

A NOTE ON THE COOK'S DISTANCE

José A. Díaz-García and Graciela González-Farías

Comunicación Técnica No I-01-14/25-07-2001
(PE/CIMAT)



A Note on the Cook's Distance

José A. Díaz-García*

Departamento de Estadística y Cálculo
Universidad Autónoma Agraria Antonio Narro
25350 Buenavista, Saltillo, Coahuila, MÉXICO.
jadiaz@narro.uaaan.mx

and

Graciela González-Farías

Centro de Investigación en Matemáticas A. C.
Callejón de Jalisco s/n
36240 Guanajuato, MÉXICO
farias@cimat.mx

KEY WORD AND PHRASES: Diagnostic tools, Generalized Mahalanobis Distance, Influential data.

ABSTRACT

A modification of the classical Cook's distance is proposed, providing us with a generalized Mahalanobis distance in the context of multivariate elliptical linear regression models. We establish the exact distribution of a pivotal type statistics based on this generalized Mahalanobis distance, which provides critical points for the identification of outlier data points. We illustrate the procedure with an example, in the context of multiple and multivariate linear regression.

1. INTRODUCTION

The identification problem of outliers or influential data, in the univariate or multivariate linear regression setting and under the assumption of Gaussian errors has been studied by several authors like Cook (1977), Besley et al. (1980), Cook and Weisberg (1982) and Chatterjee and Hadi (1988), just to mention a few. Most of these results has been extended to the case of elliptical contour distributions, see for example Galea et al. (1997), Liu (2000) and Díaz-García et al. (2001), among others. One way or another, in all those works, the original idea of the so called Cook's distance is mentioned as a tool for identifying one influential point or sets of influential observations. However when we use this criteria we have only critical points

*This article was written while the first author was a Visiting Professor at the Department of Mathematics of ITESM, Monterrey, México and during a research visit at the Department of Probability and Statistics of CIMAT, Guanajuato, México

provided by an approximate central \mathcal{F} distribution used as it was proposed by Cook (1977). Our purpose here is to modify this distance and derive its exact distribution. Suppose $Y \in \mathbb{R}^{n \times p}$ has an elliptical distribution with location parameter $\mu \in \mathbb{R}^{n \times p}$ and scale matrix $\Sigma \otimes \Theta \in \mathbb{R}^{np \times np}$ with $\Sigma \in \mathbb{R}^{p \times p}$, $\Sigma > 0$ and $\Theta \in \mathbb{R}^{n \times n}$, $\Theta > 0$ then the density function is given by

$$f_Y(Y) = |\Sigma|^{-p/2} |\Theta|^{-n/2} g[\text{tr}(\Sigma^{-1}(Y - \mu)^T \Theta^{-1}(Y - \mu))]$$

where $g : \mathbb{R} \mapsto [0, \infty)$ is such that $\int_0^\infty u^{(np-2)/2} g(u) du < \infty$ being a density kernel. Let us denote this fact as $Y \sim \mathcal{E}l_{n \times p}(\mu, \Sigma \otimes \Theta, g)$.

This distribution family has been studied by different authors, see for example, Fang and Zhang (1990), Fang and Anderson (1990), Gupta and Varga (1993), among others. The elliptical distribution family includes subfamily distribution functions such as Gaussian, Pearson type VII and Logistic distributions, just to mention some.

Consider the multivariate linear regression model:

$$Y = X\beta + \epsilon, \quad (1)$$

where $Y \in \mathbb{R}^{n \times p}$ is the response matrix, $X \in \mathbb{R}^{n \times q}$, with $r(X) = q$, $\beta \in \mathbb{R}^{q \times p}$ the matrix of the unknown parameters and $\epsilon \in \mathbb{R}^{n \times p}$ is an error matrix, such that $\epsilon \sim \mathcal{E}l_{n \times p}(0, \Sigma \otimes I_n, g)$. This model is known as multivariate linear elliptical regression model. If g is a continuous and decreasing function, the maximum likelihood estimators for β and Σ are given by, see Fang and Zhang [pp. 129, 1990],

$$\hat{\beta} = (X^T X)^{-1} X^T Y = X^- Y \quad \text{and} \quad \hat{\Sigma} = u_0 (Y - X\hat{\beta})^T (Y - X\hat{\beta}),$$

where X^- is the Moore-Penrose inverse of X , and u_0 maximize the function

$$h(u) = u^{-np} g(p/u), u \geq 0$$

We consider a multivariate linear elliptical regression and propose an extension and modification of the Cook's distance. It will allow us to derive the exact distribution for the new distance providing a critical point to decide if a particular observation (or set of observations) behaves as an outlier.

