



iJRASET

International Journal For Research in
Applied Science and Engineering Technology



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 6 Issue: IV Month of publication: April 2018

DOI: <http://doi.org/10.22214/ijraset.2018.4199>

www.ijraset.com

Call:  08813907089

E-mail ID: ijraset@gmail.com

Effect of Additional Variables on Regression Diagnostics

Disha K. Rank¹

¹Department of Statistics, Saurashtra University, Rajkot

Abstract: We obtain the effect of adding concomitant variables in a linear regression model on estimators of parametric functions of interest and on some commonly used regression diagnostics.

Keywords: Regression model, Regression diagnostics, linear function, DFFIT, DFBETA.

I. INTRODUCTION

Regression diagnostics and detection of influential observations play an important role in statistical modelling (see for example Belsley, Kuh, Welsch (1980), Chatterjee and Haldi (1988)). In this paper, we consider the usual linear regression model $(Y, X\beta, \sigma^2 I)$ where $X'X$ could be singular and obtain the effect of adding more independent variables to this model on some of the most commonly used regression diagnostics DFBETA and DFFIT. This is important in usual regression set up as one might like to look into this aspect also while doing a step-wise regression. Also the stability of the regression diagnostics as the regressors increase would help in identifying influential observations in a more meaningful manner. The results of this paper can be interpreted in the design set up also. Consider the original model as the ANOVA model and the latter model as the analysis of covariance model. We are usually interested in inferences on linear functions of a subset of β parameters (for example treatment contrasts). The present study enables us to assess the influence of an observation on the inferences concerning the desired parametric functions in the light of concomitant variables. A detailed study in this regard is underway and will be reported elsewhere. The following notations will be used in the paper. For a matrix A : $\rho(A)$, $C(A)$, $R(A)$ and A' denote the rank, column space, row space, and transpose, respectively. A^- denotes a g-inverse of A . For a matrix B , b_{i*} and b_{*j} denote the i^{th} row and the j^{th} column of B , respectively, and $B_{(i)}$ denotes the matrix obtained from B by dropping the i^{th} row. P_B denotes the orthogonal projector projecting into $C(B)$. (The inner product considered is the Euclidean inner product.)

II. MAIN RESULTS

We consider four models

$$a) (Y, X\beta, \sigma^2 I) \tag{2.1}$$

$$b) (Y_{(i)}, X_{(i)}\beta, \sigma^2 I) \tag{2.2}$$

$$c) (Y, X\beta + C\gamma, \sigma^2 I) \tag{2.3}$$

$$d) (Y_{(i)}, X_{(i)}\beta + C_{(i)}\gamma, \sigma^2 I) \tag{2.4}$$

where X and C are fixed matrices of orders $n \times m$ and $n \times k$, respectively. In all the above models $X'X$ could be singular. Let $\hat{\beta}^{(j)}$, $j = 1, 2$ and $\begin{pmatrix} \hat{\beta}^{(j)} \\ \hat{\beta}^{(j)} \end{pmatrix}$, $j = 3, 4$, satisfy the normal equations and $R_0^{2(j)}$, $j = 1, \dots, 4$, denote the residual sum of squares for the models (2.1), (2.2), (2.3) and (2.4) respectively.

We start with the definitions of a few regression diagnostics for the model (2.1). $(\hat{\beta}^{(1)} - \hat{\beta}^{(2)})$ is called DFBETA. Let $\lambda'\beta$ be estimable under models (2.1) and (2.2) then

DFBETA _{λ} is the quantity $\lambda'(\hat{\beta}^{(1)} - \hat{\beta}^{(2)})$, namely the difference in BLUEs of $\lambda'\beta$ under the models (2.1) and (2.2). DFBETA _{λ, x'_j} is called DFFIT _{ij} . The matrix $H = D(X\beta^{(1)})/\sigma^2 = X(X'X)^-X' = P_X$ is called the Hat matrix and its diagonal elements are called the hat matrix diagonals (or leverage values).

A. *Lemma 2.1.* Let A be a symmetric matrix and A^- be a g-inverse of A . Let $u \in C(A)$. Then for any nonzero real number k ,

$$(A + kuu') = \begin{cases} A^- & \text{if } u^{A^-}u = \frac{1}{k} \\ A^- - \frac{A^-uu'A^-}{\left(\frac{1}{k}\right) + u'A^-u} & \text{otherwise.} \end{cases}$$

This lemma is well known. (See, for example, Rao and Mitra (1971) p. 40).

