

Multiple Linear Regression

Hong Kong Baptist University

Fall 2021

Multiple Linear Regression

Multiple linear regression model:

$$Y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_k x_{ki} + \epsilon_i, \quad i = 1, \dots, n, \quad (1)$$

where $n \geq k + 1$, $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$ are i.i.d, and x_{ji} is the i^{th} observation of the j^{th} independent variable.

General linear model is linear in its parameters. The independent variables can be higher order terms like x^2 or $\log x$. That is, the model

$$Y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 \log x + \epsilon$$

is still a linear model; it is linear in its parameters. An example of nonlinear model is

$$Y = \beta_0 + \beta_1 x^{\beta_2} + \epsilon,$$

which is not linear in its parameters.

The model given in (1) is in fact a system of linear equations, and can be expressed in matrix terms as follows:

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad (2)$$

where

$$\mathbf{Y} = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{21} & \cdots & x_{k1} \\ 1 & x_{12} & x_{22} & \cdots & x_{k2} \\ \vdots & \vdots & & & \vdots \\ 1 & x_{1n} & x_{2n} & \cdots & x_{kn} \end{bmatrix}, \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix}, \quad \boldsymbol{\varepsilon} = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}.$$

Note that x_{ij} is the i^{th} row j^{th} column element of the design matrix \mathbf{X} . We have a column of 1's in \mathbf{X} because we include β_0 in $\boldsymbol{\beta}$ so that in equation (2) we do not need a separate intercept term.

Revision of matrix algebra

Let us denote a matrix by

$$\mathbf{A}_{n \times m} = \mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1m} \\ a_{21} & a_{22} & \cdots & a_{2m} \\ \vdots & \vdots & & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nm} \end{bmatrix} = [a_{ij}]_{n \times m}.$$

Matrix addition/subtraction is done elementwise:

$$\mathbf{A} \pm \mathbf{B} = [a_{ij} \pm b_{ij}],$$

and so is scalar multiplication:

$$k\mathbf{A} = [ka_{ij}],$$

but matrix multiplication is more complicated:

$$\mathbf{A}_{n \times m} \mathbf{B}_{m \times p} = \left[\sum_{k=1}^m a_{ik} b_{kj} \right]_{n \times p},$$

and so in general $\mathbf{AB} \neq \mathbf{BA}$.

The transpose of a matrix $A = [a_{ij}]_{n \times m}$ is denoted by A' , which is defined by

$$A' = [a_{ji}]_{m \times n},$$

and $(AB)' = B'A'$. Obviously $(A')' = A$.

If $A' = A$, then A is symmetric; it is of course a square matrix, in which $n = m$. A simple example of a symmetric matrix is $A'A$ because $(A'A)' = A'A$.

The identity matrix I is the square matrix that the diagonal elements are all 1, whilst all others are zero, and so $IA = AI = A$.

The inverse of a square matrix, denoted by A^{-1} , if exists, is the unique matrix satisfying $AA^{-1} = A^{-1}A = I$. If A^{-1} exists, then we say that A is non-singular.

A system of linear equations, expressed in matrix terms, is

$$\mathbf{A}\mathbf{x} = \hat{\boldsymbol{\beta}},$$

and the solution is simply

$$\mathbf{x} = \mathbf{A}^{-1}\hat{\boldsymbol{\beta}}.$$

Let $\mathbf{x} = [x_1, x_2, \dots, x_n]'$ (a column vector) and y be a real-valued function of x_1, x_2, \dots, x_n , then define the derivative of y with respect to the column vector \mathbf{x} to be the column vector of partial derivatives:

$$\frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \vdots \\ \frac{\partial y}{\partial x_n} \end{bmatrix}_{n \times 1}.$$

If $\mathbf{a} = [a_1, a_2, \dots, a_n]'$ (a vector of constants; constants mean they are not functions of \mathbf{x}) and $y = \mathbf{a}'\mathbf{x} = a_1x_1 + a_2x_2 + \dots + a_nx_n$, then we have

$$\frac{\partial \mathbf{a}'\mathbf{x}}{\partial \mathbf{x}} = \frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \vdots \\ \frac{\partial y}{\partial x_n} \end{bmatrix} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} = \mathbf{a}.$$

Suppose \mathbf{A} is a matrix of constants:

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1m} \\ a_{21} & a_{22} & \cdots & a_{2m} \\ \vdots & \vdots & & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nm} \end{bmatrix}_{n \times m},$$

and

$$\begin{aligned} \mathbf{Y} = \mathbf{x}'\mathbf{A} &= [x_1, x_2, \dots, x_n] \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1m} \\ a_{21} & a_{22} & \cdots & a_{2m} \\ \vdots & \vdots & & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nm} \end{bmatrix} \\ &= [(a_{11}x_1 + a_{21}x_2 + \cdots + a_{n1}x_n), \dots, (a_{1m}x_1 + a_{2m}x_2 + \cdots + a_{nm}x_n)] \end{aligned}$$

Then for the row vector $Y = x'A$, its partial derivative with respect to the column vector x will be a matrix:

$$\begin{aligned} \frac{\partial x'A}{\partial x} &= \frac{\partial Y}{\partial x} = \left[\frac{\partial y_1}{\partial x}, \frac{\partial y_2}{\partial x}, \dots, \frac{\partial y_m}{\partial x} \right] \\ &= \begin{bmatrix} \frac{\partial Y}{\partial x_1} \\ \frac{\partial Y}{\partial x_2} \\ \vdots \\ \frac{\partial Y}{\partial x_n} \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1m} \\ a_{21} & a_{22} & \cdots & a_{2m} \\ \vdots & \vdots & & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nm} \end{bmatrix} = A. \end{aligned} \quad (3)$$

