{bmatrix}$$

and the $X'y$ vector is

$$X'y = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ 2 & 8 & \cdots & 5 \\ 50 & 110 & \cdots & 400 \end{bmatrix} \begin{bmatrix} 9.95 \\ 24.45 \\ \vdots \\ 21.15 \end{bmatrix} = \begin{bmatrix} 725.82 \\ 8,008.47 \\ 274,816.71 \end{bmatrix}$$

The least squares estimates are found from Equation 12-13 as

$$\hat{\beta} = (X'X)^{-1}X'y$$

or

$$\begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix} = \begin{bmatrix} 25 & 206 & 8,294 \\ 206 & 2,396 & 77,177 \\ 8,294 & 77,177 & 3,531,848 \end{bmatrix}^{-1} \begin{bmatrix} 725.82 \\ 8,008.47 \\ 274,811.31 \end{bmatrix}$$

$$= \begin{bmatrix} 0.214653 & -0.007491 & -0.000340 \\ -0.007491 & 0.001671 & -0.000019 \\ -0.000340 & -0.000019 & +0.0000015 \end{bmatrix} \begin{bmatrix} 725.82 \\ 8,008.47 \\ 274,811.31 \end{bmatrix}$$

$$= \begin{bmatrix} 2.26379143 \\ 2.74426964 \\ 0.01252781 \end{bmatrix}$$

Therefore, the fitted regression model with the regression coefficients rounded to five decimal places is

$$\hat{y} = 2.26379 + 2.74427x_1 + 0.01253x_2$$

This is identical to the results obtained in Example 12-1.

ISyE 6739, Regression

Example (continued)

Table 12-3 Observations, Fitted Values, and Residuals for Example 12-2

Observation Number	y_i	\hat{y}_i	$e_i = y_i - \hat{y}_i$	Observation Number	y_i	\hat{y}_i	$e_i = y_i - \hat{y}_i$
1	9.95	8.38	1.57	14	11.66	12.26	-0.60
2	24.45	25.60	-1.15	15	21.65	15.81	5.84
3	31.75	33.95	-2.20	16	17.89	18.25	-0.36
4	35.00	36.60	-1.60	17	69.00	64.67	4.33
5	25.02	27.91	-2.89	18	10.30	12.34	-2.04
6	16.86	15.75	1.11	19	34.93	36.47	-1.54
7	14.38	12.45	1.93	20	46.59	46.56	0.03
8	9.60	8.40	1.20	21	44.88	47.06	-2.18
9	24.35	28.21	-3.86	22	54.12	52.56	1.56
10	27.50	27.98	-0.48	23	56.63	56.31	0.32
11	17.08	18.40	-1.32	24	22.13	19.98	2.15
12	37.00	37.46	-0.46	25	21.15	21.00	0.15
13	41.95	41.46	0.49				

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n e_i^2}{n-p} = \frac{SS_E}{n-p} \quad (12-16)$$

Covariance Matrix Estimation

Unbiased estimators:
$$\begin{aligned} E(\hat{\beta}) &= E[(X'X)^{-1}X'Y] \\ &= E[(X'X)^{-1}X'(X\beta + \epsilon)] \\ &= E[(X'X)^{-1}X'X\beta + (X'X)^{-1}X'\epsilon] \\ &= \beta \end{aligned}$$

Covariance Matrix:
$$\text{cov}(\hat{\beta}) = \sigma^2(X'X)^{-1} = \sigma^2 C$$

$$\begin{aligned} V(\hat{\beta}_j) &= \sigma^2 C_{jj}, & j = 0, 1, 2 \\ \text{cov}(\hat{\beta}_i, \hat{\beta}_j) &= \sigma^2 C_{ij}, & i \neq j \end{aligned} \quad C = (X'X)^{-1} = \begin{bmatrix} C_{00} & C_{01} & C_{02} \\ C_{10} & C_{11} & C_{12} \\ C_{20} & C_{21} & C_{22} \end{bmatrix}$$

Confidence Interval for Regression Coefficients

Mean and variance of the slope estimator

$$V(\hat{\beta}_j) = \sigma^2 C_{jj}, \quad E(\hat{\beta}) = \beta$$

A $100(1 - \alpha) \%$ confidence interval on the regression coefficient β_j , $j = 0, 1, \dots, k$ in the multiple linear regression model is given by

$$\hat{\beta}_j - t_{\alpha/2, n-p} \sqrt{\hat{\sigma}^2 C_{jj}} \leq \beta_j \leq \hat{\beta}_j + t_{\alpha/2, n-p} \sqrt{\hat{\sigma}^2 C_{jj}} \quad (12-35)$$

$$\hat{\sigma}^2 = \frac{SS_E}{n - p}$$

$$SS_E = SS_T - SS_R = \mathbf{y}'\mathbf{y} - \hat{\beta}'\mathbf{X}'\mathbf{y}$$