B. Lemma 2.2. (Rohde (1965), Bhimasankaram (1971)). Let Σ be a nonnegative definite (nnd) matrix. Partition Σ as

$$\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$$

Where Σ_{11} and Σ_{22} are square matrices. let Σ_{11}^- be a g - inverse of Σ_{11} . Denote $F = \Sigma_{22} - \Sigma_{21}\Sigma_{11}^-\Sigma_{12}$. Then one choice of a g-inverse of Σ is

$$\begin{pmatrix} \Sigma_{11}^- & 0 \\ 0 & 0 \end{pmatrix} + \begin{pmatrix} \Sigma_{11}^- & \Sigma_{12} \\ & -I \end{pmatrix} F^- \begin{pmatrix} \Sigma_{21} & \Sigma_{11}^- & -I \end{pmatrix}$$

where F^- is any g-inverse of F .

C. Lemma 2.3. Consider the model (2.1). Let $x_{i^*} \notin R(X(i))$. Then the following hold :

- 1) The error space is the same for the models (2.1) and (2.2).
- 2) For all $p \in C(X'(i))$ the BLUES of $p'\beta$ are the same in the models (2.1) and (2.2) and $x_i'\hat{\beta} = Y_i$
- 3) The usual estimator of σ^2 is the same under the models (2.1) and (2.2).

Proof : (a) $l'Y$ belongs to the error space $\Rightarrow l'X = 0$

$l'_{(i)}X_{(i)} + l_i x_{i^*} = 0 \Rightarrow l_i = 0$ since, $x_{i^*} \notin R(X'(i))$. $\Rightarrow l'Y = l'_{(i)}Y_{(i)}$. Thus error space of the model (2.1) is contained in the error space of the model (2.2). Thus other inclusion is trivial. Proofs of (b) and (c) are easy in view of (a) and are omitted. (These are also available in Bhimasankaram and Jammalamadaka (1990), henceforth denoted as B & J 1990)).

D. Lemma 2.4 (B&J (1990)). Consider the model (2.1). Let $x_{i^*} \in R(X(i))$ and let

$$e_i = y_i - x_{i^*}'\hat{\beta}^{(1)}. \text{ Then the following hold :}$$

$$(a) \hat{\beta}^{(2)} = \hat{\beta}^{(1)} - \frac{(X'X)^- x_{i^*}' e_i}{1-h_i}$$

$$(b) (R_0^2)^2 = (R_0^{(1)})^2 - \frac{e_i^2}{1-h_i}$$

$$\text{where } h_i = x_{i^*}'(X'X)x_{i^*}.$$

E. Lemma 2.5. $P_{(xx)} = P_x + P_{(I-P_x)C}$.

The proof when $(X : C)$ is of full column rank is given in Chatterjee and Hadi (1988). The proof when $(X : C)$ is not of full rank follows from lemma (2.2) along the same lines.

F. Lemma 2.6. Consider the models (2.1) and (2.3). Let θ_r denote a solution of the normal equations for the Model $(c_{*p}, X\theta, \sigma^2I)$. Let $s_{pq} = (c_{*p} - X\hat{\theta}_p)'(c_{*q} - X\hat{\theta}_q)$, $p, q = 1, \dots, k$ and $s_{p0} = (c_{*p} - X\hat{\theta}_p)'(Y - X\hat{\beta}^{(1)})$. Write $T = (\hat{\theta}_1, \dots, \hat{\theta}_k)$, $S_k = ((Spq))$ and $S_0 = (S_{10}, \dots, S_{k0})'$. Then $(\hat{\beta}^{(3)}, \hat{\gamma}^{(3)})'$ can be obtained as follows : $\hat{\gamma}^{(3)}$ is a solution of $s_{\gamma} = S_0$ and $\hat{\beta}^{(3)} = \hat{\beta}^{(1)} - T\hat{\gamma}^{(3)}$. $(R_0^{(3)})^2 = (R_0^{(1)})^2 - S_0'\hat{\gamma}^{(3)}$.

Proof. If S is singular then the system of equations $s_{\gamma} = S_0$ will be consistent because $S = [C'(I - P_x)] [C'(I - P_x)]'$ and $S_0 = C'(I - P_x)Y$ and the normal equations for the model (2.3) are

$$X'X\beta + X'C\gamma = X'Y \tag{1}$$

$$C'X\beta + C'C\gamma = C'Y \tag{2}$$

and the reduced normal equations for γ is $C'(I - P_x)Cy = C'(I - P_x)Y$. Therefore β and γ will be estimated along the same lines of Rao's method.