2. MODIFIED DISTANCE : ONE OBSERVATION

Consider the modified multivariate linear elliptical regression model,

$$Y_{(i)} = X_{(i)} \beta^* + \epsilon_{(i)}, \quad \epsilon_{(i)} \sim \mathcal{E}l_{(n-1) \times p}(0, \Sigma^* \otimes I_{n-1}, g_{(i)}), \quad (2)$$

we get this model from (1) deleting the i th row from Y , X and ϵ , that is, deleting the i th observation.

For the modified model we have:

$$\hat{\beta}_{(i)} = (X_{(i)}^T X_{(i)})^{-1} X_{(i)}^T Y_{(i)} = X_{(i)}^- Y_{(i)} \quad \text{and} \quad \hat{\Sigma}_{(i)} = u_{0(i)} (Y_{(i)} - X_{(i)} \hat{\beta}_{(i)})^T (Y_{(i)} - X_{(i)} \hat{\beta}_{(i)}),$$

First of all we need to work out a simple representation for $\hat{\beta} - \hat{\beta}_{(i)}$. For that, consider the following partition matrices:

$$Y = \begin{pmatrix} Y_1^T \\ Y_2^T \\ \vdots \\ Y_n^T \end{pmatrix}, \quad Y_i \in \mathbb{R}^p \quad \epsilon = \begin{pmatrix} \epsilon_1^T \\ \epsilon_2^T \\ \vdots \\ \epsilon_n^T \end{pmatrix}, \quad \epsilon_i \in \mathbb{R}^p \quad X = \begin{pmatrix} X_1^T \\ X_2^T \\ \vdots \\ X_n^T \end{pmatrix}, \quad X_i \in \mathbb{R}^q.$$

therefore

$$X^T X = (X_1 X_2 \cdots X_n) \begin{pmatrix} X_1^T \\ X_2^T \\ \vdots \\ X_n^T \end{pmatrix} = \sum_{k=1}^n X_k X_k^T = \sum_{k \neq i} X_k X_k^T + X_i X_i^T = X_{(i)}^T X_{(i)} + X_i X_i^T,$$

and

$$X^T Y = (X_1 X_2 \cdots X_n) \begin{pmatrix} Y_1^T \\ Y_2^T \\ \vdots \\ Y_n^T \end{pmatrix} = \sum_{k=1}^n X_k Y_k^T = \sum_{k \neq i} X_k Y_k^T + X_i Y_i^T = X_{(i)}^T Y_{(i)} + X_i Y_i^T.$$

Note that if e_i^n is the i th vector of the canonical base in \mathbb{R}^n , that is, the unit vector, $e_i^n = (0 \cdots 0 \ 1 \ 0 \cdots 0)^T$, then: $e_i^{n^T} Y = Y_i^T$, $e_i^{n^T} X = X_i^T$ and $e_i^{n^T} \epsilon = \epsilon_i^T$.

By Rao [pp. 33, 1973], if A is nonsingular, v and u are two arbitrary vector, then

$$(A - uv^T)^{-1} = A^{-1} + \frac{A^{-1}uv^T A^{-1}}{1 - v^T A^{-1}u}$$

therefore, if we defined $A = X^T X$ and $u = v = X_i$, we get

$$(X^T X - X_i X_i^T)^{-1} = (X_{(i)}^T X_{(i)})^{-1} = (X^T X)^{-1} + \frac{(X^T X)^{-1} X_i X_i^T (X^T X)^{-1}}{(1 - h_i)}, \quad (3)$$

with $h_i = X_i^T (X^T X)^{-1} X_i$.

