

LINEAR REGRESSION ANALYSIS

MODULE – II

Lecture - 5

Simple Linear Regression Analysis

Dr. Shalabh

Department of Mathematics and Statistics

Indian Institute of Technology Kanpur

Joint confidence region for β_0 and β_1

A joint confidence region for β_0 and β_1 can also be found. Such region will provide a $100(1-\alpha)\%$ confidence that both the estimates of β_0 and β_1 are correct. Consider the centered version of the linear regression model

$$y_i = \beta_0^* + \beta_1(x_i - \bar{x}) + \varepsilon_i$$

where $\beta_0^* = \beta_0 + \beta_1\bar{x}$. The least squares estimators of β_0^* and β_1 are $b_0^* = \bar{y}$ and $b_1 = \frac{s_{xy}}{s_{xx}}$ respectively.

Using the results that $E(b_0^*) = \beta_0^*$,

$$E(b_1) = \beta_1,$$

$$Var(b_0^*) = \frac{\sigma^2}{n},$$

$$Var(b_1) = \frac{\sigma^2}{s_{xx}}.$$

When σ^2 is known, then the statistic

$$\frac{b_0^* - \beta_0^*}{\sqrt{\frac{\sigma^2}{n}}} \sim N(0,1)$$

and

$$\frac{b_1 - \beta_1}{\sqrt{\frac{\sigma^2}{s_{xx}}}} \sim N(0,1).$$

Moreover, both the statistics are independently distributed. Thus

$$\left(\frac{b_0^* - \beta_0^*}{\sqrt{\frac{\sigma^2}{n}}} \right)^2 \sim \chi_1^2$$

and

$$\left(\frac{b_1 - \beta_1}{\sqrt{\frac{\sigma^2}{s_{xx}}}} \right)^2 \sim \chi_1^2$$

are also independently distributed because b_0^* and b_1 are independently distributed. Consequently sum of these two

$$\frac{n(b_0^* - \beta_0^*)^2}{\sigma^2} + \frac{s_{xx}(b_1 - \beta_1)^2}{\sigma^2} \sim \chi_2^2.$$

Since

$$\frac{SS_{res}}{\sigma^2} \sim \chi_{n-2}^2$$

and SS_{res} is independently distributed of b_0^* and b_1 , so the ratio

$$\frac{\left(\frac{n(b_0^* - \beta_0^*)^2}{\sigma^2} + \frac{s_{xx}(b_1 - \beta_1)^2}{\sigma^2} \right) / 2}{\left(\frac{SS_{res}}{\sigma^2} \right) / (n-2)} \sim F_{2, n-2}.$$

Substituting $b_0^* = b_0 + b_1\bar{x}$ and $\beta_0^* = \beta_0 + \beta_1\bar{x}$, we get

$$\left(\frac{n-2}{2}\right)\left[\frac{Q_f}{SS_{res}}\right]$$

where

$$Q_f = n(b_0 - \beta_0)^2 + 2\sum_{i=1}^n x_i(b_0 - \beta_0)(b_1 - \beta_1) + \sum_{i=1}^n x_i^2(b_1 - \beta_1)^2.$$

Since

$$P\left[\left(\frac{n-2}{2}\right)\frac{Q_f}{SS_{res}} \leq F_{2,n-2}\right] = 1 - \alpha$$

holds true for all values of β_0 and β_1 , so the $100(1-\alpha)\%$ confidence region for β_0 and β_1 is

$$\left(\frac{n-2}{2}\right)\frac{Q_f}{SS_{res}} \leq F_{2,n-2;\alpha}.$$

This confidence region is an ellipse which gives the $100(1-\alpha)\%$ probability that β_0 and β_1 are contained simultaneously in this ellipse.

Analysis of variance

The technique of analysis of variance is usually used for testing the hypothesis related to equality of more than one parameters, like population means or slope parameters. It is more meaningful in case of multiple regression model when there are more than one slope parameters. This technique is discussed and illustrated here to understand the related basic concepts and fundamentals which will be used in developing the analysis of variance in the next module in multiple linear regression model where the explanatory variables are more than two.

A test statistic for testing $H_0 : \beta_1 = 0$ can also be formulated using the analysis of variance technique as follows.

On the basis of the identity $y_i - \hat{y}_i = (y_i - \bar{y}) - (\hat{y}_i - \bar{y})$,

the sum of squared residuals is

$$\begin{aligned} S(b) &= \sum_{i=1}^n (y_i - \hat{y}_i)^2 \\ &= \sum_{i=1}^n (y_i - \bar{y})^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 - 2 \sum_{i=1}^n (y_i - \bar{y})(\hat{y}_i - \bar{y}). \end{aligned}$$

Further consider

$$\begin{aligned} \sum_{i=1}^n (y_i - \bar{y})(\hat{y}_i - \bar{y}) &= \sum_{i=1}^n (y_i - \bar{y})b_1(x_i - \bar{x}) \\ &= b_1^2 \sum_{i=1}^n (x_i - \bar{x})^2 \\ &= \sum_{i=1}^n (\hat{y}_i - \bar{y})^2. \end{aligned}$$

Thus we have

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2.$$

The term $\sum_{i=1}^n (y_i - \bar{y})^2$ is called the **sum of squares about the mean** or **corrected sum of squares** of y (i.e., SS corrected) or total sum of squares denoted as s_{yy} .