Hypothesis Tests on Regression Coefficients

$$H_0: \beta_j = \beta_{j0}$$

$$H_1: \beta_j \neq \beta_{j0}$$

(12-24)

The test statistic is

$$T_0 = \frac{\hat{\beta}_j - \beta_{j0}}{\sqrt{\hat{\sigma}^2 C_{jj}}} = \frac{\hat{\beta}_j - \beta_{j0}}{se(\hat{\beta}_j)} \quad (12-25)$$

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

Example

EXAMPLE 12-4 Wire Bond Strength Coefficient Test

Consider the wire bond pull strength data, and suppose that we want to test the hypothesis that the regression coefficient for x_2 (die height) is zero. The hypotheses are

$$H_0: \beta_2 = 0$$

$$H_1: \beta_2 \neq 0$$

The main diagonal element of the $(X'X)^{-1}$ matrix corresponding to β_2 is $C_{22} = 0.0000015$, so the t -statistic in Equation 12-25 is

$$t_0 = \frac{\hat{\beta}_2}{\sqrt{\hat{\sigma}^2 C_{22}}} = \frac{0.01253}{\sqrt{(5.2352)(0.0000015)}} = 4.477$$

We will construct a 95% confidence interval on the parameter β_1 in the wire bond pull strength problem. The point estimate of β_1 is $\hat{\beta}_1 = 2.74427$, and the diagonal element of $(X'X)^{-1}$ corresponding to β_1 is $C_{11} = 0.001671$. The estimate of σ^2 is $\hat{\sigma}^2 = 5.2352$, and $t_{0.025,22} = 2.074$. Therefore, the 95% CI on β_1 is computed from Equation 12-35 as

$$2.74427 - (2.074)\sqrt{(5.2352)(.001671)} \leq \beta_1 \leq 2.74427 + (2.074)\sqrt{(5.2352)(.001671)}$$

which reduces to

$$2.55029 \leq \beta_1 \leq 2.93825$$

Hypothesis Tests On Multiple Coefficients

$$H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$$

$$H_1: \beta_j \neq 0 \quad \text{for at least one } j \quad (12-18)$$

SS of Total = SS of Regression + SS of Error

The test statistic is based on ANOVA

$$F_0 = \frac{SS_R/k}{SS_E/(n-p)} = \frac{MS_R}{MS_E} \quad (12-19)$$

Table 12-9 Analysis of Variance for Testing Significance of Regression in Multiple Regression

Source of Variation	Sum of Squares	Degrees of Freedom	Mean Square	F_0
Regression	SS_R	k	MS_R	MS_R/MS_E
Error or residual	SS_E	$n - p$	MS_E	
Total	SS_T	$n - 1$		

Example

EXAMPLE 12-3 Wire Bond Strength ANOVA

We will test for significance of regression (with $\alpha = 0.05$) using the wire bond pull strength data from Example 12-1. The total sum of squares is

$$SS_T = y'y - \frac{\left(\sum_{i=1}^n y_i\right)^2}{n} = 27,178.5316 - \frac{(725.82)^2}{25} = 6105.9447$$

The regression or model sum of squares is computed from Equation 12-20 as follows:

$$SS_R = \hat{\beta}'X'y - \frac{\left(\sum_{i=1}^n y_i\right)^2}{n} = 27,063.3581 - \frac{(725.82)^2}{25} = 5990.7712$$

and by subtraction

$$SS_E = SS_T - SS_R = y'y - \hat{\beta}'X'y = 115.1716$$

Table 12-10 Test for Significance of Regression for Example 12-3

Source of Variation	Sum of Squares	Degrees of Freedom	Mean Square	f_0	P -value
Regression	5990.7712	2	2995.3856	572.17	1.08E-19
Error or residual	115.1735	22	5.2352		
Total	6105.9447	24			

