

## **Study of Outliers Under Variance-inflation Model in Experimental Designs**

Lalmohan Bhar and V.K. Gupta

*Indian Agricultural Statistics Research Institute, New Delhi-110012*

(Received : January, 2001)

### **SUMMARY**

For detecting outliers in linear regression model, Cook [2] developed a statistic under 'mean-shift' model. This statistic is appropriately developed under 'variance-inflation' model for detecting outliers in the general linear model ( $\mathbf{y}, \mathbf{X}\theta, \sigma^2\mathbf{I}$ ) for experimental designs. The simplified form of this statistic is obtained for both one-way and two-way elimination of heterogeneity settings. The relationship between Cook-statistic under 'mean-shift' model and Cook-statistic under 'variance-inflation' model for a single outlier is established and it has been shown that an observation which is an outlier under 'mean shift' model may not be so under 'variance-inflation' model. Through an example it has been shown that heteroscedasticity in the data may not influence the estimation of parameters.

*Key words* : Outlier, General linear model, Variance-inflation model, Cook-statistic.

### *1. Introduction*

A common approach to modelling outliers in the fixed effects linear model is to assume that outliers result from slippages in the expected values of contaminated observations. However, this assumption may not hold good in every occasion; instead the variances of all the observations may not remain constant which may severely affect the estimation of parameters. It may be a common occurrence in agricultural experiments where all the plots in a block may not be uniform, causing inflated variances for some of the observations. We may term these observations as outliers. In case of 'mean-shift' single outlier model, the test-statistic for the presence of an outlier is a monotonic function of the largest absolute Studentized residual (Srikantan [10]). This is not always true for 'variance-inflation' model. Cook *et al.* [4] showed that if the largest absolute residual corresponds to the largest absolute Studentized residual, then the estimated outlier is same for both the models.

For detecting outliers in linear regression model, Cook [2] developed a statistic under 'mean-shift' model. In designed experiments, the experimenter is

generally interested only in the estimation of a subset of parameters, rather than the whole set of parameters. One may, therefore, be interested to assess the influence of an outlying observation on the estimation of such function of subset of parameters. Since the design matrix in experimental designs is deficient in rank, the statistic developed by Cook [2] cannot be applied to it. Bhar and Gupta [1] modified this test-statistic for detecting outliers in experimental designs. However, no diagnostic measures seem to be developed so far for dealing with outliers under 'variance-inflation' model. Cook *et al.* [4] studied such outliers in linear regression model. Their study, however, was confined to the behaviour of the estimated outlier from 'mean-shift' model under 'variance-inflation' model. In the present paper we develop Cook-statistic under 'variance-inflation' model for detecting outliers in experimental designs. It is important to note that 'variance-inflation' model is actually a heteroscedastic model. However, our interest here is not in the estimation of parameters under a heteroscedastic model; instead we are interested to see whether inflated variances of some of the observations cause a severe effect on the estimation of a subset of parameters assuming that the model is homoscedastic.

In Section 2 this statistic has been developed under a general linear model set up of experimental designs. Section 3 considers the application of this statistic to some specific designs. The paper concludes with a Section on Discussion.

Throughout we use the following notations. All matrices and vectors are real, vectors being written as column vectors. We denote an  $n$ -component vector of all unities by  $\mathbf{I}_n$ , by  $\mathbf{J}_n$  an  $n \times n$  matrix of all unities and by  $\mathbf{I}_n$  an identity matrix of order  $n$ . Matrices  $\mathbf{A}'$ ,  $\xi(\mathbf{A})$ ,  $\mathbf{A}^-$  and  $\mathbf{A}^+$  will respectively denote the transpose, column space (range), a generalized inverse (g-inverse) and the Moore-Penrose inverse of  $\mathbf{A}$ .

## 2. Development of Cook – statistic

### 2.1 Cook-statistic in Mean-shift Model

Consider the general linear model

$$\mathbf{y} = \mathbf{X}\theta + \mathbf{e}; E(\mathbf{e}) = \mathbf{0}; D(\mathbf{e}) = \sigma^2\mathbf{I}_n; \sigma^2 > 0 \quad (2.1)$$

where  $\mathbf{y}$  is an  $n \times 1$  vector of observations,  $\mathbf{X}$  is an  $n \times p$  matrix of known constants with full column rank  $p$ ,  $\theta$  is a  $p \times 1$  vector of unknown parameters, and  $\mathbf{e}$  is an  $n \times 1$  vector of independent random variables each with zero mean and common variance  $\sigma^2 (>0)$ .