By (3), we get,

$$\begin{aligned} \hat{\beta} - \hat{\beta}_{(i)} &= (X^T X)^{-1} X^T Y - (X_{(i)}^T X_{(i)})^{-1} X_{(i)}^T Y_{(i)} \\ &= \left((X_{(i)}^T X_{(i)})^{-1} - \frac{(X^T X)^{-1} X_i X_i^T (X^T X)^{-1}}{(1 - h_i)} \right) X^T Y - (X_{(i)}^T X_{(i)})^{-1} X_{(i)}^T Y_{(i)} \\ &= (X_{(i)}^T X_{(i)})^{-1} (X^T Y - X_{(i)}^T Y_{(i)}) - \frac{(X^T X)^{-1} X_i X_i^T (X^T X)^{-1} X^T Y}{(1 - h_i)} \\ &= (X_{(i)}^T X_{(i)})^{-1} X_i Y_i^T - \frac{(X^T X)^{-1} X_i X_i^T (X^T X)^{-1} X^T Y}{(1 - h_i)}. \end{aligned} \quad (4)$$

Using (3) on, the first part of (4) we have,

$$\begin{aligned}
(X_{(i)}^T X_{(i)})^{-1} X_i Y_i^T &= \left((X^T X)^{-1} + \frac{(X^T X)^{-1} X_i X_i^T (X^T X)^{-1}}{(1 - h_i)} \right) X_i Y_i^T \\
&= (X^T X)^{-1} X_i Y_i^T + \frac{h_i (X^T X)^{-1} X_i Y_i^T}{(1 - h_i)} \\
&= \frac{(X^T X)^{-1} X_i Y_i^T}{(1 - h_i)}.
\end{aligned} \tag{5}$$

Substituting (5) in (4)

$$\begin{aligned}
\hat{\beta} - \hat{\beta}_{(i)} &= \frac{(X^T X)^{-1} X_i Y_i^T - (X^T X)^{-1} X_i X_i^T (X^T X)^{-1} X^T Y}{(1 - h_i)} \\
&= \frac{(X^T X)^{-1} X_i}{(1 - h_i)} (Y_i^T - X_i^T (X^T X)^{-1} X^T Y)
\end{aligned} \tag{6}$$

Now, since $\hat{\epsilon} = (Y - X\hat{\beta}) = (I - XX^{-})Y = (I - P)Y$, where P is the orthogonal projector over the image of X . Then $e_i^{n^T} \hat{\epsilon} = \hat{\epsilon}_i^T = e_i^{n^T} (Y - X\hat{\beta}) = Y_i^T - X_i^T (X^T X)^{-1} X^T Y$, so, we obtain:

$$\hat{\beta} - \hat{\beta}_{(i)} = \frac{(X^T X)^{-1} X_i \hat{\epsilon}_i^T}{(1 - h_i)} \tag{7}$$

Under the assumption of the elliptical distribution having moments, we propose the following modification to the Cook's distance, called \mathcal{D}_m :

$$\mathcal{D}_m = \text{vec}(\hat{\beta} - \hat{\beta}_{(i)})^T \widehat{\text{Cov}}(\text{vec}(\hat{\beta} - \hat{\beta}_{(i)}))^{-1} \text{vec}(\hat{\beta} - \hat{\beta}_{(i)}) \tag{8}$$

Note 1. Expression (8) is no more than an extension for the squared Mahalanobis generalized distance, see Rao and Mitra [203-206, 1971].

The second step is to find a simple expression for the variance covariance matrix $\text{Cov}(\text{vec}(\hat{\beta} - \hat{\beta}_{(i)}))$, under the assumption that the error distribution has moments.

Recall that $Y \sim \mathcal{E}l_{n \times p}(\mu, \Sigma \otimes \Theta, g)$, then its characteristic function is given by

$$\Psi_Y(T) = \text{etr}(i\mu T^T) \phi(\text{tr}(\Sigma T^T \Theta T))$$

and, $E(Y) = \mu$, $\text{Cov}(\text{vec}(Y)) = c_0(\Sigma \otimes \Theta)$, with $c_0 = -2\phi'(0)$, see Gupta and Varga [pp. 33, 1993].

