

The Methods for examining Observations in the transformed Regression

By Byung-Yub Cho & Ju-Sun Ahn

1. Introduction

In recent years there has been considerable interest in finding the outlying or the influential observations in regression. The plots of the residuals and the examination of standardized residuals or of studentized residuals have been used to find those observations.

Behnken & Draper (1972) also have illustrated the estimated variances of the residuals contain more relevant informations than that furnished by residual plots or studentized residuals.

The quantities used to find such observations have been proposed by Cook (1977) and Andrews & Pregibon (1978). By these methods and quantities we may not detect which explanatory variables can provide the greatest influence for the residuals and which observations can have the same properties clustering them.

It is the purpose of this paper to describe and exemplify the methods which detect such outliers and explanatory variables and which find clusters of the observations.

The residuals and statistics are discussed to compare the methods of this paper with that of other paper in section 2. The transformations of variables and parameters are introduced in section 3. The methods for displaying the transformed regression and for examining observations is discussed in section 4. And we illustrate the plots of the transformed regression as an example whose data have two explanatory variables in section 5.

2. Residuals and statistics

In the following sections a linear model

$$\underline{Y} = \underline{X}\underline{\beta} + \underline{e}$$

is discussed where \underline{Y} is an $n \times 1$ vector of response datas (or dependent variables), X is an $n \times p$ full rank matrix of known constants (or independent variables), $\underline{\beta}$ is a $p \times 1$ vector of unknown parameters and \underline{e} is an $n \times 1$ vector of randomly distributed errors (or noise) such that $E(\underline{e}) = \underline{0}$ and $V(\underline{e}) = I\sigma^2$.

The least squares estimate of $\underline{\beta}$ is given by

$$\underline{b} = (X'X)^{-1} X'Y$$

The corresponding residual vector is

$$\begin{aligned} \underline{R} = (R_k) &= \underline{Y} - \hat{\underline{Y}} = \underline{Y} - X\underline{b}, \quad k=1, \dots, n \\ &= [I-M]\underline{Y} = [I-M]\underline{e} \\ \text{where } M &= X(X'X)^{-1}X'. \end{aligned}$$

The expectations of \underline{Y} , \underline{b} and \underline{R} are respectively

$$\begin{aligned} E(\underline{Y}) &= X\underline{\beta} \text{ (or } E(\hat{\underline{Y}}) = X\underline{b} = \hat{\underline{Y}}), \\ E(\underline{b}) &= E[(X'X)^{-1}X'Y] = \underline{\beta}, \\ \text{and } E(\underline{R}) &= \underline{0}. \end{aligned}$$

The covariance matrices of $\hat{\underline{Y}}$, \underline{b} and \underline{R} are respectively

$$\begin{aligned} V(\hat{\underline{Y}}) &= X(X'X)^{-1}X'\sigma^2 = M\sigma^2, \\ V(\underline{b}) &= (X'X)^{-1}\sigma^2, \\ \text{and } V(\underline{R}) &= E(\underline{R}\underline{R}') = (I-M)\sigma^2. \end{aligned}$$

For the k th observation, it follows that

$$\begin{aligned} \hat{y}_k &= \underline{x}_k' \underline{b}, \\ V(\hat{y}_k) &= \underline{x}_k'(X'X)^{-1}\underline{x}_k\sigma^2 = v_k\sigma^2, \\ \text{and } V(R_k) &= (1 - v_k)\sigma^2. \end{aligned}$$

where \underline{x}_k is the vector of explanatory variables for the k th observation and $v_k = \underline{x}_k'(X'X)^{-1}\underline{x}_k$.

Since the average value of the variances of the predicted values at the observation points is given by

$$\overline{V(\hat{y})} = \sum_{k=1}^n V(\hat{y}_k)/n = \text{tr}\{X(X'X)^{-1}X'\sigma^2\}/n = \frac{p\sigma^2}{n} \text{ where } X \text{ is an } n \times p \text{ matrix, the}$$

average variance of the residuals is

$$\overline{V(\underline{R})} = \sum_{k=1}^n V(R_k)/n = \frac{(n-p)\sigma^2}{n}$$

The standardized residuals N_k and the studentized residuals T_k are defined as

$$N_k = \frac{R_k}{\{(n-p)V(R_k)\}^{1/2}} = \frac{R_k}{S(1-v_k)^{1/2}}$$

$$T_k = \frac{R_k}{\{V(R_k)\}^{1/2}} = \sqrt{(n-p)}N_k$$

where $S^2 = R'R$, the residual sum of squares and v_k is the k th diagonal element of $X(X'X)^{-1}X'[b]$.

If the variances of the residuals varied a great deal, it would be worthwhile to examine the $\frac{R_k}{(1-v_k)^{1/2}}$ instead of R_k in the usual residuals plots and to use N_k instead of $\frac{R_k}{S}$ as the "normal deviate form" of the residuals.

Behnken and Draper (1972) have recommended the studentized residuals T_k as more appropriate than the standardized residuals.

A wide variation in the $V(R_k)$ reflects a peculiarity of the X matrix, namely a nonhomogeneous spacing of the observations and will thus often direct attention to data deficiencies. From the standpoint of detection outliers, it is always important to recognize the changing magnitude of $V(R_k)$.