Cook [2] defined a statistic which measures the effect of an outlying observation on the estimation of parameters of the model as

$$D_i = \frac{\left( \hat{\theta} - \hat{\theta}_{(i)} \right)' \left[ D(\hat{\theta}) \right]^{-1} \left( \hat{\theta} - \hat{\theta}_{(i)} \right)}{\text{Rank} \left[ D(\hat{\theta}) \right]} \quad (2.2)$$

where  $\hat{\theta}$  and  $\hat{\theta}_{(i)}$  are the least-square estimates of  $\theta$  with and without the  $i$ -th data point respectively and  $D(x)$  denotes the dispersion matrix of  $x$ . The statistic provides a measure of distance between  $\hat{\theta}$  and  $\hat{\theta}_{(i)}$  in terms of descriptive levels of significance, because  $D_i$  is actually  $(1 - \alpha) \times 100\%$  confidence ellipsoid for the vector  $\theta$  under normal theory, which satisfies  $D_i \leq F(p, n - p, (1 - \alpha))$ . Extension of  $D_i$  to more than one outlier is straight forward. For usual interpretation of Cook-statistic, see Cook ([2], [3]).

## 2.2 Cook-statistic in Variance-inflation Model

Now consider the general linear model (2.1), for an experimental design  $d$  (say) with the design matrix  $X$  deficient in rank, i.e.,  $\text{Rank}(X) = m (< p)$ . Let  $\theta = (\theta_1' \ \theta_2')'$ , where  $\theta_1$  be a  $v$ -component vector containing all parameters of interest to the experimenter and  $\theta_2$  be a  $(p - v)$  component vector containing the set of nuisance parameters in the model that are not of much interest to the experimenter. Thus

$$y = (X_1 \ X_2) \begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix} + e \quad (2.3)$$

where  $X$  is partitioned in conformity with the partition of the parameter vector. We also assume that the ranks of  $X_1$  and  $X_2$  are  $m_1$  and  $m_2$  respectively. From the normal equations  $X'X\theta = X'y$  we obtain, on eliminating  $\theta_2$ , the equations involving only  $\theta_1$  as

$$C_{\theta_1} \theta_1 = Q_{\theta_1} \quad (2.4)$$

where  $C_{\theta_1} = X_1' B X_1$ ,  $Q_{\theta_1} = X_1' B y$  and  $B = I_n - X_2 (X_2' X_2)^{-1} X_2'$

The matrix  $B$  is symmetric and idempotent and the matrix  $C_{\theta_1}$  is symmetric. Since the linear model considered here is for experimental designs, it is not unrealistic to assume that the column space of  $X_1$  and  $X_2$  contain the vector  $\mathbf{1}$ , thus  $C_{\theta_1} \mathbf{1} = \mathbf{0}$ .

Suppose that  $t \{t \leq \min(m_1, m_2)\}$  of the  $n$  observations are suspected to be outliers in the sense that the variance of each of the  $t$  observations is shifted from the error variance of others. Then we have the variance-inflation model as

$$y = (X_1 \ X_2) \begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix} + z \tag{2.5}$$

where  $E(z) = 0$ ,  $D(z) = \sigma^2 \begin{bmatrix} W & 0' \\ 0 & I_{n-t} \end{bmatrix}$  and  $W = \text{diag}(w_1, \dots, w_t)$  with

$\sigma^2 > 0$ , and  $w_i > 1$ ;  $\forall i = 1, \dots, t$  are unknown. We also assume that the distribution of  $z$  is normal with mean vector  $0$  and dispersion matrix as given above. The model (2.5) can alternatively be written as

$$y = (X_1 \ X_2) \begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix} + U\delta + e \tag{2.6}$$

where  $U = (u_1, u_2, \dots, u_t)$ ;  $h_i^2 = (0, \dots, 1(i\text{-th}), 0, \dots, 0)$ ;  $i = 1, \dots, t$

and  $\delta = (\delta_1, \dots, \delta_t)'$

Now  $\delta_i$ 's are distributed normally with mean 0 and common variance  $\sigma^2 (w_i - 1)$ ;  $\forall i = 1, \dots, t$ ; and the components of  $e$  are as usual distributed normally with mean 0 and variance  $\sigma^2$ . Thus from (2.5) it follows that

$$D(y) = \sigma^2 \begin{bmatrix} W & 0' \\ 0 & I_{n-t} \end{bmatrix} = \sigma^2 [I_n - U(I_t - W)U'] = \Sigma \text{ (say)} \tag{2.7}$$

