



Multiple Linear Regression

1

What you are going to learn

- Concept of Multiple Linear Regression
- Special Cases
- Properties of the Residual
- Multiple Linear Regression in Matrix Notation
- Identification of the Coefficients
- The Coefficient of Determination
- Multiple Linear Quasi-Regression
- Statistical Models of Multiple Linear Regressions
- The General Linear Model



Multiple Linear Regression: Definition

2

Let Y and X_1, \dots, X_m be numerical random variables on a common probability space with finite expected values, positive, finite variances and a non-singular covariance matrix S_{xx} . The regression $E(Y | X_1, \dots, X_m)$ is called *linear in* (X_1, \dots, X_m) , if

$$E(Y | X_1, \dots, X_m) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_m X_m.$$



Multiple Linear Regression: Special Case I

3

Note that in the definition presented on the last transparency it is not assumed that the regressors X_1, \dots, X_m are defined independently of each other. This is only necessary if we consider the conditional regressions $E_{X_2=x_2, \dots, X_m=x_m}(Y | X_1)$.

We have already seen in Chap. 9, that a simple quadratic regression $E(Y | X)$ is linear in (X, X^2) , if

$$E(Y | X) = \beta_0 + \beta_1 X + \beta_2 X^2.$$



Multiple Linear Regression: Special Case II

4

A second special case of the multiple linear regression is the following:

$$E(Y | X_1, X_2) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 \cdot X_2.$$

Note that the regression $E(Y | X_1, X_2)$ is linear in $(X_1, X_2, X_1 \cdot X_2)$ and that

$$E(Y | X_1, X_2) = E(Y | X_1, X_2, X_1 \cdot X_2).$$



Multiple Linear Regression: Special Case III

5

If the regressor X has n values x_1, \dots, x_n , the cell means model

$$E(Y|X) = \beta_0 + \beta_1 \cdot I_1 + \beta_2 \cdot I_2 + \dots + \beta_n \cdot I_n, \text{ with } \beta_0 = 0,$$

is also a special case of a multiple linear regression. The regression $E(Y|X)$ is always linear in (I_1, \dots, I_n) , even if the regression $E(Y|X)$ is not linear in (X) .



The Properties of the Residual

6

The residual is defined by

$$\mathbf{e} := Y - E(Y|X_1, \dots, X_m)$$

It has the following properties

$$E(\mathbf{e}) = 0,$$

$$E(\mathbf{e} | X_1, \dots, X_m) = 0,$$

$$\text{Cov}[\mathbf{e}, f(X_1, \dots, X_m)] = 0,$$

$$\text{Cov}(\mathbf{e}, X_i) = 0 \quad \text{for } i = 1, \dots, m,$$

where $f(X_1, \dots, X_m)$ can be any numerical function of the regressors.



Matrix Notation I

7

If you gather all regressors in one row vector $\mathbf{x}' = (X_1 \dots X_m)$ and the regression coefficients β_1, \dots, β_m in an m -dimensional column vector $\mathbf{b} = (\beta_1 \dots \beta_m)'$, then you can write the multiple linear regression in matrix notation:

$$E(\mathbf{y} | \mathbf{x}) = \beta_0 + \mathbf{x}' \mathbf{b} = \beta_0 + (X_1 \dots X_m) \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_m \end{pmatrix}.$$

The regressand $\mathbf{y} = (Y)$ is a vector with only one component, Y . Therefore β_0 is a real number.



Matrix Notation II

8

If we define the row vector $\mathbf{z}' := (1 \ X_1 \dots X_m)$ and the column vector $\mathbf{g} := (\beta_0 \ \beta_1 \dots \beta_m)'$, the last equation simplifies to:

$$E(\mathbf{y} | \mathbf{x}) = E(\mathbf{y} | \mathbf{z}) = \mathbf{z}' \mathbf{g} = (1 \ X_1 \dots X_m) \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_m \end{pmatrix}.$$



Identification of the Coefficients I

9

In order to identify β_0 and the components of $\mathbf{b} = (\beta_1 \dots \beta_m)'$ we need the expected values of the regressand and the regressors as well as the covariance matrices S_{xx} and S_{xy} . The identification equation of the constant of the regression β_0 is:

$$\beta_0 = E(y) - E(\mathbf{x})' \mathbf{b} = E(Y) - [E(X_1) \dots E(X_m)] \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_m \end{pmatrix}.$$



Identification of the Coefficients II

10

Pre-multiplying both sides of the following equation

$$\begin{aligned} S_{xy} = \text{Cov}(\mathbf{x}, y) &= \text{Cov}(\mathbf{x}, \beta_0 + \mathbf{x}' \mathbf{b} + e) = \text{Cov}(\mathbf{x}, \beta_0 + \mathbf{b}' \mathbf{x} + e) \\ &= \text{Cov}(\mathbf{x}, \mathbf{x}) \mathbf{b} = S_{xx} \mathbf{b} \end{aligned}$$

by the inverse S_{xx}^{-1} yields

$$S_{xx}^{-1} S_{xx} \mathbf{b} = S_{xx}^{-1} S_{xy}.$$

Because $S_{xx}^{-1} S_{xx} = \mathbf{I}$ is the identity matrix, this simplifies to

$$\mathbf{b} = S_{xx}^{-1} S_{xy},$$

which is the general formula for the identification of the coefficients of the regression.