The term $\sum_{i=1}^n (y_i - \hat{y}_i)^2$ describes the deviation: observation minus predicted value, viz., the residual sum of squares, i.e.:

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

whereas the term $\sum_{i=1}^n (\hat{y}_i - \bar{y})^2$ describes the proportion of variability explained by regression,

$$SS_{reg} = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2.$$

If all observations y_i are located on a straight line, then in this case $\sum_{i=1}^n (y_i - \hat{y}_i)^2 = 0$ and thus $SS_{corrected} = SS_{reg}$.

Note that SS_{reg} is completely determined by b_1 and so has only one degrees of freedom. The total sum of squares

$s_{yy} = \sum_{i=1}^n (y_i - \bar{y})^2$ has $(n - 1)$ degrees of freedom due to constraint $\sum_{i=1}^n (y_i - \bar{y}) = 0$ and SS_{res} has $(n - 2)$ degrees of freedom as it depends on b_0 and b_1 .

All sums of squares are mutually independent and distributed as χ_{df}^2 with df degrees of freedom if the errors are normally distributed.

The mean square due to regression is

$$MS_{reg} = \frac{SS_{reg}}{1}$$

and mean square due to residuals is

$$MSE = \frac{SS_{res}}{n - 2}.$$

The test statistic for testing $H_0 : \beta_1 = 0$ is

$$F_0 = \frac{MS_{reg}}{MSE}.$$

If $H_0 : \beta_1 = 0$ is true, then MS_{reg} and MSE are independently distributed and thus $F_0 \sim F_{1,n-2}$.

The decision rule for $H_1 : \beta_1 \neq 0$ is to reject H_0 if $F_0 > F_{1,n-2;1-\alpha}$

at α level of significance. The test procedure can be described in an Analysis of Variance table.

Analysis of variance for testing $H_0 : \beta_1 = 0$

Source of variation	Sum of squares	Degrees of freedom	Mean Square
Regression	SS_{reg}	1	MS_{reg}
Residual	SS_{res}	$n - 2$	MSE
Total	s_{yy}	$n - 1$	

Some other forms of SS_{reg} , SS_{res} and s_{yy} can be derived as follows:

The sample correlation coefficient then may be written as

$$r_{xy} = \frac{s_{xy}}{\sqrt{s_{xx}} \sqrt{s_{yy}}}.$$

Moreover, we have $b_1 = \frac{s_{xy}}{s_{xx}} = r_{xy} \sqrt{\frac{s_{yy}}{s_{xx}}}$.

The estimator of σ^2 in this case may be expressed as

$$\begin{aligned} s^2 &= \frac{1}{n-2} \sum_{i=1}^n e_i^2 \\ &= \frac{1}{n-2} SS_{res}. \end{aligned}$$

Various alternative formulations for SS_{res} are in use as well:

$$\begin{aligned}
 SS_{res} &= \sum_{i=1}^n [y_i - (b_0 + b_1 x_i)]^2 \\
 &= \sum_{i=1}^n [(y_i - \bar{y}) - b_1 (x_i - \bar{x})]^2 \\
 &= s_{yy} + b_1^2 s_{xx} - 2b_1 s_{xy} \\
 &= s_{yy} - b_1^2 s_{xx} \\
 &= s_{yy} - \frac{(s_{xy})^2}{s_{xx}}.
 \end{aligned}$$

Using this result, we find that

$$SS_{corrected} = s_{yy}$$

and

$$\begin{aligned}
 SS_{reg} &= s_{yy} - SS_{res} \\
 &= \frac{(s_{xy})^2}{s_{xx}} \\
 &= b_1^2 s_{xx} \\
 &= b_1 s_{xy}.
 \end{aligned}$$

Goodness of fit of regression

It can be noted that a fitted model can be said to be good when residuals are small. Since SS_{res} is based on residuals, so a measure of quality of fitted model can be based on SS_{res} . When intercept term is present in the model, a measure of goodness of fit of the model is given by

$$R^2 = 1 - \frac{SS_{res}}{s_{yy}}$$

$$= \frac{SS_{reg}}{s_{yy}}.$$

This is known as the **coefficient of determination**. This measure is based on the concept that how much variation in y 's stated by s_{yy} is explainable by SS_{reg} and how much unexplainable part is contained in SS_{res} . The ratio SS_{reg} / s_{yy} describes the proportion of variability that is explained by regression in relation to the total variability of y . The ratio SS_{res} / s_{yy} describes the proportion of variability that is not covered by the regression.

It can be seen that

$$R^2 = r_{xy}^2.$$

where r_{xy} is the simple correlation coefficient between x and y . Clearly $0 \leq R^2 \leq 1$, so a value of R^2 closer to one indicates the better fit and value of R^2 closer to zero indicates the poor fit.