Even if the residual variance is reasonably constant, it will still be useful to consider the magnitude of R_k relative to $S \left\{ \frac{n-p}{n} \right\}^{1/2}$ rather than to S (the estimate of σ) in cases where n is not large relative to p .

Andrews-Pregibon statistic is defined as

$$R_{ij\dots}^k(X^*) = \frac{D_{ij\dots}^k |X^* X^*|}{|X^* X^*|}$$

where $|X^* X^*| = R'R \cdot |X'X|$ and $D_{ij\dots}^k |X^* X^*|$ are the operator which deletes the elements associated with the K observations from $|X^* X^*|$.

If the deletion of an observation has a large (small) effect on $|X'X|$, the observation has a large (small) influence on the resulting estimates. Also the more a particular observation deviates from the fitted values, the more the deletion of this observation will reduce the residual sum of squares $R'R$. By combining these two ideas, small values of $R_{ij\dots}^k(X^*)$ are associated with deviant and/or influential observations.

Cook statistic is defined as

$$D_k = \frac{T_k^2}{p} \cdot \frac{V(\hat{Y}_k)}{V(R_k)}$$

Note that D_k depends on three relevant quantities the number of parameters, p , the studentized residual, T_k which is a measure of the degree to consider k th observation as an outlier from the assumed model, and the ratio $V(\hat{Y}_k)/V(R_k)$ which measure the relative sensitivity of the estimate, \underline{b} , to potential outlying values at each data point.

3. Transformations

Suppose that the predicted values are given as follows

$$\hat{y}_k = \underline{x}'_k \underline{b}$$

where $\underline{x}'_k = (x_{1k}, x_{2k}, \dots, x_{pk})$, $k=1, 2, \dots, n$ and $\underline{b}' = (b_1, \dots, b_p)$.

The corresponding residual is

$$R_k = y_k - \underline{x}'_k \underline{b}$$

Let suppose that $b_j > 0$ for all j and $x_{jk} \neq x_{j'k'}$ for some $k \neq k'$.

Transform the given y_k and x'_k by $p+1$ real valued functions f_j , $j=0, 1, \dots, p$ as follows:

$$w_{0k} = f_0(y_k) = \frac{y_k - y_l}{y_u - y_l} \pi \quad k=1, 2, \dots, n \quad \dots (1)$$

$$w_{jk} = f_j(x_{jk}) = \frac{x_{jk} - x_{jl}}{x_{ju} - x_{jl}} \pi \quad \dots (2)$$

where $y_u = \max_{1 \leq k \leq n} y_k$, $y_l = \min_{1 \leq k \leq n} y_k$, $x_{ju} = \max_{1 \leq k \leq n} x_{jk}$ and $x_{jl} = \min_{1 \leq k \leq n} x_{jk}$

Consequently, the f_j ($f_j=0, 1, \dots, p$) satisfy the following conditions

$$(A) \quad 0 \leq f_j(x_{jk}) \leq \pi \quad \dots (3)$$

(B) f_j is a strictly monotone function in x_{jk}

Transform the given $\underline{b}' = (b_1, \dots, b_p)$ by d_j , $j=1, 2, \dots, p$, as follows:

$$d_j = b_j / \sum_{j=1}^p b_j \quad \dots (4)$$

$$\text{such that } \sum_{j=1}^p d_j = 1 \quad \dots (5)$$

From the result of the above transformations (1), (2), and (4), complex numbers C'_k and C_k are corresponded to y_k and $\hat{y}_k = \underline{x}'_k \underline{b}$, as follows

$$C'_k = \exp(iw_{0k}) \quad \dots (6)$$

$$C_k = \sum_{j=1}^p d_j \exp(iw_{jk}) = \sum_{j=1}^p d_j (\cos w_{jk} + i \sin w_{jk}) \quad k=1, 2, \dots, n \quad \dots (7)$$

where $i = \sqrt{-1}$ and d_j is the weight to be assigned to j th variable and satisfies (4). From (3) and (5), we can easily see arguments and absolutes of C'_k and C_k as follows:

$$(C) \quad 0 \leq \arg(C'_k) \leq \pi, \quad 0 \leq \arg(C_k) \leq \pi \quad \dots\dots(8)$$

$$(D) \quad |C'_k| = 1, \quad |C_k| \leq 1 \quad \dots\dots(9)$$

where $|C_k| = 1$ if and only if $w_{jk} = w_{j'k}$ for all j and j' .

4. The method for examining observations

From (6), (7), (8), and (9), we now establish the methods for displaying the transformed regression.

[Method 1] : Display C'_k on the circumference of the upper half of the unit circle in the complex plane according to (6) and C_k in and on it according to (7) where C'_k is displayed on it if and only if w_{jk} 's have same angles for all j . In this displayed figure, we can find clusters of the observations and examine the distance between the two points C'_k and C_k instead of R_k , N_k or T_k .

[Method 2] : Display C'_k on the circumference of the upper half of the unit circle, and then from this point display C_k to the opposite direction of C'_k .