If $y = \mathbf{x}'\mathbf{A}\mathbf{x} = \sum_{i=1}^n \sum_{j=1}^n a_{ij}x_ix_j$, then applying the product rule:

$$\frac{\partial \mathbf{u}'\mathbf{A}\mathbf{v}}{\partial \mathbf{x}} = \frac{\partial \mathbf{u}'}{\partial \mathbf{x}}\mathbf{A}\mathbf{v} + \frac{\partial \mathbf{v}'}{\partial \mathbf{x}}\mathbf{A}'\mathbf{u},$$

for vector-valued functions \mathbf{u} , \mathbf{v} of \mathbf{x} , we have

$$\frac{\partial \mathbf{x}'\mathbf{A}\mathbf{x}}{\partial \mathbf{x}} = \frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial \sum \sum a_{ij}x_ix_j}{\partial x_1} \\ \vdots \\ \frac{\partial \sum \sum a_{ij}x_ix_j}{\partial x_n} \end{bmatrix} = \frac{\partial \mathbf{x}'}{\partial \mathbf{x}}\mathbf{A}\mathbf{x} + \frac{\partial \mathbf{x}'}{\partial \mathbf{x}}\mathbf{A}'\mathbf{x} = (\mathbf{A} + \mathbf{A}')\mathbf{x}.$$

In particular, if \mathbf{A} is symmetric so that $\mathbf{A} = \mathbf{A}'$, then we have

$$\frac{\partial \mathbf{x}'\mathbf{A}\mathbf{x}}{\partial \mathbf{x}} = \frac{\partial y}{\partial \mathbf{x}} = 2\mathbf{A}\mathbf{x}. \quad (4)$$

Least Squares Estimation

Now, the least squares estimation for the model parameter β of the multiple regression model

$$Y = X\beta + \varepsilon$$

can be obtained straightforwardly. Consider

$$\varepsilon = Y - X\beta.$$

The sum of squares of errors (which is a scalar) is equal to

$$Q = \varepsilon'\varepsilon = (Y - X\beta)'(Y - X\beta) = Y'Y - 2\beta'X'Y + \beta'X'X\beta$$

(notice that $\beta'X'Y$ and $Y'X\beta$ are the same because they are scalars).

To minimize Q we take the partial derivative with respect to β . Using equations (3) and (4) above, we get

$$\frac{\partial Q}{\partial \beta} = -2\mathbf{X}'\mathbf{Y} + 2\mathbf{X}'\mathbf{X}\beta.$$

The least squares estimator $\hat{\beta}$ for the parameter β should satisfy the *normal equation*

$$\left. \frac{\partial Q}{\partial \beta} \right|_{\beta=\hat{\beta}} = \mathbf{0} \Rightarrow \boxed{\mathbf{X}'\mathbf{Y} = \mathbf{X}'\mathbf{X}\hat{\beta}},$$

and its solution is

$$\boxed{\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}}. \quad (5)$$

From this elegant expression, we can see immediately that if ε has a normal distribution, then \mathbf{Y} has a normal distribution, and so does $\hat{\beta}$ [which is a non-random matrix $(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ times the normally distributed \mathbf{Y}].

Multivariate Normal Distribution

We say a column vector has a multivariate normal distribution, denoted by

$$\mathbf{X}_{n \times 1} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}),$$

if and only if the joint probability density function $f_{\mathbf{X}}$ is given by

$$f_{\mathbf{X}}(\mathbf{x}) = \left(\frac{1}{\sqrt{2\pi}} \right)^n \frac{1}{\sqrt{|\boldsymbol{\Sigma}|}} e^{-\frac{(\mathbf{x}-\boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{x}-\boldsymbol{\mu})}{2}}, \quad \text{for all } \mathbf{x} \in \mathbb{R}^n,$$

where $|\boldsymbol{\Sigma}|$ denotes the determinant of $\boldsymbol{\Sigma}$.

Maximum Likelihood Estimation

Now, if $\varepsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ (i.e. ε_i are i.i.d. normal with mean zero and variance σ^2 ; the independence between ε_i results in a diagonal covariance matrix because the covariance between ε_i and ε_j is the (i, j) -element in Σ ; the identical distribution assumption results in one common value σ^2 for all diagonal elements in Σ), then the likelihood is just

$$\ell(\boldsymbol{\beta}) = \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right)^n e^{-\frac{\boldsymbol{\varepsilon}'\boldsymbol{\varepsilon}}{2\sigma^2}} = \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right)^n e^{-\frac{(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})}{2\sigma^2}}.$$

So maximizing ℓ is equivalent to minimizing $(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})$, i.e.

$$\text{MLE} = \text{LSE}. \quad (6)$$

Note that the equivalence given in (6) is true if we have normally distributed errors. If the errors follow another distribution, it may not be true. For example, if the errors are i.i.d., following the Laplace distribution, i.e.

$$f(\varepsilon) = \frac{1}{2\sigma} e^{-\frac{|\varepsilon|}{\sigma}}, \quad \text{for all } \varepsilon \in \mathbb{R},$$

then the likelihood is given by

$$\ell(\beta) = \left(\frac{1}{2\sigma}\right)^n e^{-\frac{\sum_i |\varepsilon_i|}{\sigma}},$$

and so maximizing the likelihood ℓ is equivalent to minimizing $\sum_i |\varepsilon_i|$ (least absolute deviations, or LAD for short).

Under LSE, each deviation is squared, i.e. each of them is weighted by itself, while under LAD, each deviation carries the same weight. Therefore, an outlier will have a stronger influence on the LSE than on the LAD, as a larger deviation will play a more dominant role in LSE than in LAD.

Geometry of the LSE

Recall that two vectors p and q are *orthogonal* if $p'q = q'p = 0$. Now, consider the fitted value \hat{Y} and residual e :

$$\hat{Y} = X\hat{\beta} = X(X'X)^{-1}X'Y \equiv HY, \quad e = Y - \hat{Y} = Y - X\hat{\beta} = (I - H)Y,$$

where $H = X(X'X)^{-1}X'$ is the hat matrix. The actual observation Y is in \mathbb{R}^n . Because

$$\hat{Y} = X\hat{\beta} = \hat{\beta}_0 \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} + \hat{\beta}_1 \begin{bmatrix} x_{11} \\ x_{12} \\ \vdots \\ x_{1n} \end{bmatrix} + \cdots + \hat{\beta}_k \begin{bmatrix} x_{k1} \\ x_{k2} \\ \vdots \\ x_{kn} \end{bmatrix}$$

the fitted value \hat{Y} is lying in the so-called *estimation space*, which is the space generated by the $p = k + 1$ columns of X (we use p to represent the number of parameters, including the intercept β_0), whilst the residual e is lying in the *error space*, which has $n - p = n - k - 1$ dimensions.

