

REGRESSION MODELLING AND LEAST-SQUARES

The aim of **regression modelling** is to explain the observed variation in a **response** variable, Y , by relating it to some collection of **predictor** variables, X_1, X_2, \dots, X_D . The variation is decomposed into a **systematic** component, that is, a deterministic function of the predictor variables, and a **random** component that is the result of measurement procedures (or unmeasurable variability). The simplest model is the **Normal Linear Model**, where the systematic component is a function of the predictors and some model parameters, β , and the random variation is assumed to be the result of additive normally distributed random error terms. This model is explained in section 1.

1 THE NORMAL LINEAR MODEL

We assume that the variables to be modelled are as follows; we will observe paired data, with response data y_i paired to predictor variables stored in vector form $x_i = (x_{i1}, \dots, x_{iD})^\top$, and our aim is to explain the variation in (y_1, \dots, y_n) . We achieve this by modelling the conditional distribution of response variable Y_i given the observed value of predictor variable $X_i = x_i$. Specifically, we may write

$$Y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_D x_{iD} + \varepsilon_i = \beta_0 + \sum_{j=1}^D \beta_j x_{ij} + \varepsilon_i \quad (1)$$

where $\varepsilon_i \sim N(0, \sigma^2)$ for $i = 1, \dots, n$ are independent and identically distributed random error terms. Note that this implies

$$Y_i | X_i = x_i \sim N\left(\beta_0 + \sum_{j=1}^D \beta_j x_{ij}, \sigma^2\right) \quad \therefore \quad E_{f_{Y|X}} [Y_i | X_i = x_i] = \beta_0 + \sum_{j=1}^D \beta_j x_{ij}. \quad (2)$$

In vector notation, this model can be re-written $Y_i = x_i^\top \beta + \varepsilon_i$, where $x_i = (1, x_{i1}, x_{i2}, \dots, x_{iD})^\top$, and thus, for vector $Y = (Y_1, \dots, Y_n)^\top$ we have

$$Y = \mathbf{X}\beta + \varepsilon$$

where \mathbf{X} is a $n \times (D + 1)$ matrix called the **design** matrix

$$\mathbf{X} = \begin{bmatrix} 1 & x_{11} & \cdots & x_{1D} \\ 1 & x_{21} & \cdots & x_{2D} \\ 1 & x_{31} & \cdots & x_{3D} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n1} & \cdots & x_{nD} \end{bmatrix}$$

and to mimic (2)

$$Y \sim N_n(\mathbf{X}\beta, \sigma^2 I_n) \quad (3)$$

where I_n is the $n \times n$ identity matrix, giving a joint pdf for Y given \mathbf{X} of the form

$$f_Y(y; \beta, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left\{-\frac{1}{2\sigma^2}(y - \mathbf{X}\beta)^\top (y - \mathbf{X}\beta)\right\} \quad (4)$$

This joint probability model will form the basis of inference.

NOTE: we will also see that a model-free estimation procedure can be constructed on the basis of a **goodness of fit** criterion.

1.1 THE EXTENDED LINEAR MODEL

The formulation of the linear model above can be extended to allow for more general dependence on the predictors. Suppose that g_1, g_2, \dots, g_K are K (potentially non-linear) functions of the D original predictors, that is

$$g_k(x_i) = g_k(x_{i1}, \dots, x_{iD})$$

is some scalar function, for example, we could have

- $g_k(x_{i1}, \dots, x_{iD}) = g_k(x_{i1}) = x_{i1}$ (the identity function)
- $g_k(x_{i1}, \dots, x_{iD}) = g_k(x_{i1}) = a_k \sqrt{x_{i1}}$
- $g_k(x_{i1}, \dots, x_{iD}) = g_k(x_{i1}) = a_k \log x_{i1}$
- $g_k(x_{i1}, \dots, x_{iD}) = g_k(x_{i1}, x_{i2}) = a_k x_{i1} + b_k x_{i2}$

and so on. This reformulation does not effect our probabilistic definition of the model in (3); we can simply redefine design matrix \mathbf{X} as

$$\mathbf{X} = \begin{bmatrix} 1 & g_1(x_1) & \cdots & g_K(x_1) \\ 1 & g_1(x_2) & \cdots & g_K(x_2) \\ 1 & g_1(x_3) & \cdots & g_K(x_3) \\ \vdots & \vdots & \vdots & \vdots \\ 1 & g_1(x_n) & \cdots & g_K(x_n) \end{bmatrix}$$

now an $n \times (K + 1)$ matrix. In the discussion below, we will regard the **transformed** variables $(g_1(X), g_2(X), \dots, g_K(X))$ as the predictors and drop the dependence on the transformation functions. Hence we have