Clearly  $\Sigma$  is non-singular and  $\delta$  and  $e$  are independent by assumption.

$$\text{Now, } \Sigma^{-1} = \frac{1}{\sigma^2} I_n + \frac{1}{\sigma^2} U W_0 U' \tag{2.8}$$

$$\text{where } W_0 = \text{diag} \left( \frac{1-w_1}{w_1}, \dots, \frac{1-w_t}{w_t} \right) \tag{2.9}$$

Now from (2.5) the normal equations for estimating the linear functions of the parameters  $\theta_1$  are

$$C_{\theta_1}^H \theta_1 = Q_{\theta_1}^H \tag{2.10}$$

$$\text{where } C_{\theta_1}^H = [X_1' \Sigma^{-1} X_1] - [X_1' \Sigma^{-1} X_2] [X_2' \Sigma^{-1} X_2]^{-1} [X_2' \Sigma^{-1} X_1] \tag{2.11}$$

$$\text{and } Q_{\theta_1}^H = [X_1' \Sigma^{-1} y] - [X_1' \Sigma^{-1} X_2] [X_2' \Sigma^{-1} X_2]^{-1} [X_2' \Sigma^{-1} y] \tag{2.12}$$

Now  $C_{\theta_1}^H$  and  $Q_{\theta_1}^H$  can alternatively be written as

$$C_{\theta_1}^H = (C_{\theta_1} + LF^{-1}L')/\sigma^2 \quad (2.13)$$

$$Q_{\theta_1}^H = (Q_{\theta_1} + LF^{-1}U_*'By)/\sigma^2 \quad (2.14)$$

where  $L = X_1'BU_*$ ,  $F = I + U_*'X_2(X_2'X_2)^-X_2'U_*$ ,  $U_* = UW_*^{1/2}$  and  $C_{\theta_1}$  and  $Q_{\theta_1}$  are the corresponding matrices under the homoscedastic model. Clearly the matrix  $F$  is non-singular.

Now we have the following theorem

*Theorem 2.1.* (i)  $E(Q_{\theta_1}^H) = C_{\theta_1}^H \theta_1$

(ii)  $D(Q_{\theta_1}^H) = C_{\theta_1}^H$

We assume that the design  $d$  is connected and the rank of  $C_{\theta_1}^H$  is  $(v-1)$ . This can be checked by the connectedness property of the design under homoscedastic model, since a design that is connected under homoscedastic model is also connected under heteroscedastic model (see Gupta [7]). Let  $P\theta_1$  be the set of all  $(v-1)$  orthonormalized contrasts of the parameters  $\theta_1$ . The  $(v-1) \times v$  matrix  $P$  is such that  $PP' = I_{v-1}$ ,  $P'P = I_v - \frac{1}{v}J_v$  and the least square estimator of  $P\theta_1$  is given by  $P\hat{\theta}_1$ , where  $\hat{\theta}_1$  is any solution of the normal equation (2.10). Since the least square estimator is BLUE, the set is also BLUE. Then we have

$$P\hat{\theta}_1 = P(C_{\theta_1}^H)^- \Gamma_{P_3} = P(C_{\theta_1} + LF^{-1}L')^+ (Q_{\theta_1} + LF^{-1}U_*'By)$$

Now taking the Moore-Penrose inverse of  $C_{\theta_1} + LF^{-1}L'$ , we get

$$P\hat{\theta}_1 = P(C_{\theta_1}^+ - C_{\theta_1}^+LE^{-1}L'C_{\theta_1}^+) (Q_{\theta_1} + LF^{-1}U_*'By) \quad (2.15)$$

where  $E = F + L'C_{\theta_1}^+L$  and clearly it is non-singular.

Now we give the following lemma.

*Lemma 2.1.*  $D(P\hat{\theta}_1) = [P(C_{\theta_1} + LF^{-1}L')P']^{-1} \sigma^2$

Now we consider another model in which the outliers are actually omitted

$$y_{(t)} = \begin{pmatrix} X_{1(t)} & X_{2(t)} \end{pmatrix} \begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix} + z_{(t)} \tag{2.16}$$

where  $(t)$  denotes the removal of  $t$  rows from the corresponding vector and matrices. Note that the dispersion matrix  $\Sigma_{(t)}$  now becomes  $\sigma^2 I_{n-t}$ . The normal equations for estimating the contrasts of the parameters  $\theta_1$  under the model (2.16) are given by

$$C_{\theta_{1t}}^H \theta_{1t} = Q_{\theta_{1t}}^H \tag{2.17}$$

where  $C_{\theta_{1t}}^H = X'_{1(t)} X_{1(t)} - X'_{1(t)} X_{2(t)} (X'_{2(t)} X_{2(t)})^{-1} X'_{2(t)} X_{1(t)}$  (2.18)

and  $Q_{\theta_{1t}}^H = X'_{1(t)} y_{1(t)} - X'_{1(t)} X_{2(t)} (X'_{2(t)} X_{2(t)})^{-1} X'_{2(t)} y_{1(t)}$  (2.19)

