

Choosing the best model for prediction

Calculating the PRESS residuals $\delta_i^* = \hat{y}_i - y_{i,-i}$ is an example of doing leave one out cross validation. These residuals measure how well the i -th observation of the response can be estimated or predicted in terms of the other observations. The model with the best prediction ability should be the one that minimizes $\sum_{i=1}^m |\delta_i^*|$ or $\sum_{i=1}^m (\delta_i^*)^2 = \text{PRESS}$, so eventually the one that maximizes $R_{\text{pred}}^2 = 1 - \frac{\sum_{i=1}^m (\delta_i^*)^2}{SST}$

Categorical or Indicator variables

Suppose $y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 z_i + \epsilon_i$ where z is a categorical variable with l categories. Such regression problems are normally best solved by introducing $l-1$ indicator variables. For instance, with three different catalyst the regression model can be:

$$\cancel{y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 z_{i1} + \beta_4 z_{i2} + \epsilon_i, \quad i=1, 2, \dots, m}$$

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 z_{i1} + \beta_4 z_{i2} + \epsilon_i, \quad i=1, 2, \dots, m$$

Where the columns Z_1 and Z_2 are

Z_1	Z_2
1	0
1	0
1	0
0	1
0	1
0	1
0	1
0	1
0	1

The $(1,0)$ combination correspond to analyst 1, the $(0,1)$ to catalyst 2 and the $(0,0)$ to catalyst 3.

The model is then,

$$Y_i = (\beta_0 + \beta_3) + \beta_1 X_{i1} + \beta_2 X_{i2} + \epsilon_i \text{ for catalyst 1}$$

$$Y_i = (\beta_0 + \beta_4) + \beta_1 X_{i1} + \beta_2 X_{i2} + \epsilon_i \quad -u - 2$$

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2}$$

It may happen that the coefficients for the continuous variables depend on the ~~level of the~~ category variables.

Suppose we construct the model

~~$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \beta_3 Z_{i1} + \beta_4 Z_{i2} + \beta_5 X_{i1} Z_{i1} + \beta_6 X_{i2} Z_{i2}$$~~

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \beta_3 Z_{i1} + \beta_4 Z_{i2} + \beta_5 X_{i1} Z_{i1} + \beta_6 X_{i2} Z_{i2} + \epsilon_i$$

$$i = 1, 2, \dots, n$$

Then we have:

$$Y_i = (\beta_0 + \beta_3) + (\beta_1 + \beta_5) X_{i1} + \beta_2 X_{i2} + \epsilon_i \text{ catalyst 1}$$

$$Y_i = (\beta_0 + \beta_4) + \beta_1 X_{i1} + (\beta_2 + \beta_6) X_{i2} + \epsilon_i \text{ catalyst 2}$$

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \epsilon_i \text{ catalyst 3.}$$

Orthogonal columns in the Design matrix

Let the design matrix $\underline{X} = [1, \underline{x}_1, \underline{x}_2, \dots, \underline{x}_k]$

If $\underline{x}_p^T \underline{x}_q = 0$ i.e. $\sum_{i=1}^n x_{ip} x_{iq} = 0$, $p \neq q$ \underline{x}_p and \underline{x}_q are orthogonal

If all columns are orthogonal $\underline{X}^T \underline{X} =$

$$\begin{bmatrix} n & 0 & \dots & 0 \\ 0 & \underline{x}_1^T \underline{x}_1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \underline{x}_k^T \underline{x}_k \end{bmatrix}$$

and $\hat{\beta} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{Y} = \begin{bmatrix} 1/n & 0 & \dots & 0 \\ 0 & (\underline{x}_1^T \underline{x}_1)^{-1} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & (\underline{x}_k^T \underline{x}_k)^{-1} \end{bmatrix} \begin{bmatrix} \sum_i y_i \\ \underline{x}_1^T \underline{Y} \\ \vdots \\ \underline{x}_k^T \underline{Y} \end{bmatrix}$

The estimator for β_j will only depend on \underline{x}_j and \underline{Y}

$$\begin{aligned} \text{and } SS_R &= \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = \sum_{i=1}^n (b_0 + b_1 x_{i1} + \dots + b_k x_{ik} - \bar{y})^2 \\ &\simeq b_1^2 \sum_{i=1}^n x_{i1}^2 + b_2^2 \sum_{i=1}^n x_{i2}^2 + \dots + b_k^2 \sum_{i=1}^n x_{ik}^2 \\ &= SS_R(\beta_1 | \beta_0) + SS_R(\beta_2 | \beta_0) + \dots + SS_R(\beta_k | \beta_0) \end{aligned}$$

Testing a general linear hypothesis

Suppose $E[\underline{Y}] = \beta_0 + \beta_1 \underline{x}_1 + \beta_2 \underline{x}_2 + \beta_3 \underline{x}_3$ and we want to test,

$$H_0: \begin{bmatrix} \beta_1 - \beta_2 \\ \beta_2 - \beta_3 \end{bmatrix} = \underline{0} \quad \text{or} \quad H_0: \beta_1 = \beta_2 = \beta_3 = \beta.$$

Can be formulated as: $H_0: C \underline{\beta} = \underline{0} \quad H_1: C \underline{\beta} \neq \underline{0}$

where $C = \begin{bmatrix} 0 & 1 & -1 & 0 \\ 0 & 0 & 1 & -1 \end{bmatrix}$ and $\underline{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix}$, $C \underline{\beta} = \underline{0} \Leftrightarrow \begin{bmatrix} \beta_1 - \beta_2 \\ \beta_2 - \beta_3 \end{bmatrix} = \underline{0}$

We have n (in this case 2) constraints on $\underline{\beta}$.

$$\hat{\beta} \sim N(\beta, \sigma^2 (\underline{x}^T \underline{x})^{-1}) \Rightarrow \underline{C}\hat{\beta} \sim N(\underline{C}\beta, \sigma^2 \underline{C} (\underline{x}^T \underline{x})^{-1} \underline{C}^T)$$

$$\text{Hence under } H_0 \quad (\underline{C}\hat{\beta})^T (\sigma^2 \underline{C} (\underline{x}^T \underline{x})^{-1} \underline{C}^T)^{-1} (\underline{C}\hat{\beta}) \sim \chi^2_n$$

Since $\hat{\beta}$ and SSE (full model) are independent and $\frac{SSE}{\sigma^2} \sim \chi^2_{(n-k-1)}$

we get.

$$F = \frac{(\underline{C}\hat{\beta})^T (\underline{C} (\underline{x}^T \underline{x})^{-1} \underline{C}^T)^{-1} (\underline{C}\hat{\beta}) / \sigma^2 n}{SSE / \sigma^2 (n-k-1)} \sim F_{k, n-k-1}$$

Under H_0 . Reject if F is large.

$$\text{Testing } H_0: C\beta = d \quad H_1: C\beta \neq d$$

we just substitute $C\hat{\beta}$ with $C\hat{\beta} - d$ in F