Since $\hat{\epsilon}_i^T = e_i^{n^T} (Y - X\hat{\beta}) = e_i^{n^T} (I - P)Y$, it is clear that $\text{vec} \hat{\epsilon}_i^T = (I_p \otimes e_i^{n^T} (I - P)) \text{vec} Y$, therefore

$$\begin{aligned}
\text{vec}(\hat{\beta} - \hat{\beta}_{(i)}) &= \frac{(I_p \otimes (X^T X)^{-1} X_i)}{(1 - h_i)} (I_p \otimes e_i^{n^T} (I - P)) \text{vec} Y \\
&= \frac{(I_p \otimes (X^T X)^{-1} X_i e_i^{n^T} (I - P))}{(1 - h_i)} \text{vec} Y \\
&= \frac{(I_p \otimes (X^T X)^{-1} X_i P_i^T)}{(1 - h_i)} \text{vec} Y
\end{aligned} \tag{9}$$

where $P_i^T = e_i^{n^T} (I - P)$ is the i th row of the matrix $(I - P)$. Then

$$\begin{aligned}
\text{Cov}(\text{vec}(\hat{\beta} - \hat{\beta}_{(i)})) &= \frac{(I_p \otimes (X^T X)^{-1} X_i P_i^T)}{(1 - h_i)} \text{Cov}(\text{vec } Y) \frac{(I_p \otimes (X^T X)^{-1} X_i P_i^T)^T}{(1 - h_i)} \\
&= \frac{(I_p \otimes (X^T X)^{-1} X_i P_i^T)}{(1 - h_i)^2} (\Sigma^* \otimes I) (I_p \otimes P_i X_i^T (X^T X)^{-1}) \\
&= \frac{\|P_i\|^2 (\Sigma^* \otimes (X^T X)^{-1} X_i X_i^T (X^T X)^{-1})}{(1 - h_i)^2}
\end{aligned} \tag{10}$$

where $\Sigma^* = c_0 \Sigma$.

Note that,

$$\begin{aligned}
\|P_i\|^2 &= e_i^{n^T} (I - P)(I - P)e_i^n \\
&= e_i^{n^T} (I - P)e_i^n \\
&= e_i^{n^T} e_i^n - e_i^{n^T} X (X^T X)^{-1} X^T e_i^n \\
&= 1 - X_i (X^T X)^{-1} X_i^T \\
&= 1 - h_i
\end{aligned} \tag{11}$$

Substituting (11) in (10) we get,

$$\text{Cov}(\text{vec}(\hat{\beta} - \hat{\beta}_{(i)})) = \frac{(\Sigma^* \otimes (X^T X)^{-1} X_i X_i^T (X^T X)^{-1})}{(1 - h_i)} \tag{12}$$

Let $S_1 = \hat{\Sigma}/(u_0(n - q))$ and observe that $E(S_1) = \Sigma^*$, see Fang and Zhang [pp. 138, 1990], then

$$\widehat{\text{Cov}}(\text{vec}(\hat{\beta} - \hat{\beta}_{(i)})) = \frac{(S_1 \otimes (X^T X)^{-1} X_i X_i^T (X^T X)^{-1})}{(1 - h_i)} \tag{13}$$

Let $r_i = (X^T X)^{-1} X_i$. Based on the following standard results:

1. For $a \in \mathbb{R}^n$, $a^- = a^T / \|a\|^2$,
 2. Given $A \in \mathbb{R}^{p \times q}$, $(AA^T)^- = A^{T-} A^-$ with $A^{-1} = A^-$ if A is non singular,
 3. Given the matrices A and B , $(A \otimes B)^- = A^- \otimes B^-$,
- we get,

$$\begin{aligned}
(\widehat{\text{Cov}}(\text{vec}(\hat{\beta} - \hat{\beta}_{(i)})))^- &= \widehat{\text{Cov}}(\text{vec}(\hat{\beta} - \hat{\beta}_{(i)}))^- \\
&= \left(\frac{(S_1 \otimes r_i r_i^T)}{(1 - h_i)} \right)^- \\
&= \frac{(1 - h_i)}{\|r_i\|^4} (S_1^{-1} \otimes r_i r_i^T)
\end{aligned}$$