The matrices  $C_{\theta_{1t}}^H$  and  $Q_{\theta_{1t}}^H$  can alternatively be written as

$$C_{\theta_{1t}}^H = C_{\theta_1} - L_* F_*^{-1} L_*' \tag{2.20}$$

and  $Q_{\theta_{1t}}^H = Q_{\theta_1} - L_* F_*^{-1} U' B y$  (2.21)

where  $L_* = X'_1 B U$  and  $F_* = U' B U$  (2.22)

We also assume that the residual design obtained after deleting  $t$  outlying observations remains connected. This fact can be checked by the connectedness property of the residual design under homoscedasticity, since it is known that if a design obtained after deleting  $t$  observations is connected under homoscedastic model it is also connected under heteroscedastic model (see Lal *et al.* [8]). The set of all orthonormalized contrasts for the parameters  $\theta_1$  is given by  $P\theta_{1t}$ . If

$\hat{\theta}_{1t}$  is any solution of the normal equations (2.17), then

$$P\hat{\theta}_{1t} = P(C_{\theta_{1t}}^H)^{-1} Q_{\theta_{1t}}^H = P(C_{\theta_1} - L_* F_*^{-1} L_*')^{-1} (Q_{\theta_1} - L_* F_*^{-1} U' B y) \tag{2.23}$$

Now taking the Moore-Penrose inverse of  $C_{\theta_1} - L_* F_*^{-1} L_*'$  we get

$$P\hat{\theta}_{1t} = P(C_{\theta_1}^+ + C_{\theta_1}^+ L_* E_*^{-1} L_*' C_{\theta_1}^+) (Q_{\theta_1} - L_* F_*^{-1} U' B y) \tag{2.24}$$

where  $E_* = F_* - L_*' C_{\theta_1}^+ L_*$  (2.25)

*Lemma 2.2.* The difference between the estimators of the contrasts under the model (2.6) and (2.16) can be expressed as

$$P\hat{\theta}_1 - P\hat{\theta}_{1t} = - P C_{\theta_1}^+ G B y$$

where  $G = (L E L' + L_* E_* L_*') C_{\theta_1}^+ X'_1 - L F^{-1} U' - L_* F_*^{-1} U' + L E L' C_{\theta_1}^+ L F^{-1} U' - L_* E_* L_*' C_{\theta_1}^+ L_* F_*^{-1} U'$  (2.26)

*Cook-statistic:* From (2.2) and using Lemma 2.2 we get Cook-statistic under variance-inflation model as

$$\begin{aligned} D_t^H &= \mathbf{y}'\mathbf{B}\mathbf{G}'\mathbf{C}_{\theta_1}^+\mathbf{P}\left[\mathbf{D}(\mathbf{P}\hat{\theta}_1)\right]^{-1}\mathbf{P}\mathbf{C}_{\theta_1}^+\mathbf{G}\mathbf{B}\mathbf{y}/(v-1) \\ &= \mathbf{y}'\mathbf{B}\mathbf{G}'\mathbf{C}_{\theta_1}^+[\mathbf{C}_{\theta_1} + \mathbf{L}\mathbf{F}^{-1}\mathbf{L}']\mathbf{C}_{\theta_1}^+\mathbf{G}\mathbf{B}\mathbf{y}/(v-1)\hat{\sigma}^2 \end{aligned} \quad (2.27)$$

where  $\hat{\sigma}^2$  is substituted for  $\sigma^2$ .