- Y as a $n \times 1$ column vector
- \mathbf{X} as a $n \times (K + 1)$ matrix with i th row $(1, g_1(x_i), \dots, g_K(x_i))$
- β as a $(K + 1) \times 1$ column vector

1.2 LEAST SQUARES ESTIMATION IN THE LINEAR MODEL

Equation (4) illustrates a way that parameter estimates can be obtained. For any parameter vector β , the fitted value is $\mathbf{X}\beta$, and thus the term in the exponent

$$S(\beta) = (y - \mathbf{X}\beta)^\top (y - \mathbf{X}\beta)$$

is a measure of goodness of fit; if $S(\beta)$ is small, then y and the fitted values are closely located. It can be shown (below) that the minimum value of $S(\beta)$ is obtained when

$$\mathbf{X}^\top \mathbf{X} \beta = \mathbf{X}^\top y.$$

yielding the **Ordinary Least Squares** solution

$$\tilde{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top y.$$

However, without further distributional assumptions, we cannot proceed to understand the uncertainty in the estimation.

1.3 MAXIMUM LIKELIHOOD ESTIMATION IN THE LINEAR MODEL

Maximum likelihood estimation for the normal linear model is straightforward. If $\theta = (\beta, \sigma^2)$ then the mle $\hat{\theta}$ is given by

$$\hat{\theta}_{ML} = \arg \max_{\theta \in \Theta} f_{Y|\beta, \sigma^2}(y; \beta, \sigma^2) = \arg \max_{\theta \in \Theta} L(\beta, \sigma^2; y, x)$$

where parameter space $\Theta \equiv \mathbb{R}^K \times \mathbb{R}^+$. Taking logs in (4) gives

$$\log L(\beta, \sigma^2; y, x) = -\frac{n}{2} \log \sigma^2 - \frac{n}{2} \log 2\pi - \frac{1}{2\sigma^2} (y - \mathbf{X}\beta)^\top (y - \mathbf{X}\beta) \quad (5)$$

and considering the maximization for β indicates

$$\arg \max_{\beta \in \mathbb{R}^K} \log L(\beta, \sigma^2; y, x) = \arg \min_{\beta \in \mathbb{R}^K} (y - \mathbf{X}\beta)^\top (y - \mathbf{X}\beta)$$

and thus,

$$\begin{aligned} S(\beta) &= (y - \mathbf{X}\beta)^\top (y - \mathbf{X}\beta) \\ &= y^\top y - y^\top \mathbf{X}\beta - \beta^\top \mathbf{X}^\top y + \beta^\top \mathbf{X}^\top \mathbf{X}\beta \\ &= y^\top y - 2y^\top \mathbf{X}\beta + \beta^\top \mathbf{X}^\top \mathbf{X}\beta. \end{aligned}$$

Using vector/matrix differentiation

$$\frac{d}{d\beta} \{y^\top \mathbf{X}\beta\} = y^\top \mathbf{X} \quad \frac{d}{d\beta} \{\beta^\top \mathbf{X}^\top \mathbf{X}\beta\} = 2\mathbf{X}^\top \mathbf{X}\beta \quad (6)$$

and so if $\hat{\beta}$ is the solution of

$$\frac{dS(\beta)}{d\beta} = -y^\top \mathbf{X} + \mathbf{X}^\top \mathbf{X}\beta = 0$$

then it follows that $\hat{\beta}$ satisfies

$$\mathbf{X}^\top \mathbf{X}\hat{\beta} = \mathbf{X}^\top y. \quad (7)$$

If the matrix $\mathbf{X}^\top \mathbf{X}$ is non-singular, then we have the ML estimates of β as