In view of lemma 2.3, it is clear that the case $x_{i*} \notin R(X_{(i)})$ is of no practical interest in connection with the regression diagnostics. We shall consider the case $x_{i*} \in R(X_{(i)})$ and prove

G. Theorem 2.7. Consider the models (2.1) and (2.3). Let $x_{i*} \in R(X_{(i)})$. Then the following hold :

Let $\lambda'\beta$ be estimable under both the models; then the change in $DFBETA_{i,\lambda}$

$$\text{i.e. } \lambda'((\hat{\beta}^{(3)} - \hat{\beta}^{(4)}) - ((\hat{\beta}^{(1)} - \hat{\beta}^{(2)}))) \text{ is given by } X'(W\hat{\gamma}^{(4)} - T\hat{\gamma}^{(3)}),$$

the change in the residual sum of squares is given by $V_o'\hat{\gamma}^{(4)} - s_o'\hat{\gamma}^{(3)}$ where W and V_o correspond to T and S_o respectively when the models (2.1) and (2.3) are replaced by (2.2) and (2.4), respectively, in lemma 2.6.

The proof follows by applying lemma 2.6 twice for models (2.1) and (2.3) and then for models (2.2) and (2.4).

H. Theorem 2.8. Consider the models (2.1) and (2.3). Let $x_{i*} \in R(X_{(i)})$. Then the following hold :

(a) $x_{i*}\beta$ is estimable under all the models (2.1), (2.2), (2.3) and (2.4). Then the change in

$$DFFIT_{ij} \text{ i.e. } (\hat{y}_j^{(3)} - \hat{y}_j^{(4)}) - (\hat{y}_j^{(1)} - \hat{y}_j^{(2)}) \text{ is given by}$$

$$x_{j*}(\widehat{W}\hat{\gamma}^{(4)} - \widehat{T}\hat{\gamma}^{(3)}) + \frac{c_{*j}'s^- \eta(e_i - \eta)\hat{\gamma}^{(3)}}{1 - h_i - \eta's^- \eta}$$

Where $\eta' = (x_{i*}T - c_{i*})$ and W is the same as above.

The proof follows by straight forward computations using lemmas 2.6 and 2.4. (b)The effect on hat matrix for $(X : C)H_{(x:c)}$

$$H_{(x:c)} = (X : C)[(X : C)'(X : C)]^{-1}(X : C)'$$

Now

$$[(X : C)'(X : C)]^{-1} = \begin{pmatrix} (X'X)^{-1} & 0 \\ 0 & -I \end{pmatrix} + \begin{pmatrix} (X'X)^{-1}X'C \\ -I \end{pmatrix} F^{-1}(C'X(X'X)^{-1} - I)$$

where $F = C'(I - P_x)C$.

Hence $H(x:c) = Hx + (I - Hx)C[C'(I - Hx)C]^{-1}C'(I - Hx)$. For single additional variable discussion is given in Chatterjee and Haldi (1988).

The variance of an estimable parametric function $\lambda'\hat{\beta}^{(3)}$, assuming X and C to be non-stochastic, is $\sigma^2\lambda'[(X'X)^{-1} + T\bar{S}T']\lambda$, which will be needed to compute $DFBETAS$ and $DFFITS$ etc.

Now let us consider one additional variable c case. Then the change in $DFBETA_{i,\lambda}$ using Rao's method of estimating β and γ is

$$\frac{\delta i [e_i ((1 - h_i)\hat{\theta}^{(3)} - \delta i V) + \hat{\gamma}^{(3)}V]}{(1 - h_i)(\delta'\delta - \delta_i^2)}$$

Where $V = (X'X)^{-1}x_{i*}$ and $\delta = c - X\hat{\theta}$, and the other regression diagnostics change correspondingly.

REFERENCES

- [1] Belsley, D.A., Kuh, E. and Welsch, R.E., (1980), "Regression Diagnostics", Wiley, New York
- [2] P. Bhimasankaram, D. Sengupta and S. Ramanathan (1995). Recursive inference in the general linear model, Sankhya Series A, 57(2), 227-25
- [3] Chatterjee, S. and Hadi, A.S. (1988), "Sensitivity Analysis in Linear Regression", Wiley, New York.
- [4] Rao, C.R. (1973), "Linear Statistical Inference and its Applications", John Wiley, New York.
- [5] Rao, C.R. & Mitra, S. (1971) : Generalized inverse of matrices and its applications, John Wiley and Sons.
- [6] Rohde, C.A. (1965) : Generalized inverses of partitioned matrices. J. SIAM., 13, 1033-1035.
- [7] Bhimasankaram, P. (1971) : On generalized inverses of partitioned matrices. Sankhya, Series , 33, 311-314.
- [8] D. Sengupta and S.R. Jammalamadaka (2003). Linear Models: An Integrated Approach, World Scientific Publishing Company, River Edge, NJ,



10.22214/IJRASET



45.98



IMPACT FACTOR:
7.129



IMPACT FACTOR:
7.429



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24*7 Support on Whatsapp)