*Remark 2.1.* Suppose that only a single outlier is present. Without loss of generality we assume that the 1st observation is an outlier. Thus  $\mathbf{U} = \mathbf{u}_1$ . Then Cook-statistic takes the following form

$$D_1^H = d \left( \mathbf{L}'\mathbf{C}_{\theta_1}^+\mathbf{Q}_{\theta_1} - \mathbf{u}'_1\mathbf{B}\mathbf{y} \right)^2 \left( \mathbf{L}'\mathbf{C}_{\theta_1}^+\mathbf{L}_* \right) \quad (2.28)$$

where 
$$d = d_0 \left\{ 1 + \left( \frac{1}{1-w} - \mathbf{u}'_1\mathbf{B}\mathbf{u}_1 \right)^{-1} \mathbf{L}'\mathbf{C}_{\theta_1}^+\mathbf{L}_* \right\}$$

$$d_0 = \left\{ \left( \frac{1}{1-w} - \alpha \right)^{-1} + \alpha^{-1} \right\}^2$$

$$\alpha = (\mathbf{u}'_1\mathbf{B}\mathbf{u}_1 - \mathbf{L}'\mathbf{C}_{\theta_1}^+\mathbf{L}_*) \text{ and } \mathbf{L}_* = \mathbf{X}'_1\mathbf{B}\mathbf{u}_1 \quad (2.29)$$

*Remark 2.2.* If all the  $w_i$ 's are equal to 1, then  $w_* = 0$  and the Cook-statistic can be written as

$$D_1^H = \frac{\mathbf{y}'\mathbf{V}\mathbf{U}(\mathbf{U}'\mathbf{V}\mathbf{U})^{-1}\mathbf{U}'\mathbf{B}\mathbf{X}_1\mathbf{C}_{\theta_1}^+\mathbf{X}'_1\mathbf{B}\mathbf{U}(\mathbf{U}'\mathbf{V}\mathbf{U})^{-1}\mathbf{U}'\mathbf{V}\mathbf{y}}{(v-1)\hat{\sigma}^2} = D_1 \quad (2.30)$$

Thus variance-inflation model is transformed into mean-shift model (Bhar and Gupta [1]).

*Remark 2.3.* Cook-statistic under 'mean-shift' model, i.e.,  $D_t$  actually represents the volume of confidence ellipsoid (see Cook and Weisberg [5]) and thus  $D_t \leq F(1-\alpha, p, n_e)$ , where  $p$  is the rank of the corresponding dispersion matrix and  $n_e$  is the degrees of freedom (df) for estimating  $\sigma^2$ . Cook and Weisberg [5] suggest to compare  $D_t$  with critical value of  $F$ . Now  $D_t^H$  also represents confidence ellipsoid under heteroscedastic model. Thus the critical value of  $F$  can be taken as good approximation for  $D_t^H$ . In the literature it is found that  $\sigma^2$  under heteroscedastic model is estimated with the same df as it is estimated under homoscedastic model (see, e.g. Sen and Srivastava [9]). Moreover, here  $w_i$ 's are estimated either assuming that they are equal or

proportional to some character. Thus they may be treated as adhoc estimates and df may be taken as it is taken under homoscedastic model. We, therefore, suggest that  $D_i^H$  as given in (2.27) may be compared with the critical value of F with  $v-1$  and  $n_e$  df.

*Estimation of w*

In practice  $w_i$ 's are unknown. These are to be estimated from the data. Here we obtain the estimates of  $w_i$ 's by assuming that all  $w_i$ 's are equal. The extra sum of squares due to the parameters  $\delta$  under the model (2.6) is given by  $s_d^2$ , where  $s_d^2 = \mathbf{y}'\mathbf{V}\mathbf{U}(\mathbf{U}'\mathbf{V}\mathbf{U})^{-1}\mathbf{U}'\mathbf{V}\mathbf{y}$  and  $\mathbf{V} = \mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ . Under the variance-inflation model  $s_d^2$  has the expectation

$$E(s_d^2) = \sigma^2 + \text{tr}(\mathbf{U}'\mathbf{V}\mathbf{U})\sigma_d^2, \text{ where } \sigma_d^2 = \sigma^2(w - 1) \tag{2.31}$$

where  $\text{tr}(\cdot)$  denotes the trace of a square matrix. Since, the variance-covariance matrix of  $\mathbf{y}$  is now

$D(\mathbf{y}) = [\mathbf{I}_n + (w - 1)\mathbf{U}\mathbf{U}'] \sigma^2 = \sigma^2 \mathbf{I}_n + \mathbf{U}\mathbf{U}' \sigma_d^2$ , an estimate of  $\sigma_d^2$  is obtained as

$$\hat{\sigma}_d^2 = \frac{(s_d^2 - \hat{\sigma}^2)}{\text{tr}(\mathbf{U}'\mathbf{V}\mathbf{U})}, \text{ i.e., } \hat{\sigma}^2(w - 1) = \frac{(s_d^2 - \hat{\sigma}^2)}{\text{tr}(\mathbf{U}'\mathbf{V}\mathbf{U})} \tag{2.32}$$

$$\text{Finally, } \hat{w} = 1 + \frac{(s_d^2 - \hat{\sigma}^2)}{\hat{\sigma}^2 \text{tr}(\mathbf{U}'\mathbf{V}\mathbf{U})} \tag{2.33}$$

It is indeed possible that  $\hat{\sigma}_d^2$  is negative and  $\hat{w}$  is negative. However, we take the absolute value of  $\hat{w}$  as an estimate of  $w$ .