$$\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top y \quad (8)$$

and substituting back into (5) gives

$$\hat{\sigma}^2 = \frac{1}{n} (y - \mathbf{X}\hat{\beta})^\top (y - \mathbf{X}\hat{\beta}) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (9)$$

where $\hat{y}_i = x_i^\top \hat{\beta}$ is the **fitted value**, and $y_i - \hat{y}_i$ is the **residual**. Note that $\mathbf{X}^\top \mathbf{X}$ is a symmetric matrix. The expression $(y - \mathbf{X}\hat{\beta})^\top (y - \mathbf{X}\hat{\beta})$ is termed the **residual sum of squares** (or **RSS**). A common **adjusted** estimate is

$$\hat{\sigma}_{ADJ}^2 = \frac{1}{n - K - 1} (y - \mathbf{X}\hat{\beta})^\top (y - \mathbf{X}\hat{\beta}) \quad (10)$$

the justification for this result depends on the sampling distribution of the estimator.

If $K = 1$, with identity function $g(t) = t$

$$\mathbf{X}^T \mathbf{X} = \begin{bmatrix} n & \sum_{i=1}^n x_i \\ \sum_{i=1}^n x_i & \sum_{i=1}^n x_i^2 \end{bmatrix} \quad (\mathbf{X}^T \mathbf{X})^{-1} = \frac{1}{n^2 S_x^2} \begin{bmatrix} \sum_{i=1}^n x_i^2 & -\sum_{i=1}^n x_i \\ -\sum_{i=1}^n x_i & n \end{bmatrix}$$

where

$$S_x^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - (\bar{x})^2$$

and so

$$\begin{aligned} \hat{\beta} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T y = \frac{1}{n^2 S_x^2} \begin{bmatrix} \sum_{i=1}^n x_i^2 & -\sum_{i=1}^n x_i \\ -\sum_{i=1}^n x_i & n \end{bmatrix} \begin{bmatrix} \sum_{i=1}^n y_i \\ \sum_{i=1}^n x_i y_i \end{bmatrix} \\ &= \frac{1}{n^2 S_x^2} \begin{bmatrix} \sum_{i=1}^n x_i^2 \sum_{i=1}^n y_i - \sum_{i=1}^n x_i \sum_{i=1}^n x_i y_i \\ n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i \end{bmatrix} \end{aligned}$$

For $K > 1$ calculations can proceed in this fashion, but generally the matrix form

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

is easier to work with.

NOTE: Most good computing packages (R, MATLAB) have pre-written functions that implement this form of linear model in a straightforward fashion.

1.4 FITTED VALUES

Fitted values are readily obtained from this model. The fitted value \hat{y} is obtained as

$$\begin{aligned} \hat{y} &= \mathbf{X} \hat{\beta} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T y \\ &= H y \end{aligned}$$

and the measure of misfit is

$$\begin{aligned} S(\hat{\beta}) &= (y - \hat{y})^T (y - \hat{y}) \\ &= (y - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T y)^T (y - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T y) \\ &= y^T (I_n - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T) y \end{aligned}$$

1.5 PROPERTIES OF THE ML ESTIMATORS

By elementary properties of random variables, the properties of ML estimator $T = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T Y$

$$\begin{aligned} E_{Y|X,\beta,\sigma^2} [T] &= E_{Y|X,\beta,\sigma^2} \left[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T Y \right] = ((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T) E_{Y|X,\beta,\sigma^2} [Y] \\ &= ((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T) \mathbf{X} \beta = (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X}) \beta = \beta \end{aligned}$$

so that T is **unbiased** for β , and

$$\begin{aligned} \text{Var}_{Y|X,\beta,\sigma^2} [T] &= \text{Var}_{Y|X,\beta,\sigma^2} \left[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T Y \right] \\ &= ((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T) \text{Var}_{Y|X,\beta,\sigma^2} [Y] ((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T)^T \\ &= ((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T) \sigma^2 I_n (\mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}) \\ &= \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X}) (\mathbf{X}^T \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}. \end{aligned}$$

Note that, in fact, given β and σ^2

$$Y \sim N_n(\mathbf{X}\beta, \sigma^2 I_n) \implies T = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T Y \sim N_{K+1}(\beta, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}). \quad (11)$$