Therefore the modified Cook's distance can be rewritten as:

$$\begin{aligned}
\mathcal{D}_m &= \text{vec}(\hat{\beta} - \hat{\beta}_{(i)})^T \widehat{\text{Cov}}(\text{vec}(\hat{\beta} - \hat{\beta}_{(i)}))^{-1} \text{vec}(\hat{\beta} - \hat{\beta}_{(i)}) \\
&= \left(\frac{(I_p \otimes (X^T X)^{-1} X_i P_i^T) \text{vec} Y}{(1 - h_i)} \right)^T \frac{(1 - h_i)(S_1^{-1} \otimes r_i r_i^T)}{\|r_i\|^4} \left(\frac{(I_p \otimes (X^T X)^{-1} X_i P_i^T) \text{vec} Y}{(1 - h_i)} \right) \\
&= \frac{(1 - h_i)^{-1}}{\|r_i\|^4} \text{vec}^T Y (S_1^{-1} \otimes P_i r_i^T r_i r_i^T P_i^T) \text{vec} Y \\
&= (1 - h_i)^{-1} \text{vec}^T Y (S_1^{-1} \otimes P_i P_i^T) \text{vec} Y. \tag{14}
\end{aligned}$$

Alternatively, since $\text{tr} B X^T C X D = \text{vec}^T X (B^T D^T \otimes C) \text{vec} X = \text{vec}^T X (D B \otimes C^T) \text{vec} X$, for matrices of the correct sizes, we can write \mathcal{D}_m as

$$\mathcal{D}_m = (1 - h_i)^{-1} \text{tr} S_1^{-1} Y^T P_i P_i^T Y.$$

On the other hand, since $\hat{\epsilon}_i^T = e_i^{nT} (Y - X\hat{\beta}) = P_i Y$, then,

$$\begin{aligned}
\mathcal{D}_m &= (1 - h_i)^{-1} \text{tr} S_1^{-1} \hat{\epsilon}_i \hat{\epsilon}_i^T \\
&= (1 - h_i)^{-1} \text{tr} \hat{\epsilon}_i^T S_1^{-1} \hat{\epsilon}_i \\
&= (1 - h_i)^{-1} \hat{\epsilon}_i^T S_1^{-1} \hat{\epsilon}_i
\end{aligned}$$

In this way we have the following alternative expressions for the square of the modified Cook's distance:

$$\mathcal{D}_m = \begin{cases} \text{vec}(\hat{\beta} - \hat{\beta}_{(i)})^T \widehat{\text{Cov}}(\text{vec}(\hat{\beta} - \hat{\beta}_{(i)}))^{-1} \text{vec}(\hat{\beta} - \hat{\beta}_{(i)}) \\ (1 - h_i)^{-1} \text{vec}^T Y (S_1^{-1} \otimes P_i P_i^T) \text{vec} Y \\ (1 - h_i)^{-1} \text{tr} S_1^{-1} Y^T P_i P_i^T Y \\ (1 - h_i)^{-1} \hat{\epsilon}_i^T S_1^{-1} \hat{\epsilon}_i \end{cases} \tag{15}$$

Note 2. According with Chatterjee and Hadi [pp. 124, 1988], we could replace the matrix S_1 by one obtained using the reduced sample $(n - 1)$, denoted by S_{1_i} .

Note 3. Cook (1977), Chatterjee and Hadi [pp. 117, 1988], Díaz- García et. al (2001), and many others use the variance covariance matrix of $\text{vec}(\hat{\beta})$ to construct the distance measure. The reformulation we proposed is based on the replacement of that variance-covariance matrix by the variance covariance matrix of $\text{vec}(\hat{\beta} - \hat{\beta}_{(i)})$. We can also find this idea on Chatterjee and Hadi [pp. 150, 1988], for the univariate case, but for the evaluation of influential data on a particular regression coefficient, only the variance of one coefficient is used instead of the variance of the difference. The problem when this idea is extended to the multivariate case is that such matrix is singular, so we need to consider the Moore-Penrose inverse for the variance covariance matrix of $\text{vec}(\hat{\beta} - \hat{\beta}_{(i)})$.