*3 Outliers in Some Specific Designs*

In the last section, we have derived Cook-statistic for outlier detection in general linear models. In the present section, we apply this test statistic to designs for both one-way and two-way elimination of heterogeneity settings.

**3.1 Outliers in Designs for One-Way Elimination of Heterogeneity**

**3.1.1 Some Preliminaries**

Consider the usual intra-block model of  $n$ -observations

$$\mathbf{y} = \mu\mathbf{1} + \Delta' \boldsymbol{\tau} + \mathbf{D}' \boldsymbol{\beta} + \mathbf{e} \tag{3.1}$$

with  $E(\mathbf{e}) = \mathbf{0}$  and  $D(\mathbf{e}) = \sigma^2 \mathbf{I}_n$

Here  $\Delta'$  is an  $n \times v$  ( $0 - 1$ ) design matrix for treatments,  $\mathbf{D}'$  is an  $n \times b$  ( $0 - 1$ ) design matrix for blocks,  $\mu$  is the general mean,  $\boldsymbol{\tau}$  is a



$v$ -component vector of treatment effects and  $\beta$  is a  $b$ -component vector of block effects. Also,  $\Delta'1 = 1$ ,  $D'1 = 1$ ,  $\Delta 1 = r$ ,  $D1 = k$ , where  $r = (r_1, \dots, r_v)'$  and  $k = (k_1, \dots, k_b)'$  are the vectors of replications and block sizes, respectively. We also define  $N = \Delta D' = ((n_{ij}))$ , where the non-negative integers  $n_{ij}$  denote the number of times  $i$ -th treatment appears in the  $j$ -th block;  $\sum_j n_{ij} = r_i$ ,

$\forall i = 1, \dots, v$ ;  $\sum_i n_{ij} = k_j, \forall j = 1, \dots, b$ . Here the parameter vector of interest is

$\tau$ , and the normal equations for  $\tau$  is obtained from (2.4) as

$$C_\tau \tau = Q_\tau \tag{3.2}$$

where  $C_\tau = \Delta B \Delta'$ ,  $Q_\tau = \Delta B y$  and  $B = I - D(D'D)^{-1}D'$  (3.3)

### 3.1.2 Single Outlier in Binary Block Designs

*Cook-statistic:* Without any loss of generality we assume that the observation pertaining to the first treatment in the first block is an outlier. Then the incidence matrix  $N$  can be written as

$$N = \begin{bmatrix} 1 & \epsilon' \\ f & N_0 \end{bmatrix} \tag{3.4}$$

where,  $f$  is a  $(v - 1)$  component  $(0 - 1)$  vector of incidence of remaining  $(v - 1)$  treatments in the first block,  $\epsilon$  is a  $(b - 1)$  component  $(0 - 1)$  vector of incidence of the first treatment in the remaining  $(b - 1)$  blocks and  $N_0$  is the incidence matrix of the remaining  $(v - 1)$  treatments in the remaining  $(b - 1)$  blocks.

The matrix  $L_*$  given in (2.22) becomes in the present case as

$$L_* = \Delta B u_1 = \frac{1}{k_1} \begin{pmatrix} k_1 - 1 \\ -f \end{pmatrix} = \begin{pmatrix} k_1 - 1 \\ k_1 \end{pmatrix}^{\frac{1}{2}} U_0 \tag{3.5}$$

where  $U_0 = \{k_1(k_1 - 1)\}^{-1/2} \begin{pmatrix} k_1 - 1 \\ -f \end{pmatrix}$  (3.6)

Thus  $L_*' C_{\theta_1}^+ L_* = \begin{pmatrix} k_1 - 1 \\ k_1 \end{pmatrix} U_0' C_\tau^+ U_0$

$$u_1' B u_1 = \frac{k_1 - 1}{k_1} \text{ and } \left( L_*' C_{\theta_1}^+ Q_{\theta_1} - u_1' B y \right) = u_1' v y = -r_1^*$$

where  $r_1^*$  is the first ordinary residual. Thus from (2.28) we get

$$D_1^H = d_0 \frac{U_0' C_\tau^+ U_0}{1 - U_0' C_\tau^+ U_0} \frac{t_1^2}{v-1} = d_0 D_1 \tag{3.7}$$

where  $t_1^2$  is the Studentized residual for the first observation

$$d_0 = \left( \frac{k_1}{wk_1 - w + 1} \right)^2 \left( \frac{1}{1 + \frac{1}{wk_1 - w + 1} U_0' C_\tau^+ U_0} \right) \tag{3.8}$$

and  $D_1$  is the corresponding Cook-statistic under the homoscedastic model.