In this displayed figure, we can examine the distance of the point C_k from the origin of the complex plane. [Method 2] is equivalent to display $r_k = C'_k - C_k$ which move C_k of $\overrightarrow{C'_k C_k}$ to the origin.

In case of $b_j < 0$ for some j , since $b_j x_{jk} = (-b_j) (-x_{jk})$, $k=1, \dots, n$, we may change the sign of all values of j th variable and then use the procedure of section 3 and this section.

In case of $x_{jk} = x_{j'k}$ for all k and k' , this cases occur in intercept models, without loss of generality we may change such intercept models to nonintercept models by taking the difference $y_k - b_j x_{jk}$ instead of response data y_k , $k=1, \dots, n$. Therefore, we can pick out the intercept term because of transformation (1).

5. Example

This example uses data given by Hald[3]. A fitted first order model of this data is $Y=52.577+1.468 X_1+0.662 X_2$. Table 1 lists x_{1k} , x_{2k} , $y_k-52.577$, R_k , R_k/s , T_k and D_k .

(Table 1) The 13 observations, residuals and Cook statistics

observation (k)	x_{1k}	x_{2k}	$y_k-52.577$	R_k	R_k/s	T_k	D_k
1	7	26	26.068	-1.6	-0.66	-0.77	0.014
2	1	29	21.666	1.0	0.42	0.48	0.008
3	11	56	51.72	-1.5	-0.62	-0.66	0.026
4	11	31	34.97	-1.7	-0.71	-0.81	0.05
5	7	52	43.3	-1.4	-0.58	-0.61	0.018
6	11	55	56.558	4.0	1.66	1.77	0.186
7	8	71	57.446	-1.3	-0.54	-0.68	0.053
8	1	31	19.89	-2.1	-0.87	-1.00	0.04
9	2	54	40.484	1.8	0.75	0.83	0.094
10	21	47	63.342	1.4	0.58	0.87	0.536
11	1	40	31.248	3.3	1.37	1.52	0.174
12	11	60	56.768	0.9	0.37	0.42	0.014
13	10	68	56.796	-2.9	-1.20	-1.36	0.154

Figure 1 is to display C'_k and C_k corresponding to these given observations and the representing way of observation # 5 by [Method1]. From figure 1, it is seen that the observation # 6 which has the greatest R_k (or $|T_k|$) has the largest distance of C_k from C'_k and since $m(\angle OQR)$ is greater than $m(\angle QC_sP)$, we can see that the variable X_1 provides more influence than the variable X_2 about the observation #5. Also, we can see that the observation #10 which takes the greatest value in the Cook's statistics D_k 's deviates from the other observations.

Figure 2 is to display r_k by [Method2] and here we may find the magnitude of residuals.

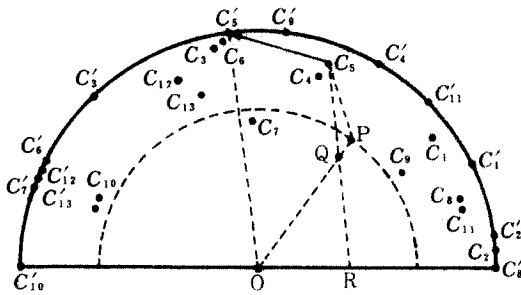


Fig.1 The Graph of C'_k, C_k by Method 1.

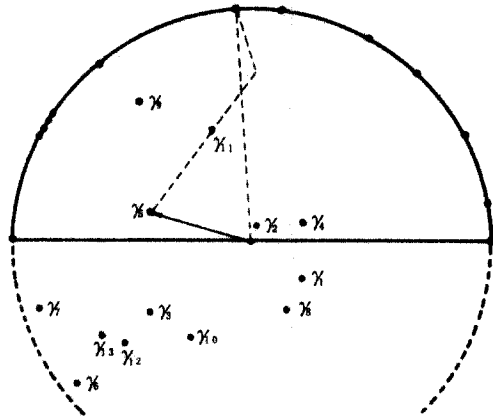


Fig.2 The Graph of γ_k by Method 2.

References

- [1] Andrews, D. F. and Pregibon, D (1978). *Finding the outliers that matter*. *J. Roy. Statist. Soc. Ser. B*40, 85 – 93
- [2] Atkinson, A. C. (1981). *Two graphical displays for outlying and influential observations in regression*. *Biometrika* 68, 1, 13-20
- [3] Behnken, D. W. and Draper, N. R. (1972). *Residuals and their variance patterns*. *Technometrics* 14, 102-111
- [4] Cook, R. D. (1977). *Detection of influential observations in linear regression*. *Technometrics* 19, 15-18
- [5] Doornbos, R. (1981). *Testing for a single outlier in a linear model*. *Biometrics* 37, 705-711
- [6] Draper, N. R. and John, J. A. (1981). *Influential observations and outliers in regression*. *Technometrics* 23, 21-28
- [7] Searle, S. R. (1971). *Linear models*. John Wiley & Sons
- [8] Wakimoto, K. and Taguri, M. (1978). *Constellation graphical method for representing multi-dimensional data*. *Ann. Inst. Statist. Math.* 30, A97-104