Identification of the Coefficients: Two Regressors I 11

Let's say you have two regressors X_1 and X_2 and apply these formulas. This leads to the same identification equations given in Chapter 9:

$$\beta_0 = E(Y) - E(x)' b = E(Y) - [\beta_1 E(X_1) + \beta_2 E(X_2)] = E(Y) - \beta_1 E(X_1) - \beta_2 E(X_2).$$

The covariance matrix and the variance covariance matrix are

$$S_{xx} = \begin{pmatrix} \text{Var}(X_1) & \text{Cov}(X_1, X_2) \\ \text{Cov}(X_1, X_2) & \text{Var}(X_2) \end{pmatrix} \text{ and } S_{xy} = \begin{pmatrix} \text{Cov}(X_1, Y) \\ \text{Cov}(X_2, Y) \end{pmatrix}.$$

Cramer's Rule as given in the last chapter yields

$$S_{xx}^{-1} = \frac{1}{\text{Var}(X_1)\text{Var}(X_2) - \text{Cov}(X_1, X_2)^2} \begin{pmatrix} \text{Var}(X_2) & -\text{Cov}(X_1, X_2) \\ -\text{Cov}(X_1, X_2) & \text{Var}(X_1) \end{pmatrix}.$$



Identification of the Coefficients: Two Regressors II 12

Multiplying S_{xx}^{-1} with S_{xy} yields

$$\beta_1 = \frac{\text{Var}(X_2) \text{Cov}(X_1, Y) - \text{Cov}(X_2, Y) \text{Cov}(X_1, X_2)}{\text{Var}(X_1) \text{Var}(X_2) - \text{Cov}(X_1, X_2)^2},$$

$$\beta_2 = \frac{\text{Var}(X_1) \text{Cov}(X_2, Y) - \text{Cov}(X_1, Y) \text{Cov}(X_1, X_2)}{\text{Var}(X_1) \text{Var}(X_2) - \text{Cov}(X_1, X_2)^2}.$$



Coefficient of Determination

13

The variance $Var[E(y|x)]$ of the regression is:

$$\begin{aligned} Var[E(y|x)] &= Var(\beta_0 + \mathbf{x}'\mathbf{b}) = Var(\mathbf{x}'\mathbf{b}) = Var(\mathbf{b}'\mathbf{x}) = \mathbf{b}' Var(\mathbf{x}) \mathbf{b} = \\ &= \mathbf{b}' S_{xx} \mathbf{b}, \end{aligned}$$

where $Var(\mathbf{x}) = S_{xx}$ is the $(m \times m)$ variance-covariance matrix of the regressors X_1, \dots, X_m and \mathbf{b} is the m -dimensional column vector of the regression coefficients β_1, \dots, β_m . Hence, the multiple coefficient of determination is:

$$R_{Y|X_1, \dots, X_m}^2 = [\mathbf{b}' S_{xx} \mathbf{b}] / Var(Y).$$



Coefficient of Determination: Special Cases

14

If you have only two regressors X_1 and X_2 , the last equation yields

$$R_{Y|X_1, X_2}^2 = [\beta_1^2 Var(X_1) + \beta_2^2 Var(X_2) + 2\beta_1 \beta_2 Cov(X_1, X_2)] / Var(Y),$$

which is identical to the equation of Chapter 9.