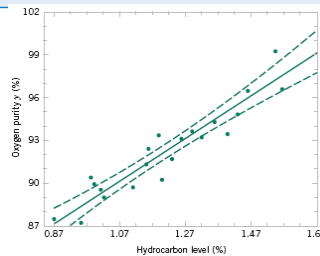
Confidence Interval on Mean Response

The mean response at a point x_0 is estimated by $\hat{\mu}_{Y|x_0} = x_0'\hat{\beta}$

The variance of the estimated mean response is $V(\hat{\mu}_{Y|x_0}) = \sigma^2 x_0'(X'X)^{-1}x_0$

For the multiple linear regression model, a $100(1 - \alpha)\%$ confidence interval on the mean response at the point $x_{01}, x_{02}, \dots, x_{0k}$ is

$$\hat{\mu}_{Y|x_0} - t_{\alpha/2, n-p} \sqrt{\hat{\sigma}^2 x_0'(X'X)^{-1}x_0} \leq \mu_{Y|x_0} \leq \hat{\mu}_{Y|x_0} + t_{\alpha/2, n-p} \sqrt{\hat{\sigma}^2 x_0'(X'X)^{-1}x_0} \quad (12-39)$$



Example

EXAMPLE 12-8 Wire Bond Strength Confidence Interval on the Mean Response

The engineer in Example 12-1 would like to construct a 95% CI on the mean pull strength for a wire bond with wire length $x_1 = 8$ and die height $x_2 = 275$. Therefore,

$$\mathbf{x}_0 = \begin{bmatrix} 1 \\ 8 \\ 275 \end{bmatrix}$$

The estimated mean response at this point is found from Equation 12-36 as

$$\hat{\mu}_{Y|\mathbf{x}_0} = \mathbf{x}_0' \hat{\boldsymbol{\beta}} = \begin{bmatrix} 1 & 8 & 275 \end{bmatrix} \begin{bmatrix} 2.26379 \\ 2.74427 \\ 0.01253 \end{bmatrix} = 27.66$$

The variance of $\hat{\mu}_{Y|\mathbf{x}_0}$ is estimated by

$$\begin{aligned} \hat{\sigma}^2 \mathbf{x}_0' (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_0 &= 5.2352 \begin{bmatrix} 1 & 8 & 275 \end{bmatrix} \\ &\times \begin{bmatrix} .214653 & -.007491 & -.000340 \\ -.007491 & .001671 & -.000019 \\ -.000340 & -.000019 & .0000015 \end{bmatrix} \begin{bmatrix} 1 \\ 8 \\ 275 \end{bmatrix} \\ &= 5.2352 (0.0444) = 0.23244 \end{aligned}$$

Therefore, a 95% CI on the mean pull strength at this point is found from Equation 12-39 as

$$27.66 - 2.074 \sqrt{0.23244} \leq \mu_{Y|\mathbf{x}_0} \leq 27.66 + 2.074 \sqrt{0.23244}$$

which reduces to

$$26.66 \leq \mu_{Y|\mathbf{x}_0} \leq 28.66$$

Prediction Interval for New Observations

A point estimate of the future observation Y_0 is $\hat{y}_0 = \mathbf{x}_0' \hat{\boldsymbol{\beta}}$

A $100(1 - \alpha)\%$ prediction interval for this future observation is

$$\begin{aligned} \hat{y}_0 - t_{\alpha/2, n-p} \sqrt{\hat{\sigma}^2 (1 + \mathbf{x}_0' (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_0)} \\ \leq Y_0 \leq \hat{y}_0 + t_{\alpha/2, n-p} \sqrt{\hat{\sigma}^2 (1 + \mathbf{x}_0' (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_0)} \end{aligned} \quad (12-41)$$

Adequacy of Regression Model

Simple linear regression assumptions:

1. Errors are uncorrelated random variables with mean zero;
2. Errors have constant variance; and,
3. Errors be normally distributed. $\varepsilon_i \sim NID(0, \sigma^2)$

- The analyst should always consider the validity of these assumptions to be doubtful and conduct analyses to examine the adequacy of the model

• The residuals from a regression model are $e_i = y_i - \hat{y}_i$, where y_i is an actual observation and \hat{y}_i is the corresponding fitted value from the regression model.

• Analysis of the residuals is frequently helpful in checking the assumption that the errors are approximately normally distributed with constant variance, and in determining whether additional terms in the model would be useful.