Since

$$\hat{Y}'e = (\mathbf{X}\hat{\beta})'e = \hat{\beta}'\mathbf{X}'(\mathbf{Y} - \mathbf{X}\hat{\beta}) = \hat{\beta}'(\mathbf{X}'\mathbf{Y} - \mathbf{X}'\mathbf{X}\hat{\beta}) = 0,$$

(note that $\mathbf{X}'\mathbf{Y} = \mathbf{X}'\mathbf{X}\hat{\beta}$ is the normal equation), we can see that the fitted value \hat{Y} and residual e are orthogonal, and hence the name *normal* equation, because under which the residual e is a normal vector (i.e. a vector perpendicular) to the estimation space.

Therefore, the least squares fitting procedure splits the space \mathbb{R}^n into two orthogonal spaces; every vector in the estimation space is orthogonal to every vector in the error space, and the fitted value \hat{Y} is the orthogonal projection of the observed Y to the estimation space ($\hat{Y} = \mathbf{H}\mathbf{Y}$, $\mathbf{H}\hat{Y} = \mathbf{H}\mathbf{H}\mathbf{Y} = \mathbf{H}\mathbf{Y} = \hat{Y}$, $e'\hat{Y} = 0$), so that no other vector in the estimation space is closer (measured in the Euclidean distance) to Y than \hat{Y} is.

The unbiasedness of $\hat{\beta}$ (which will be shown below) implies

$$\mathbb{E}(e) = \mathbb{E}(Y - X\hat{\beta}) = \mathbb{E}(X\beta) + \mathbb{E}(\varepsilon) - \mathbb{E}(X\hat{\beta}) = X\beta + \mathbf{0}_{n \times 1} - X\beta = \mathbf{0}_{n \times 1},$$

where $\mathbf{0}_{l \times m}$ is the $l \times m$ matrix of zeros. And we can show that these two orthogonal random vectors \hat{Y} and e are uncorrelated:

$$\begin{aligned} \text{cov}(\hat{Y}, e) &= \text{cov}\{HY, (I - H)Y\} \\ &= \mathbb{E}[H(Y - X\beta)(Y - X\beta)'(I - H)'] = H\text{cov}(Y)(I - H) \\ &= H(\sigma^2 I)(I - H) = \sigma^2 H(I - H) = \sigma^2(H - H) = \mathbf{0}_{n \times n}. \end{aligned}$$

The expression given in equation (5) allow us to derive the statistical properties of LSE. The mean and the covariance matrix of $\hat{\beta}$ are

$$\begin{aligned} \mathbb{E}(\hat{\beta}) &= \mathbb{E}\{(X'X)^{-1}X'Y\} = (X'X)^{-1}X'\mathbb{E}(Y) = (X'X)^{-1}X'X\beta = \beta \\ \text{cov}(\hat{\beta}) &= \text{cov}\{(X'X)^{-1}X'Y\} =: \text{cov}(AY) = A\text{cov}(Y)A' \\ &= A(\sigma^2 I)A' = \sigma^2 AA' = \sigma^2\{(X'X)^{-1}X'\}\{X(X'X)^{-1}'\} \\ &= \sigma^2(X'X)^{-1}(X'X)(X'X)^{-1} = \sigma^2(X'X)^{-1}, \end{aligned}$$

where for ease of presentation we wrote $A := (X'X)^{-1}X'$.

Finally, to work out the distribution of the estimator $\hat{\beta}$, we just have to note that $\hat{\beta} = \mathbf{A}\mathbf{Y}$ (where $\mathbf{A} := (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$) is random because \mathbf{Y} is random, and $\mathbf{Y} = \mathbf{X}\beta + \varepsilon$ is random because ε is random. Thus, to know the distribution of $\hat{\beta}$, we have to know the distribution of ε . If $\varepsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2\mathbf{I})$, then $\mathbf{Y} \sim \mathcal{N}(\mathbf{X}\beta, \sigma^2\mathbf{I})$ and consequently, we can see

$$\hat{\beta} = \mathbf{A}\mathbf{Y} \sim \mathcal{N}(\mathbf{A}\mathbf{X}\beta, \sigma^2\mathbf{A}\mathbf{A}') = \mathcal{N}(\beta, \sigma^2(\mathbf{X}'\mathbf{X})^{-1}).$$

Let me put the formula and its distribution together in one line for your quick reference:

$$\boxed{\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y} \sim \mathcal{N}(\beta, \sigma^2(\mathbf{X}'\mathbf{X})^{-1})}.$$

Consequently, because $\hat{\mathbf{Y}}$ is simply the product of a constant matrix and $\hat{\beta}$, it also has the normal distribution:

$$\hat{\mathbf{Y}} = \mathbf{X}\hat{\beta} \sim \mathcal{N}(\mathbf{X}\beta, \sigma^2\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}').$$

For \hat{Y}_0 at some x_0 for a given fixed vector x_0 containing specified values of the independent variables

$$\mathbf{x}_0 = (1, x_{10}, x_{20}, \dots, x_{k0})'$$

using the same argument, we have

$$\hat{Y}_0 = \mathbf{x}'_0 \hat{\boldsymbol{\beta}} \sim \mathcal{N}(\mathbf{x}'_0 \boldsymbol{\beta}, \sigma^2 \mathbf{x}'_0 (\mathbf{X}' \mathbf{X})^{-1} \mathbf{x}_0).$$

Mathematically, we can see that if $\mathbf{X}'\mathbf{X}$ is singular, then $\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$ is not well defined, because $(\mathbf{X}'\mathbf{X})^{-1}$ does not exist. The implication of the singularity of $\mathbf{X}'\mathbf{X}$ is that we cannot estimate β uniquely.