If all regressors X_1, \dots, X_m are *pairwise uncorrelated*, then the equation simplifies to

$$R_{Y|X_1, \dots, X_m}^2 = \left(\sum_{i=1}^m \beta_i^2 Var(X_i) \right) \frac{1}{Var(Y)}.$$



Multiple Linear Quasi-Regression: Definition

15

Definition 14.2. Given the assumptions of Definition 14.1, we define the *multiple linear Quasi-Regression*, denoted $Q(Y | X_1, \dots, X_m)$ or $Q(y | \mathbf{x})$, as that linear combination $\beta_0 + \beta_1 X_1 + \dots + \beta_m X_m = \beta_0 + \mathbf{b}' \mathbf{x}$ of the components of $\mathbf{x} = (X_1 \dots X_m)'$, for which

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_m X_m + \mathbf{n}$$

holds with

$$E(\mathbf{n}) = 0,$$

and

$$\text{Cov}(\mathbf{n}, X_1) = \dots = \text{Cov}(\mathbf{n}, X_m) = 0.$$



Multiple Linear Quasi-Regression: Alternative Definition

16

Definition 14.3. Given the same assumptions as before, we can also define $Q(y | \mathbf{x})$, as the linear combination $\beta_0 + \mathbf{b}' \mathbf{x}$, that minimizes the following function of b_0 and \mathbf{b} , the *least squares criterion*

$$LS(b_0, \mathbf{b}) = E[[Y - (b_0 + \mathbf{x}' \mathbf{b})]^2].$$

The number b_0 and the vector \mathbf{b} , for which the function $LS(b_0, \mathbf{b})$ is minimal, are denoted with β_0 and \mathbf{b} , respectively. The multiple linear quasi-regression is then defined by:

$$Q(y | \mathbf{x}) = \beta_0 + \mathbf{x}' \mathbf{b}.$$



The Coefficient of Determination of the Multiple Linear Quasi-Regression

17

The identification formulas for the coefficients of the multiple linear quasi-regression are identical to the ones of the (real) multiple linear regression.

The identification for the coefficient of determination of the multiple linear quasi-regression is also identical to the one of the multiple linear regression.

It is

$$Q_{Y|X_1, \dots, X_m}^2 := \text{Var}[Q(\mathbf{y} | \mathbf{x})] / \text{Var}(Y) = [\mathbf{b} \hat{\mathbf{c}} \mathbf{S}_{xx} \mathbf{b}] / \text{Var}(Y).$$



Statistical Models for the Multiple Linear Regression

18

So far we have always considered a single unit trial: Drawing an observational unit u out of a population and registering the values of the regressand and the regressors. Statistical models, however, refer to N random experiments, in which information about the parameters to be estimated is collected.



Models With Stochastic Regressors

19

Models with stochastic regressors consider the case, in which the single unit trial is repeated N -times. This yields N vectors $(Y_i X_{i1} \dots X_{im})$, $i = 1, \dots, N$, each representing the outcome of the random experiment i . Different assumptions about the distribution of these vectors can be made, e.g. that the vectors $(Y_i X_{i1} \dots X_{im})$ are independent and each one is $(m + 1)$ -variate normally distributed.



Models With Fixed Regressors

20

Other models estimate, *within* each of the combinations x_1, \dots, x_m of the regressors X_1, \dots, X_m , the expected values $E(Y | X_1 = x_1, \dots, X_m = x_m)$ of Y . Hence, within each of these combinations the regressand Y is observed many times. The values x_1, \dots, x_m are not chosen at random but are considered to be fixed quantities, that characterize the design of the experiment. This is why these models are called models with *fixed* or *non-stochastic* regressors. The most important of these models with fixed regressors is the General Linear Model (GLM).



The General Linear Model I

21

The GLM is defined by the following assumptions:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$$
$$\mathbf{e} \sim \mathcal{N}(\mathbf{0}, s^2 \mathbf{I}),$$

where $\mathbf{y} = (Y_1 \dots Y_i \dots Y_N)'$ denotes the column vector of the “dependent” variables of a sample with size N . The *design matrix* \mathbf{X} consists of $N \times (m + 1)$ fixed numbers. Every row of \mathbf{X} consists of the vectors $\mathbf{x}_i' := (1 \ x_{i1} \dots \ x_{im})$, which are the combinations of values of the regressors, within which the Y_i is observed. The constant 1 in the first position of \mathbf{x}_i' implies that the regression constant β_0 is the first component of $\boldsymbol{\beta} = (\beta_0 \ \beta_1 \dots \ \beta_m)'$. The vector $\mathbf{e} = (e_1 \dots \ e_i \dots \ e_N)'$ consists of the residuals $e_i := Y_i - (\mathbf{x}_i' \boldsymbol{\beta})$.



The General Linear Model II

22

The second assumption $\mathbf{e} \sim \mathcal{N}(\mathbf{0}, s^2 \mathbf{I})$ implies that \mathbf{e} has a multivariate normal distribution with the vector of expected values $E(\mathbf{e}) = \mathbf{0}$ and the $N \times N$ covariance matrix $\mathbf{S}_{ee} = s^2 \mathbf{I}$. That is, the residuals e_i are uncorrelated and have all the same variances. The latter is called the assumption of homoscedasticity.