*Lemma 3.1.*  $D_1^H \leq D_1$

*Proof:*  $d_0$  given in (3.8) can alternatively be written as

$$d_0 = \left( \frac{1}{1 + (w-1) \left( 1 - \frac{1}{k_1} \right)} \right) \left( \frac{1}{\left\{ 1 + (w-1) \left( 1 - \frac{1}{k_1} \right) \right\} (1 - U_0' C_\tau^+ U_0)} \right)$$

Since the residual design after deleting the outlying observation remains connected,  $(1 - U_0' C_\tau^+ U_0) > 0$  (cf. Dey [6]). Again  $w > 1$ , therefore  $d_0 \leq 1$ . Hence the proof follows from (3.7).

*Remark 3.1.* From Lemma 3.1 it is clear that an observation which is an outlier under mean-shift model may not be so under variance-inflation model. This was also observed by Cook *et al.* [4] in case of regression model.

### 3.1.3 Multiple Outliers in Variance-inflation Model

If the magnitude of the variance of each outlying observation, *i.e.*,  $w_i$  is known then the expression for  $D_i^H$  given in (2.28) can be directly used for testing the influence of these outliers. But if the values of  $w_i$ 's are unknown, then these are to be estimated from the data, which is extremely difficult. However, a convenient formula for estimating the magnitude of inflated variance can be obtained once we assume the equality of these magnitudes, and for that the estimating procedure is given in Section 2. Thus replacing the matrices  $C_{\theta_i}$  and  $B$  in the expression (2.28) by appropriate matrices for designs for one-way elimination of heterogeneity setting we can easily test for outliers.

*Remark 3.2.* Proceeding on the same lines as for designs for one-way elimination of heterogeneity, we can easily show that the Cook-statistic for designs for two-way elimination of heterogeneity can be written as

$$D_i^H = d_0^* \frac{U_i^* C_{\theta_i}^+ U_i^*}{\left( \frac{g_2}{g_1} - U_i^* C_{\theta_i}^+ U_i^* \right)^{v-1}} \frac{t_i^2}{v-1} = d_0^* D_i \quad (3.9)$$

where  $U_i^* = \frac{1}{\sqrt{g_1}} X_1' B u_i$

$$d_0^* = \left( \frac{1}{\{1 - (1-w)g_2\}^2} \right) \left( \frac{1}{1 + \frac{(1-w)g_1}{1 - (1-w)g_2} U_i^* C_{\theta_i}^+ U_i^*} \right)$$

$g_1$  is the norm of the vector  $X_1' B u_i$ ,  $g_2$  is the  $i$ -th diagonal element of  $B$ ,  $D_i$  is the corresponding Cook-statistic under 'mean-shift' model. Here the matrices  $C_{\theta_i}^+$ ,  $X_1$  and  $B$  are to be changed appropriately with the corresponding matrices of designs for two-way elimination of heterogeneity.

#### 4. Discussion

In this section, we discuss various aspects of heteroscedasticity in data generated from designed experiments. We begin with the following data set. Table 1 contains the grain yield (in Kg.) of rice variety IR8 with six different rates of seeding, from a Randomized Complete Block Design experiment with four replications. Figure within the parentheses indicates the observation number.

Table 1. Grain yield of rice variety IR8 (in Kg.)

Treatment No.	Replication No.			
	1	2	3	4
1	5113(1)	5398(7)	5307(13)	4678(19)
2	5346(2)	5952(8)	4719(14)	4264(20)
3	5272(3)	5713(9)	5483(15)	4749(21)
4	5164(4)	4831(10)	4986(16)	4410(22)
5	4804(5)	4848(11)	4432(17)	4748(23)
6	5254(6)	4542(12)	4919(18)	4098(24)

We compute the Cook-statistic given in (3.7) for all the 24 observations. The observations are numbered block wise in the ascending order of treatments. The largest value of Cook-statistic among the 24 observations is 0.018 and it pertains to observation number 8 corresponding to the treatment number 2 in the second block. The table value of  $F$  for 5 and 16 degrees of freedom and at 5% level of significance is 0.21. It, therefore, follows that none of the 24 observations is influential due to inflated variance.