Now we will consider the case when the errors from the multivariate regression model has an elliptical distribution without moments, (for example, when the errors have a matrix Cauchy distribution, see Gupta and Varga [pp. 76, 1993]), then the proposed distance maybe be defined by (15), without taking in account that,

$$\frac{(S_1 \otimes (X^T X)^{-1} X_i X_i^T (X^T X)^{-1})}{(1 - h_i)}$$

is the variance covariance matrix of $(\text{vec}(\hat{\beta} - \hat{\beta}_{(i)}))$. Let us called the distance matrix.

Theorem 1. *Consider the elliptical regression multivariate model given by (1), were the error distribution may or may not have moments. Then a modified squared Cook's distance to detect an outlier data, it can be written as:*

$$\mathcal{D}_m = \begin{cases} \text{vec}(\hat{\beta} - \hat{\beta}_{(i)})^T \left(\frac{(S_1 \otimes (X^T X)^{-1} X_i X_i^T (X^T X)^{-1})}{(1 - h_i)} \right)^{-} \text{vec}(\hat{\beta} - \hat{\beta}_{(i)}) \\ (1 - h_i)^{-1} \text{vec}^T Y (S_1^{-1} \otimes P_i P_i^T) \text{vec} Y \\ (1 - h_i)^{-1} \text{tr} S_1^{-1} Y^T P_i P_i^T Y \\ (1 - h_i)^{-1} \hat{\epsilon}_i^T S_1^{-1} \hat{\epsilon}_i \end{cases} \quad (16)$$

Note 4. From (16) is easy to see that if we want to implement this measure for all the data points, it is enough to fit the model once and from the usual output, we can construct the modified distance for each point. Note that the expression \mathcal{D}_m , on the univariate Normal case, coincides with the analysis of studentized residuals, see Besley et al. (p. 201, 1980) and Chatterjee (p. 78, 1988).

3. MODIFIED DISTANCE : MULTIPLE OBSERVATIONS

Let $I = \{i_1, i_2, \dots, i_k\}$ a subset of size k from $\{1, 2, \dots, n\}$, such that $(n - k) \geq q$. Now, under model (1), denote by $X_{(I)}$, $Y_{(I)}$ and $\hat{\epsilon}_{(I)}$, the regression, the data and the error matrices respectively, after deleting the corresponding observations according with the subindexes on I . Let $\hat{\beta}_{(I)}$, and $\hat{\Sigma}_{(I)}$ the corresponding maximum likelihood estimator in the model

$$Y_{(I)} = X_{(I)} \beta^* + \epsilon_{(I)}, \quad \epsilon_{(I)} \sim \mathcal{E}l_{(n-k) \times p}(0, \Sigma^* \otimes I_{n-k}, g_{(I)}).$$

Based on the equality

$$(A - BCD^T)^{-1} = A^{-1} + A^{-1}B(D^{-1} - C^T A^{-1}B)^{-1}C^T A^{-1}$$

where A and D are non singular matrices of order s and m respectively, B and C matrices of order $s \times m$ and using similar procedures as those from section 2, it is easy to verify that,

$$\hat{\beta} - \hat{\beta}_{(I)} = (X^T X)^{-1} X_I (I - H_I)^{-1} \hat{\epsilon}_I,$$

with $(I - H_I) = (I_k - X_I^T (X^T X)^{-1} X_I)$ and X_I the matrix with the corresponding rows of X according with I . Observe that, $\hat{\epsilon}_I = U_I^T \hat{\epsilon} = U_I^T (I - P)Y$, where

$$U_I^T = \begin{pmatrix} e_{i_1}^{nT} \\ e_{i_2}^{nT} \\ \vdots \\ e_{i_k}^{nT} \end{pmatrix}$$

We get,

$$\text{vec}(\hat{\beta} - \hat{\beta}_{(I)}) = (I_p \otimes (X^T X)^{-1} X_I (I - H_I)^{-1} P_I) \text{vec} Y,$$

with $P_I = U_I^T(I - P)$.