Implications of these assumptions are

$$E(\mathbf{y}) = \mathbf{X}\boldsymbol{\beta},$$

and

$$\mathbf{S}_{yy} = s^2 \mathbf{I}.$$



The General Linear Model: Identification

23

If $\mathbf{X}'\mathbf{X}$ is non-singular then the vector of the regression coefficients can be estimated by:

$$\hat{\mathbf{b}} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y}.$$

This formula is derived by minimizing the least-square criterion

$$LS(\mathbf{b}) = (\mathbf{y} - \mathbf{X}\mathbf{b})'(\mathbf{y} - \mathbf{X}\mathbf{b}).$$

Of interest is also the covariance matrix

$$S_{\hat{\mathbf{b}}\hat{\mathbf{b}}} = s^2 (\mathbf{X}'\mathbf{X})^{-1}$$

of these estimates. The square roots of the components on the main diagonal of this matrix are the standard errors of the regression coefficients.



Estimating the Coefficient of Determination

24

The formula to estimate the coefficient of determination is

$$\hat{R}^2 = \frac{\mathbf{y}'\mathbf{X}\hat{\mathbf{b}} - N \cdot \bar{Y}^2}{\mathbf{y}'\mathbf{y} - N \cdot \bar{Y}^2} = \frac{\text{sum of squares of the regression}}{\text{total sum of squares}},$$

where $\bar{Y} = (1/N) \cdot \sum_{i=1}^N Y_i$.



Specifying the Null Hypothesis I

25

There are two strategies of specifying hypotheses about dependencies:

First strategy: Comparing

$$Q(Y | X_1, X_2, \dots, X_{m-p}) = \gamma_0 + \gamma_1 X_1 + \dots + \gamma_{m-p} X_{m-p}$$

with

$$E(Y | X_1, \dots, X_m) = \beta_0 + \beta_1 X_1 + \dots + \beta_m X_m.$$

The null hypothesis tested can be written in three different versions:

$$H_0: \beta_{m-p+1} = \beta_{m-p+2} = \dots = \beta_m = 0 \quad \text{version 1}$$

$$H_0: Q(Y | X_1, \dots, X_{m-p}) = E(Y | X_1, X_2, \dots, X_m). \quad \text{version 2}$$

$$H_0: R_{Y|X_1, \dots, X_m}^2 - Q_{Y|X_1, \dots, X_{m-p}}^2 = 0 \quad \text{version 3}$$



Specifying Hypothesis II

26

Second strategy: A more general way is testing the *General Linear Hypothesis* (GLH)

$$H_0: \mathbf{A} \mathbf{b} - \mathbf{d} = \mathbf{0},$$

where the matrix \mathbf{A} and the vector \mathbf{d} are specified according to the hypothesis to be tested. The matrix \mathbf{A} has to contain $p \leq m$ linear independent rows.



Test of Significance within the GLM I

27

Testing a null hypothesis in version 1

$$H_0: \beta_{m-p+1} = \beta_{m-p+2} = \dots = \beta_m = 0,$$

that some of the coefficients of the multiple linear regression are zero, is done by estimating \hat{R}_E^2 and \hat{R}_Q^2 as seen before. Again, please note that this is done by different design matrices and different regression coefficients, i.e. with m regressors for the multiple linear regression and $m - p$ regressors for the multiple linear quasi-regression. Given the assumptions of the GLM the following statistic

$$F = \frac{(\hat{R}_E^2 - \hat{R}_Q^2) / p}{(1 - \hat{R}_E^2) / (N - m - 1)},$$

is F -distributed with $df_1 = p$ and $df_2 = N - m - 1$, where N is the sample size.



Test of Significance within the GLM II

28

For the second strategy – the general linear hypothesis (GLH)

$$H_0: \mathbf{A} \mathbf{b} - \mathbf{d} = \mathbf{0},$$

the following statistic is computed:

$$F = \frac{\hat{Q}_h / p}{\hat{Q}_e / (N - m - 1)},$$

where p is the number of (linear independent) rows of the matrix \mathbf{A} of the GLH (the number of simultaneously tested single hypotheses), and further

$$\hat{Q}_h = (\mathbf{A}\hat{\mathbf{b}} - \mathbf{d})' [\mathbf{A}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{A}']^{-1} (\mathbf{A}\hat{\mathbf{b}} - \mathbf{d}) \quad \text{SSQ of the hypothesis}$$

$$\hat{Q}_e = \mathbf{y}' \mathbf{y} - \mathbf{y}' \mathbf{X} \hat{\mathbf{b}}. \quad \text{SSQ of the error term}$$

This statistic is also F -distributed with $df_1 = p$ the degrees of freedom for the numerator and $df_2 = N - m - 1$ the degrees of freedom for the denominator.