In the simple linear regression, if all observed x_i are the same, there will be infinitely many regression lines that pass through (\bar{x}, \bar{Y}) .

In the multiple regression, if the matrix \mathbf{X} is such that any of its columns can be explained as a linear combination of some other columns, this dependency will be transferred to $\mathbf{X}'\mathbf{X}$ and so $\mathbf{X}'\mathbf{X}$ will have a zero determinant and be singular.

The linear dependence arises often because the data are inadequate for fitting the model or, what is the same thing, the model is too complex for the available data. We need either more data or a simpler model for the available data.

ANOVA F -Test

The ANOVA F -test is a global test, testing the null hypothesis that $\beta_1 = \dots = \beta_k = 0$ (note that the intercept β_0 is NOT included in the null hypothesis).

ANOVA table

Source of variation	Sum of Squares	degrees of freedom	Mean square	F -ratio	p -value
Model	SSR	k	$MSE = \frac{SSR}{k}$	$\frac{MSR}{MSE}$	$\Pr(F_{k, n-k-1} \geq F\text{-ratio})$
Error	SSE	$n - k - 1$	$MSE = \frac{SSE}{n-k-1}$		
Total	SST	$n - 1$			

Here

$$SST = \sum (Y_i - \bar{Y})^2 = \sum Y_i^2 - \frac{1}{n} (\sum Y_i)^2 = \mathbf{Y}' [\mathbf{I} - \frac{1}{n} \mathbf{1}\mathbf{1}'] \mathbf{Y},$$

$$\begin{aligned} SSE &= \sum (Y_i - \hat{Y}_i)^2 = \mathbf{e}'\mathbf{e} = (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}})'(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}) \\ &= \mathbf{Y}'\mathbf{Y} - \hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{Y}, \end{aligned}$$

$$\begin{aligned} SSR &= \sum (\hat{Y}_i - \bar{Y})^2 = (\mathbf{H}\mathbf{Y} - \frac{1}{n} \mathbf{1}\mathbf{1}'\mathbf{Y})'(\mathbf{H}\mathbf{Y} - \frac{1}{n} \mathbf{1}\mathbf{1}'\mathbf{Y}) \\ &= \mathbf{Y}'\mathbf{H}\mathbf{H}\mathbf{Y} - \frac{1}{n} \mathbf{Y}'\mathbf{H}\mathbf{1}\mathbf{1}'\mathbf{Y} - \frac{1}{n} \mathbf{Y}'\mathbf{1}\mathbf{1}'\mathbf{H}\mathbf{Y} + \frac{1}{n^2} \mathbf{Y}'\mathbf{1}\mathbf{1}'\mathbf{1}\mathbf{1}'\mathbf{Y} \\ &= \mathbf{Y}'\mathbf{H}\mathbf{Y} - \frac{1}{n} \mathbf{Y}'\mathbf{1}\mathbf{1}'\mathbf{Y} - \frac{1}{n} \mathbf{Y}'\mathbf{1}\mathbf{1}'\mathbf{Y} + \frac{1}{n} \mathbf{Y}'\mathbf{1}\mathbf{1}'\mathbf{Y} \\ &= \hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{Y} - \frac{1}{n} \mathbf{Y}'\mathbf{1}\mathbf{1}'\mathbf{Y} = SST - SSE, \end{aligned}$$

in which $\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ and $\mathbf{1} = \mathbf{1}_{n \times 1} = [1, \dots, 1]'$.

T-test

For each ℓ , where $0 \leq \ell \leq k$, to test the null hypothesis that the individual parameter $\beta_\ell = 0$, we use t -test:

$$T = \frac{\hat{\beta}_\ell}{\hat{\sigma} \sqrt{c_{\ell+1, \ell+1}}} \sim t_{n-k-1} \text{ when } \beta_\ell = 0,$$

where $c_{\ell+1, \ell+1}$ is the $(\ell + 1)^{\text{st}}$ diagonal element of $(\mathbf{X}'\mathbf{X})^{-1}$. The $\sigma^2 c_{\ell+1, \ell+1}$ is the variance of $\hat{\beta}_\ell$, and the variance σ^2 of ε , when unknown, is estimated by sample variance of e_i :

$$\begin{aligned} \hat{\sigma}^2 &= \frac{\sum (Y_i - \hat{Y}_i)^2}{n - k - 1} = \frac{(\mathbf{Y} - \mathbf{X}\hat{\beta})'(\mathbf{Y} - \mathbf{X}\hat{\beta})}{n - k - 1} \\ &= \frac{\mathbf{e}'\mathbf{e}}{n - k - 1} = \frac{SSE}{n - k - 1} = MSE. \end{aligned}$$

(Recall that the value of the degrees of freedom of the t -distribution in the t -test is the same as the degrees of freedom of $(n - k - 1)\hat{\sigma}^2/\sigma^2 \sim \chi_{n-k-1}^2$. It is the same as the dimension of the error space, in which the residual vector e is lying.)

Partial F -test

The ANOVA F -test tests whether all parameters except the intercept are zero, and the t -test tests whether an individual parameter is zero. Can we test whether some (more than one, but less than k) parameters are zero?

If we want to test whether a subset of parameters are all zero, then we are in fact comparing two models.

Two models are *nested* if one model contains all the terms of the second model and at least one additional term. *Reduced* (or *restricted*) model is a special case of (nested within) the *complete* (or *full*) model.

Suppose we suspect some independent variables are insignificant and we arrange the order of x_1, \dots, x_k so that the “suspicious” variables are labelled as the $(g + 1)^{\text{st}}, \dots, k^{\text{th}}$ independent variables.

We separate the independent variables into two groups and write:

$$Y = X_1\beta_1 + X_2\beta_2 + \varepsilon,$$

where $\beta_1 = [\beta_0, \beta_1, \dots, \beta_g]'$ and $\beta_2 = [\beta_{g+1}, \beta_{g+2}, \dots, \beta_k]'$, where β_0 still denotes the intercept term. (Why β_0 must belong to the first group? Can we move β_0 to the second group?)