Adequacy of Regression Model

Analysis of Residual Patterns is useful for checking:

- Independency assumption
- Constant variance

Plot residuals (e_i) against predicted response (\hat{y}_i)

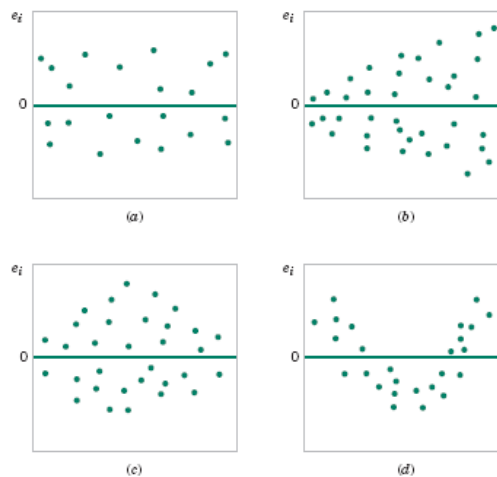
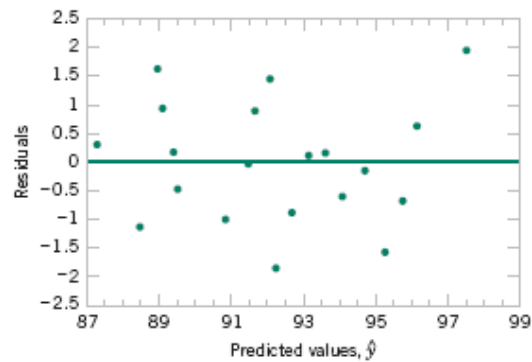


Figure 11-9 Patterns for residual plots. (a) satisfactory, (b) funnel, (c) double bow, (d) nonlinear. [Adapted from Montgomery, Peck, and Vining (2001).]

Adequacy of Regression Model

Figure 11-11 Plot of residuals versus predicted oxygen purity, \hat{y} , Example 11-7.



Adequacy of Regression Model

Histogram and Normal plot for residuals:

- Normality assumption

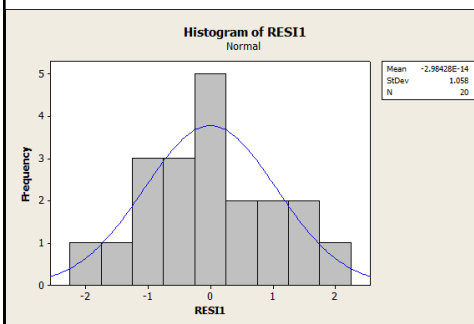
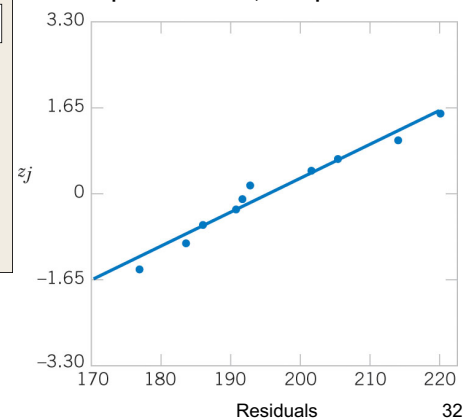


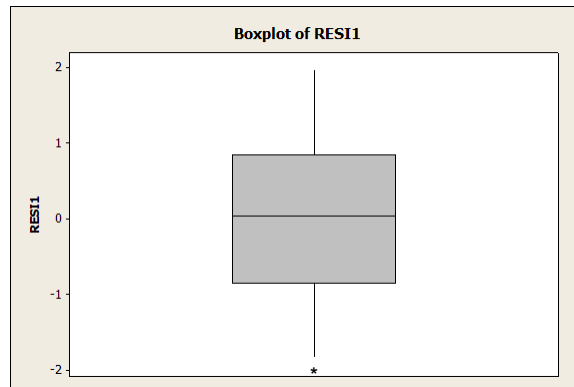
Figure 11-10 Normal probability plot of residuals, Example 11-7.



Adequacy of Regression Model

Boxplots:

- It is used to detect observations with large residuals (Outliers)



Adequacy of Regression Model

Coefficient of Determination (R^2)

R^2 is called the **coefficient of determination** and is often used to judge the adequacy of a regression model.

$$0 < R^2 < 1;$$

- We often refer (loosely) to R^2 as the amount of variability in the data explained or accounted for by the regression model.

$$R^2 = 1 - \frac{SS_E}{SS_T}$$

Adjusted Coefficient of Determination (R^2)

$$R^2_{\text{adj}} = 1 - \frac{SS_E/(n-p)}{SS_T/(n-1)} \quad (12-23)$$