It also follows that

$$(y - \mathbf{X}\beta)^T (y - \mathbf{X}\beta) = (y - \mathbf{X}\hat{\beta})^T (y - \mathbf{X}\hat{\beta}) + (\hat{\beta} - \beta)^T (\mathbf{X}^T \mathbf{X}) (\hat{\beta} - \beta)$$

or

$$S(\beta) = S(\hat{\beta}) + (\hat{\beta} - \beta)^T (\mathbf{X}^T \mathbf{X}) (\hat{\beta} - \beta)$$

where

$$S(\beta) = (y - \mathbf{X}\beta)^T (y - \mathbf{X}\beta) \quad (1)$$

$$S(\hat{\beta}) = (y - \mathbf{X}\hat{\beta})^T (y - \mathbf{X}\hat{\beta}) = (y - \hat{y})^T (y - \hat{y}) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (2)$$

$$(\hat{\beta} - \beta)^T (\mathbf{X}^T \mathbf{X}) (\hat{\beta} - \beta) = (\mathbf{X}\hat{\beta} - \mathbf{X}\beta)^T (\mathbf{X}\hat{\beta} - \mathbf{X}\beta) \quad (3)$$

are the **(1) TOTAL**, **(2) RESIDUAL** and **(3) FITTED** sum of squares (**TSS**, **RSS** and **FSS**). Therefore, by normal distribution theory, it follows that

$$\frac{S(\beta)}{\sigma^2} \sim \chi_n^2 \quad \frac{S(\hat{\beta})}{\sigma^2} \sim \chi_{n-K-1}^2 \quad (12)$$

so that

$$s^2 = \frac{S(\hat{\beta})}{(n - K - 1)} \text{ is an } \mathbf{UNBIASED} \text{ estimator of } \sigma^2$$

and the quantity

$$\frac{\hat{\beta} - \beta}{s.e.(\hat{\beta})} = \frac{\hat{\beta} - \beta}{s\sqrt{v_{ii}}} \sim \text{Student}(n - K - 1).$$

It also follows that

$$\frac{S(\beta) - S(\hat{\beta})}{\sigma^2} = \frac{(\hat{\beta} - \beta)^T (\mathbf{X}^T \mathbf{X}) (\hat{\beta} - \beta)}{\sigma^2} \sim \chi_{K+1}^2$$

so that finally

$$\frac{[S(\beta) - S(\hat{\beta})] / (K + 1)}{S(\hat{\beta}) / (n - K - 1)} \sim \text{Fisher}(K + 1, n - K - 1)$$

It follows that in this case the ML estimator is the Minimum Variance Unbiased Estimator (MVUE) and the Best Linear Unbiased Estimator (BLUE).

1.6 THE ANALYSIS OF VARIANCE

Analysis of variance or **ANOVA** is used to display the sources of variability in a collection of data samples. The ANOVA F-test compares variability **between** samples with the variability **within** samples. In the above analysis, we have that

$$S(\beta) = S(\hat{\beta}) + (\hat{\beta} - \beta)^T (\mathbf{X}^T \mathbf{X}) (\hat{\beta} - \beta) \quad \text{or} \quad TSS = RSS + FSS.$$

Now, using the distributional results above, we can construct the following **ANOVA Table** to test the hypothesis

$$H_0 : \beta_1 = \dots = \beta_K = 0$$

against the general alternative that H_0 is not true.

Source of Variation	D.F.	Sum of squares	Mean square	F
FITTED	K	FSS	FSS/K	$\frac{FSS/K}{RSS/(n - K - 1)}$
RESIDUAL	$n - K - 1$	RSS	$RSS/(n - K - 1)$	
TOTAL	$n - 1$	TSS		

This test allows a comparison of the fits of the two competing models implied by the null and alternative hypotheses. Under the null model, if H_0 is true, then the model has $Y_i \sim N(\beta_0, \sigma_0^2)$ for $i = 1, 2, \dots, n$, for some β_0 and σ_0^2 to be estimated. Under the alternative hypothesis, there are a total of $K + 1$ β parameters to be estimated using equation (8). The **degrees of freedom** column headed (D.F.) details how many parameters are used to describe the amount of variation in the corresponding row of the table; for example, for the FIT row, D.F. equals K as there are K parameters used to extend the null model to the alternative model.