To test the null hypothesis that $\beta_2 = \mathbf{0}$ against the alternative $\beta_2 \neq \mathbf{0}$, we use the **partial F -test**:

$$\begin{aligned} F &= \frac{\text{extra SS}/\#\text{paramters being tested}}{MSE_{\text{complete}}} \\ &= \frac{(SSE_{\text{reduced}} - SSE_{\text{complete}})/(k - g)}{SSE_{\text{complete}}/(n - k - 1)} \sim F_{k-g, n-k-1}, \end{aligned}$$

where the extra SS is the extra sum of squares of errors (extra unexplainable variation in Y_i) that we will have if we use the reduced model. The idea is that if using the reduced model, we will of course get a larger SSE (i.e. a smaller SSR [a smaller explainable variation]). But if the increase in SSE (i.e. the loss in SSR [loosely speaking, the loss in the information by the reduced model]) is not too large, then we can use the reduced model. This is exactly the **principle of parsimony**, which requires that in situations where two competing models have essentially the same predictive power, we choose the more parsimonious of the two; a parsimonious model is a model with a small number of parameters.

In fact, this partial F -test can test not only the reduced model with fewer parameters (i.e. some parameters are zero) but also reduced models with constraints on the parameters, e.g. $\beta_2 = \beta_4 = 0.5$ or $\beta_3 + \beta_6 = 2.3$, etc., and in such cases the value of the numerator degrees of freedom F -statistic is the number of constraints being tested, i.e.

$$F = \frac{\text{extra SS}/\#\text{constraints being tested}}{MSE_{\text{complete}}}.$$

However, note that the partial F -test cannot be used to compare two different models in general; the extra SS is nonnegative here because we know that SSE_{reduced} cannot be smaller than SSE_{complete} , as the reduced model is only a special case of the complete model and of course cannot be better than the complete model.

Coefficient of multiple determination

The **coefficient of multiple determination**

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

should be interpreted as, the same as that for the simple linear regression, how many percent of the variation in Y can be explained by the multiple regression model.

However, a larger value of R^2 computed from the sample data does not necessarily mean that the model provides a better fit of all of the data points in the whole population. You will always obtain a perfect fit $R^2 = 1$ for a sample of n points if the model contains exactly n parameters.

Adjusted coefficient of multiple determination

Thus, we introduce the **adjusted coefficient of multiple determination**:

$$\begin{aligned} R_a^2 &:= 1 - \frac{n-1}{n-k-1} \frac{SSE}{SST} = 1 - \frac{n-1}{n-k-1} (1 - R^2) \\ &\leq 1 - \frac{n-1}{n-1} (1 - R^2) = R^2. \end{aligned}$$

The idea is that R_a^2 takes into account (adjusted for) both the sample size and the number of parameters such that a model of more parameter will have a heavier penalty so that R_a^2 cannot be forced to 1 by simply adding more and more parameters.

One obvious disadvantage of R_a^2 is that its numerical value does not have a nice and easy to understand interpretation, while R^2 does.

Both R^2 and R_a^2 may be misleading when there are repeat observations. For example, if we have 100 observations in 5 groups (observations in each group have the same set of values for the independent variables), each of 20 repeats.

- A 5-parameter model will provide a perfect fit to the 5 sample means of Y and may give a very large value of R^2 or R_a^2 , especially if σ^2 is small compared with the spread of the 5 means.
- However, such a model, which passes through the 5 sample means of Y at the 5 observed values of x actually may be a very bad model if we could observe more distinct values of x .
- Even if we do not have repeat observations at the same value of x , we may still have 5 groups such that in each group the within-group variation is very small, whilst the between-group variation is large, leading again to a high R^2 or R_a^2 .

Therefore, we must make a scatter plot of the data before we do any analysis.

Coefficient of partial determination

The **coefficient of partial determination** is defined as

$$R_{Y,B|A}^2 = \frac{SSE_A - SSE_{A,B}}{SSE_A} = \frac{SSE_{\text{reduced}} - SSE_{\text{complete}}}{SSE_{\text{reduced}}},$$

where tells us the percentage of variation that cannot be explained in the reduced model but can be explained in the complete model.

For example, $R_{Y,x_2,x_3|x_1}^2$ will be the percentage of variation that cannot be explained in the reduced model $Y = \beta_0 + \beta_1x_1 + \varepsilon$ but can be explained in the complete model $Y = \beta_0 + \beta_1x_1 + \beta_2x_2 + \beta_3x_3 + \varepsilon$, i.e. it measures the contribution of x_2 and x_3 , when added to the model already containing x_1 , in terms of the model's explanatory power for the variation in Y .

Partial coefficient

Typically, we will consider the case that B contains only one independent variable, i.e. $R_{Y,x_i|A}^2$ (adding x_i to the model already containing the variables listed in A).

The square root of $R_{Y,x_i|A}^2$, with the sign (i.e. positive or negative) the same as the sign of $\hat{\beta}_i$ (the estimate of the coefficient of x_i in the model), is called the **partial correlation** between x_i and Y , given A .

Confidence Intervals

Next, consider confidence intervals. We know that

$$\hat{\beta} \sim \mathcal{N}(\beta, \sigma^2(\mathbf{X}'\mathbf{X})^{-1}) \text{ and } \hat{Y} = \mathbf{X}\hat{\beta} \sim \mathcal{N}(\mathbf{X}\beta, \sigma^2\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}').$$

These normal distributions would allow us, if we knew σ^2 , to construct confidence intervals for each individual β_j and for the estimation of the mean response and the prediction of an individual response.

In practice we would not know σ^2 , and thus we have to use the sample variance to estimate it:

$$\hat{\sigma}^2 = \frac{\sum(Y_i - \hat{Y}_i)^2}{n - k - 1} = \frac{(\mathbf{Y} - \mathbf{X}\hat{\beta})'(\mathbf{Y} - \mathbf{X}\hat{\beta})}{n - k - 1} = \frac{SSE}{n - k - 1} = MSE.$$

These upper and lower confidence limits are reported by SAS if we add the `clb` `clm` `cli` option after the `model` statement.

