

Multiple Linear Regression

When doing a linear regression we attempt to find out if the expected value of some variable Y , which we will call dependent variable, can be modeled as a linear function of a series of other variables X_j ($j = 1, 2, \dots, p$), which we call independent variables:

$$E(Y|X) = \alpha + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p. \quad (1)$$

Furthermore we might specify the variation of Y around its mean value by using the model

$$Y_i = \alpha + \beta_1 X_{1,i} + \beta_2 X_{2,i} + \dots + \beta_p X_{p,i} + \epsilon_i = E(Y|X_i) + \epsilon_i, \quad (2)$$

where ϵ_i is a Gaussian error term with mean 0 and $X_{j,i}$ is the j th coefficient (or predictor) for the i th observation ($i = 1, 2, \dots, n$). The model can be written in matrix for as

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} \quad (3)$$

or

$$\begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix} = \begin{pmatrix} 1 & X_{1,1} & X_{1,2} & \dots & X_{1,p} \\ 1 & X_{2,1} & X_{2,2} & \dots & X_{2,p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n,1} & X_{n,2} & \dots & X_{n,p} \end{pmatrix} \begin{pmatrix} \alpha \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{pmatrix}. \quad (4)$$

The least squares problem consists now in minimizing the sum of the squares of the difference between the mean value of Y and its realization values Y_i . That is, looking for the coefficients $\hat{\boldsymbol{\beta}}$ that minimize

$$F = \sum_{i=1}^n \left(Y_i - \alpha - \sum_{j=1}^p \beta_j X_{i,j} \right)^2, \quad (5)$$

which has the solution

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}. \quad (6)$$

The fitted values have then the form

$$\begin{aligned} \hat{\mathbf{Y}} &= \mathbf{X} \hat{\boldsymbol{\beta}} \\ &= \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ &= \mathbf{H} \mathbf{Y}. \end{aligned} \quad (7)$$

0.1 Coefficient of multiple determination

The coefficient of multiple determination R^2 is defined as the ratio of the sum of the squared difference between the fitted and the average values

$$\begin{aligned} SSR &= \sum_{i=1}^n \left(\hat{Y}_i - \frac{1}{n} \sum_{i=1}^n Y_i \right)^2 \\ &= \mathbf{Y}^T \left[\mathbf{H} - \frac{1}{n} \mathbf{J} \right] \mathbf{Y} \end{aligned} \quad (8)$$

where \mathbf{J} is an $n \times n$ matrix of ones, and the sum of the squared difference between the actual and average values (proportional to the variance)

$$\begin{aligned} SST &= \sum_{i=1}^n \left(Y_i - \frac{1}{n} \sum_{i=1}^n Y_i \right)^2 \\ &= \mathbf{Y}^T \left[\mathbf{I} - \frac{1}{n} \mathbf{J} \right] \mathbf{Y}, \end{aligned} \quad (9)$$

where \mathbf{I} is the $n \times n$ identity matrix. That is

$$R^2 = \frac{SSR}{SST}. \quad (10)$$

The estimated variance is

$$\hat{\sigma}^2 = \frac{SST}{n - p - 1}. \quad (11)$$

0.2 Maximum likelihood

In our model we have assumed that the random variable Y has a Gaussian distribution around its mean and, assuming that the events are independent, the joint probability distribution for the values of \mathbf{Y} obtained in the experiments given their mean and variance is

$$f(\mathbf{Y} | \mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}) = \left(\frac{1}{2\pi\sigma^2} \right)^{-n/2} e^{-\frac{1}{2\sigma^2} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})}, \quad (12)$$

and the logarithm of this (called the log likelihood) is

$$\ln[f(\mathbf{Y} | \mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})] = -\frac{1}{2\sigma^2} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) + c \quad (13)$$

where c does not depend on $\boldsymbol{\beta}$. The maximum likelihood estimate for the parameters $\tilde{\boldsymbol{\beta}}$ is the one that maximizes eq (13). After differentiating and equating to 0 we get that

$$\begin{aligned} \tilde{\boldsymbol{\beta}} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ &= \hat{\boldsymbol{\beta}}. \end{aligned} \quad (14)$$