We now consider the possibility of the presence of two outliers. We calculate Cook-statistic for two observations obtained from (2.27) for all

possible pair of observations. It is found that the maximum value of Cook-statistic is 0.201 corresponding to the pair of observations (2, 8). Again we find that two observations at a time are not influential. Naturally, a question may arise as to whether the data really exhibits heteroscedasticity? To ascertain this we construct a simple F-test taking mean square error (mse) for 2 outlying observations and mse for remaining 22 observations and perform all possible  ${}^n c_2$  F-tests. Following pairs of observations are found to be significant: (2,13), (3, 6), (3, 13) (5, 10), (7, 24), (8, 12), (8, 17), (8, 19), (8, 20), (8,22), (8, 24), (9, 17), (9, 20), (9, 22), (9, 24), (14, 19), (14, 21), (14, 23), (17, 22) and (21, 23). This result reveals one point clearly that the heteroscedasticity may not always affect the estimation of parameters. Even if heteroscedasticity is confirmed, we should not analyze the data assuming the model as heteroscedastic. Instead, we should check whether this heteroscedasticity really affects the parameter estimation or not. The present study provides a statistical tool to assess the influence by the presence of inflated variance of a few observations. If the heteroscedasticity of some observations does not affect the parameter estimation, we do not bother about their presence. We can as usual carry out our analysis assuming the model as homoscedastic.

Generally, in a particular situation we expect the presence of only a few outlying observations. For detecting  $t$  outliers through this statistic, one has to consider Cook-statistic for all  ${}^n c_t$  sets. But before analysing the data one can test for heterogeneity of error variance. Some simple tests can be constructed for this purpose. For example, if  $w_i$ 's are all equal, then one can construct an F-test as described earlier. This can also be generalized to the case when all  $w_i$ 's are not equal. Thus one can identify the subsets of data, in which heteroscedasticity is present, and Cook-statistic can be applied to that particular subset and not to all the  ${}^n c_t$  subsets.

In general if heteroscedasticity is confirmed, we take some appropriate measures such as variance stabilizing transformations and so on. But before applying transformation, one should check whether this heteroscedasticity is really influential in terms of parameter estimation. If this heteroscedasticity is not influential, we need not transform the data.

Another point worth noting is that sometimes heteroscedasticity arises due to the fact that some important covariates are not included in the model. For example inflated variance may occur due to heterogeneous plots within a block as explained in the introduction. But variance may not be inflated for all the plots. Moreover, it may not happen in all the blocks. Thus, one should identify the outlying observations in a particular block and then test whether they are influential or not. If they are influential, then one can include size or some other character of plot to ensure homoscedasticity.

## ACKNOWLEDGEMENT

The authors are grateful to the referee for making valuable suggestions that helped to improve the results considerably.

## REFERENCES

- [1] Bhar, L. and Gupta, V. K. (2001). A useful statistic for studying outliers in experimental designs. *Sankhya*, **B63**, 338 – 350.
- [2] Cook, R. D. (1977). Detection of influential observations in linear regression. *Technometrics*, **21**, 15–18.
- [3] Cook, R. D. (1979). Influential observations in linear regression. *J. Amer. Statist. Assoc.*, **74**, 169–174.
- [4] Cook, R. D., Holschuh, N. and Weisberg, S. (1982). A note on an alternative outlier model. *J. Roy. Statist. Soc.*, **B44**, 370–376.
- [5] Cook, R. D. and Weisberg, S. (1982). *Residual and Influence in Regression*. Chapman and Hall, New York.
- [6] Dey, A. (1993). Robustness of block designs against missing data. *Statistica Sinica*, **3**, 219–231.
- [7] Gupta, V. K. (1995). Universally optimal designs under heteroscedastic model. *Sankhya*, **B57**, 420–427.
- [8] Lal, K., Gupta, V. K. and Bhar, L. (2001). Robustness of designed experiments against missing data. *J. Appl. Statist.*, **28**, 63–79.
- [9] Sen, A. and Srivastava, M. (1990). *Regression Analysis: Theory, Methods, and Applications*. Springer-Verlag, New York, 120.
- [10] Srikantan, K. S. (1961). Testing for a single outlier in a regression model. *Sankhya*, **A23**, 251–260.