Note that if the model is underspecified (i.e. some true independent variables are missing), then $\hat{\sigma}^2$ is, on the average, an overestimate of σ^2 , because the variation in the calculated SSE comes from not only the errors (of variance σ^2) but also from the variation due to the missing independent variables which have not been accounted for in the regression model.

Because

$$\hat{\beta} \sim \mathcal{N}(\beta, \sigma^2(\mathbf{X}'\mathbf{X})^{-1}),$$

for a given fixed vector \mathbf{x}_0 containing specified values of the independent variables

$$\mathbf{x}_0 = (1, x_{10}, x_{20}, \dots, x_{k0})',$$

the point estimation of the mean response and the prediction of an individual response at $\mathbf{x} = \mathbf{x}_0$ are

$$\hat{Y}_0 = \mathbf{x}'_0 \hat{\beta} \sim \mathcal{N}(\mathbf{x}'_0 \beta, \sigma^2 \mathbf{x}'_0 (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_0).$$

and the unknown σ^2 again is estimated by

$\hat{\sigma}^2 = MSE = SSE/(n - k - 1)$. Thus, we have the following $100(1 - \alpha)\%$ confidence intervals for estimation of the mean response:

$$\hat{Y}_0 \pm t_{n-k-1; \alpha/2} \hat{\sigma} \sqrt{\mathbf{x}'_0 (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_0},$$

and for prediction of a future response:

$$\hat{Y}_0 \pm t_{n-k-1; \alpha/2} \hat{\sigma} \sqrt{1 + \mathbf{x}'_0 (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_0}.$$

Confidence Regions

We may want to construct a joint confidence region for more than one parameters, i.e. a region which, with $100(1 - \alpha)\%$ confidence, contains the true parameter vector.

Such a joint confidence region is a kind of *simultaneous inference* for several unknown parameters.

For the parameter vector β in the multiple regression model with normally distributed errors, because the LSE (also MLE) estimator $\hat{\beta} \sim \mathcal{N}(\beta, \sigma^2(\mathbf{X}'\mathbf{X})^{-1})$, if normalising this general multivariate normal to the standard multivariate normal (by subtracting the mean and “dividing” by the standard deviation) and then taking square, we will have the χ^2 -distribution:

$$\frac{1}{\sigma^2}(\hat{\beta} - \beta)' \mathbf{X}' \mathbf{X} (\hat{\beta} - \beta) \sim \chi_{k+1}^2.$$

Combining with the facts that SSE/σ^2 is also χ^2 -distributed and that they are independent, we have

$$\frac{\frac{\{(\hat{\beta} - \beta)' \mathbf{X}' \mathbf{X} (\hat{\beta} - \beta)\}}{\sigma^2}}{\frac{SSE}{\sigma^2} / (n - k - 1)} / (k + 1) \sim F_{k+1, n-k-1},$$

where $SSE/(n - k - 1) = MSE = \hat{\sigma}^2$.

Hence a $100(1 - \alpha)\%$ *joint confidence region* for β can be obtained from solving

$$(\hat{\beta} - \beta)' \mathbf{X}' \mathbf{X} (\hat{\beta} - \beta) \leq (k + 1) \hat{\sigma}^2 F_{k+1, n-k-1; \alpha}, \quad (7)$$

where $F_{k+1, n-k-1; \alpha}$ denotes the value such that with $100(1 - \alpha)\%$ probability, a random variable following the $F_{k+1, n-k-1}$ distribution is not larger than $F_{k+1, n-k-1; \alpha}$.

The equality given in (7) is the equation of the boundary of an elliptically (or ellipsoidally, when $k > 1$) shaped contour in \mathbb{R}^{k+1} , and the strict inequality is the interior of the ellipse (ellipsoid).

The joint confidence region takes into account the correlation between the estimators $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k$.

The individual confidence intervals are only appropriate for specifying the ranges for the individual parameters irrespective of the value of the other parameters.

If we (inappropriately) interpret these intervals for individual parameters simultaneously, i.e. wrongly regard the (hyper-)rectangle that they define as a joint confidence region, then e.g. it may be thought that a point lying outside the ellipse but inside the rectangle provide a reasonable value for β , but it is in fact not reasonable because it is not in the 95% joint confidence region.

When only two parameters are involved, construction of the confidence ellipse is not difficult. In practice, however, even for two parameters, it is rarely drawn.

Bonferroni Correction

What is more popular in simultaneous inference is the *Bonferroni* correction (also known as the Bonferroni adjustment, the Bonferroni procedure, etc.). It is based on the Bonferroni inequality, which is described as follows.

Let A_i denote the event that the confidence interval I_i of β_i does not cover the true β_i and suppose the confidence level is $100(1 - \gamma)\%$, i.e.

$$\Pr(\beta_i \notin I_i) = \Pr(A_i) = \gamma.$$

We might take the hyper-rectangle $I_0 \times I_1 \times \cdots \times I_k$ as a $100(1 - \alpha)\%$ joint confidence region if the following were true:

$$\Pr(\beta \in I_0 \times I_1 \times \cdots \times I_k) = \Pr(\bar{A}_0 \cap \bar{A}_1 \cap \cdots \cap \bar{A}_k) = 1 - \alpha,$$

where \bar{A}_i denotes the complement of A_i , i.e. $\bar{A}_i = \{\beta_i \in I_i\}$.

However, because these \bar{A}_i are not independent and the dependence structure between \bar{A}_i and \bar{A}_j is complicated, the best we can do is to be conservative, i.e. choose a γ such that

$$\Pr(\bar{A}_0 \cap \bar{A}_1 \cap \cdots \cap \bar{A}_k) \geq 1 - \alpha. \quad (8)$$

Now,

$$\begin{aligned} & \Pr(\bar{A}_0 \cap \bar{A}_1 \cap \cdots \cap \bar{A}_k) \\ &= 1 - \Pr(A_0 \cup A_1 \cup \cdots \cup A_k) \\ &= 1 - \{\Pr(A_0) + \cdots + \Pr(A_k) - \Pr(\text{parts counted more than once})\} \\ &= 1 - \{\Pr(A_0) + \cdots + \Pr(A_k)\} + \Pr(\text{parts counted more than once}) \\ &\geq 1 - \{\Pr(A_0) + \cdots + \Pr(A_k)\} = 1 - (k + 1)\gamma. \end{aligned}$$

Therefore, in order to achieve (8), we require that

$$\gamma = \frac{\alpha}{k + 1} = \frac{\alpha}{\# \text{ parameters}}.$$

That is to say, using the Bonferroni adjustment, the rectangular region formed by the Cartesian product of $k + 1$ intervals, each of them is a $100(1 - \frac{\alpha}{k+1})\%$ confidence interval of an individual parameter β_ℓ :

$$\hat{\beta}_\ell \pm t_{n-k-1; \frac{\alpha}{2(k+1)}} \cdot \hat{\sigma} \cdot \sqrt{c_{\ell+1, \ell+1}}, \quad \ell = 0, 1, \dots, k,$$

where $c_{\ell+1, \ell+1}$ is the $(\ell + 1)^{\text{st}}$ diagonal element of $(\mathbf{X}'\mathbf{X})^{-1}$, is a joint confidence region for the $k + 1$ parameters in the vector β , the confidence level of this rectangle is at least $100(1 - \alpha)\%$.

This is more popular than the exact $100(1 - \alpha)\%$ confidence ellipsoid because a hyper-rectangle is easier for understanding and interpretation for individual parameters than an ellipsoid.

Of course the Bonferroni inequality can also be applied to multiple testing problems so that if we want to perform, say, individual t -test for each of m parameters, in order to have the overall significance level not higher than α , each t -test should use α/m as its significance level.

Note that whilst using Bonferroni adjustment we can control the overall significance of multiple t -tests to be bounded above by α , but using ANOVA we can control the significance level of the F -test at exactly α .

Higher-order models

Suppose we have two independent variables x_1 and x_2 . The model containing only the *first-order terms* is

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon.$$

Under this model, for each fixed x_2 , the relationship between x_1 and Y is a straight line with slope β_1 and intercept $\beta_0 + \beta_2 x_2$.

That is to say, for different values of x_2 , we get a set of parallel lines with different intercepts; the value of x_2 does not affect the relationship (i.e. the slope) between x_1 and Y . In such a situation we say that x_1 and x_2 have no *interaction*.

However, in real applications, it is common that the independent variables interact, i.e. the relationship between x_1 and Y depends on the value of x_2 .

Imagine a situation in which for $x_2 = 10$, the straight line relationship between x_1 and Y have a positive slope, whilst that for $x_2 = 50$, x_1 and Y have a negative slope; we have nonparallel lines.

In such a situation, we say x_1 and x_2 have interaction and the model that allows interaction between two independent variables will include the cross-product term of these two independent variables.

In particular, if we have only two independent variables, then the model

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + \varepsilon \quad (9)$$

is an interaction model, where $\beta_1 + \beta_3 x_2$ is the change in $\mathbb{E}(Y)$ for every 1-unit increase in x_1 , holding x_2 fixed, and $\beta_2 + \beta_3 x_1$ is the change in $\mathbb{E}(Y)$ for every 1-unit increase in x_2 , holding x_1 fixed.

The **interaction term** $x_1 x_2$ is a **second-order** term, whilst x_1 and x_2 , called the **main effects**, are **first-order**.

When we use an interaction model, extreme care is needed in interpreting the signs and sizes of coefficients. A negative β_1 does not necessarily mean that an increase in x_1 is accompanied by a decrease in the mean of Y .

The most important parameter is the one associated with the highest order term in the model. Once an interaction is significant, do NOT test on the first-order terms of the interacting variables.

In models with more than two independent variables, there may be several higher order terms (e.g. third- or fourth-order interactions). The same principle applies to such cases: do NOT test on the lower order terms of the variables contained in the existing higher order terms in the model.

The presence of interaction implies that the lower order terms are important regardless of their p -values shown on the computer printout.

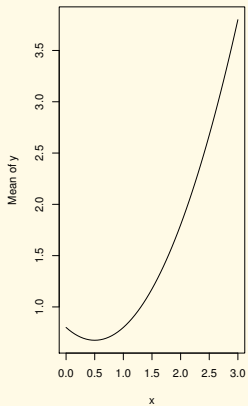
However, note that it is not a restriction required by mathematics; this is a principle that makes practical interpretation of the model easier and neater. The mathematics would not go wrong even if we excluded lower order terms but kept their higher order terms.

The model given in (9) is a second-order model. If we have only one independent variable, we still can form a second-order model, namely the **quadratic model**

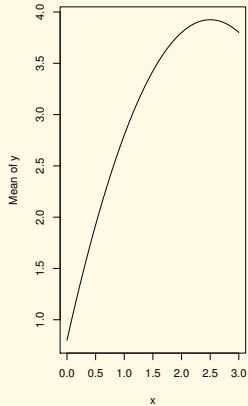
$$Y = \beta_0 + \beta_1 x + \beta_2 x^2 + \varepsilon,$$

where the parameter β_2 controls the rate of curvature so that if $\beta_2 > 0$, the curve is concave upward, whilst if $\beta_2 < 0$, the curve is concave downward.

$\beta_2 > 0$



$\beta_2 < 0$



In a quadratic model, model interpretations are not meaningful outside the range of the independent variables; in particular, $\hat{\beta}_0$ can be meaningfully interpreted only if the range of the predictor includes $x = 0$.

The parameter β_1 , in general, will not have a meaningful interpretation at all, and so rather than interpreting the numerical value of the estimate $\hat{\beta}_1$ itself, we utilize a graphical representation of the model to describe the model, which can be done in the **SAS** procedure `gplot`:

```
proc gplot;  
symbol v=plus i=rq;  
plot y*x;  
run;
```

in which `i=rq` means that the interpolation method is regression and the model is quadratic. An example of the output is given in Figure 3, which is the fitted model for Example 4.7.

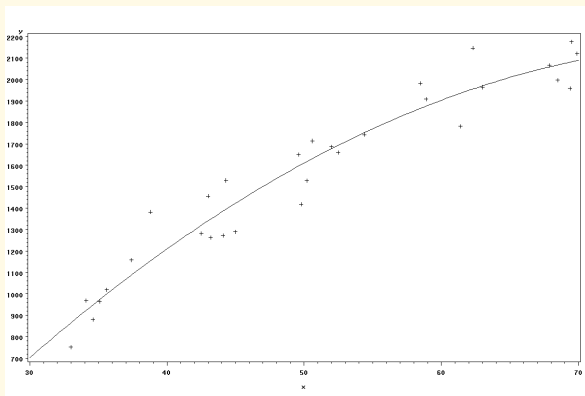


Figure 3: The scatterplot of the dataset aerobic in Example 4.7 on pages 203–207, with the fitted quadratic model.

If we have two independent variables, the *complete second-order model* is

$$Y = \beta_0 + \beta_1x_1 + \beta_2x_2 + \beta_3x_1x_2 + \beta_4x_1^2 + \beta_5x_2^2 + \varepsilon.$$

Since most relationships in the real world are curvilinear, at least to some extent, and so a good first choice would be the complete second-order model. [A curvilinear regression typically refers to a regression model containing quadratic, cubic, quartic, or higher order terms of one or more independent variables.]

If we have prior information that there is no curvature, then we may let $\beta_4 = \beta_5 = 0$, leading to the interaction model given in (9). If there is no interaction, then we may let $\beta_3 = 0$. If we do not have any prior information, we may still test whether a reduced model is as good as the complete model by the partial F -test. Which variables should be kept in the model and which should be removed from the model is the problem of *variable selection*, which will be discussed in details in Chapter 6.

Qualitative (categorical) independent variables

Up to here, we have considered only quantitative independent variables. Suppose we have qualitative (categorical) independent variables, then we have to introduce *dummy variables*. The different values of a qualitative independent variable are called its *levels*. For a two-level independent variable (say, level A and level B), we need one dummy variable:

$$x = \begin{cases} 1, & \text{if level A,} \\ 0, & \text{if level B,} \end{cases}$$

and write

$$Y = \beta_0 + \beta_1 x + \varepsilon,$$

so that

at level A: mean of the response is $\beta_0 + \beta_1$,

at level B: mean of the response is β_0 .

In general, if we have an m -level qualitative independent variable, we need $m - 1$ dummy variables

$$x_i = \begin{cases} 1, & \text{for the } i^{\text{th}} \text{ level,} \\ 0, & \text{otherwise,} \end{cases}$$

for $i = 1, \dots, m - 1$, and write

$$Y = \beta_0 + \beta_1 x_1 + \dots + \beta_{m-1} x_{m-1} + \varepsilon,$$

so that the mean responses are

at the i^{th} level: mean of the response is $\beta_0 + \beta_i$, $i = 1, \dots, m - 1$,

at the m^{th} level: mean of the response is β_0 .

But, wait, why don't we define a variable, say,

$$x = \begin{cases} 0, & \text{for the first level,} \\ 1, & \text{for the second level,} \\ 2, & \text{for the third level,} \\ \vdots & \vdots \\ m - 1 & \text{for the } m\text{th level,} \end{cases}$$

instead of getting into the trouble of introducing $m - 1$ dummy variables?

The answer is that if we do so, then we actually assume that changing x from level 1 to level 2, or changing x from level 4 to level 5, the change in the mean of Y is the same; this is because the change in the value of x is 1 in either case.

In other words, at each level of this categorical variable we get, say, a straight line and if we have no interaction, then such an artificial x will give us parallel lines that are equally apart, i.e. the vertical distance between two consecutive lines is always the same value (equal to the coefficient in front of x). This is not desirable unless there is a good reason to impose such a strong restriction by scientific theory.

To allow parallel lines that are not equally apart, each level should contribute a different β_i to the intercept, and this has to be done by using dummy variables. If using the artificial x , it means we require artificially $\beta_2 = 2\beta_1$ and $\beta_3 = \beta_2 + \beta_1 = 3\beta_1$, and so on.

Let us consider $m - 1$ dummy variables x_1, \dots, x_{m-1} for an m -level categorical independent variable, i.e.

$$Y = \beta_0 + \beta_1 x_1 + \dots + \beta_{m-1} x_{m-1} + \varepsilon. \quad (10)$$

Testing the null hypothesis that $\beta_1 = \dots = \beta_{m-1} = 0$ by the F -test in the ANOVA table of the Regression Analysis studied in this course is equivalent to testing the null hypothesis that all m means are the same $\mu_1 = \dots = \mu_m = 0$, where μ_i is the mean response at level i . The latter is the ANOVA that we have learnt in MATH2206 Prob. Stat. These two ANOVA are of course the same.

A neater expression for the model in (10) is:

$$Y_i = \beta_0 + \beta_i + \varepsilon, \quad i = 1, \dots, m,$$

where Y_i denotes the response at level i and $\beta_m = 0$ (so that the mean response at level m is just the intercept term β_0), or equivalently,

$$Y_i = \mu + \alpha_i + \varepsilon, \quad i = 1, \dots, m,$$

where μ is the overall mean and $\sum \alpha_i = 0$. The parameter α_i is interpreted as the effect of the level i on the mean, and so the mean response at level i will be the overall mean μ plus the effect α_i (some are positive and some are negative), i.e.

$$\mu_i = \mu + \alpha_i,$$

which immediately leads to the constraint we wrote above:

$$\mu = \frac{\sum \mu_i}{m} = \frac{\sum (\mu + \alpha_i)}{m} = \mu + \frac{\sum \alpha_i}{m} \Rightarrow \sum \alpha_i = 0.$$

This formulation is commonly used for ANOVA, and this will be the model for the simplest case in experimental design (MATH3815).