If we consider the case of existing moments for the distribution of errors,

$$\text{Cov}(\text{vec}(\hat{\beta} - \hat{\beta}_{(I)})) = (\Sigma^* \otimes (X^T X)^{-1} X_I (I - H_I)^{-1} X_I^T (X^T X)^{-1})$$

and

$$\widehat{\text{Cov}}(\text{vec}(\hat{\beta} - \hat{\beta}_{(I)})) = (S_1 \otimes (X^T X)^{-1} X_I (I - H_I)^{-1} X_I^T (X^T X)^{-1})$$

Under the same arguments for the non existing moments given in section 2, we have,

Theorem 2. *Consider the elliptical regression multivariate model given by (1), were the ϵ may or may not have moments. Then a modified squared Cook's distance to detect k influential data points , it can be written as*

$$\mathcal{D}_{m_I} = \begin{cases} \text{vec}(\hat{\beta} - \hat{\beta}_{(I)})^T (S_1 \otimes (X^T X)^{-1} X_I (I - H_I)^{-1} X_I^T (X^T X)^{-1})^{-1} \text{vec}(\hat{\beta} - \hat{\beta}_{(I)}) \\ \text{vec}^T Y (S_1^{-1} \otimes P_I^T (I - H_I)^{-1} P_I) \text{vec} Y \\ \text{tr} S_1^{-1} Y^T P_I^T (I - H_I)^{-1} P_I Y \\ \text{tr} S_1^{-1} \hat{\epsilon}_I^T (I - H_I)^{-1} \hat{\epsilon}_I \end{cases} \quad (17)$$

4. DISTRIBUTION FUNCTIONS ASSOCIATED WITH THE MODIFIED DISTANCES

The main reason to explore the modifications given in section 2 and 3 for the squared Cook's distance is that, instead of using an approximate \mathcal{F} distribution we will be able to derived an exact distribution for \mathcal{D}_m , which is invariant under the family of elliptical distributions. We propose a pivotal type estimator with an exact distribution for \mathcal{D}_{m_I} , the case of detection of several influential data simultaneously.

We will derive the distribution under the pivotal statistics for both cases: one influential observation at that time and multiple observations.

Theorem 3. *Under the assumptions of Theorem 1, we have,*

$$\frac{(n - q - p)}{p(n - q - 1)} \mathcal{D}_m^* \sim \mathcal{F}_{p, (n - q - p)} \quad (18)$$

where $\mathcal{F}_{p, (n - q - p)}$ denote a central \mathcal{F} distribution with p and $(n - q - p)$ degrees of freedom (df) and \mathcal{D}_m^* is given by (15) substituting S_1 by S_{1_i} , see Note 2.

Proof : It follows immediately from Theorem 5.2.2 in Anderson [p. 163, 1984] and Theorem 5.1.1 from Fang and Zhang [pp. 154, 158, 1990]. ■

From Theorem 3, given a significance level α , we may write the following decision rule: $Y_i, i = 1, 2, \dots, n$, is an outlier observation if

$$\frac{(n - q - p)}{p(n - q - 1)} \mathcal{D}_m^* \geq \mathcal{F}_{\alpha; p, (n - q - p)} \quad (19)$$

where $\mathcal{F}_{\alpha:p,(n-q-p)}$ the corresponding upper $\alpha - percentil$ from a \mathcal{F} distribution with p and $(n - q - p)$ df .

Note 5. For the univariate case, $p = 1$, the decision rule becomes: $Y_i, i = 1, 2, \dots, n$, is an outlier if

$$\mathcal{D}_m^* \geq \mathcal{F}_{\alpha:1,(n-q-1)} \quad (20)$$

where $\mathcal{F}_{\alpha:1,(n-q-1)}$ is the $\alpha - percentil$ from a \mathcal{F} distribution with 1 and $(n - q - 1)$ df , or equivalently:

$$\mathcal{D}_m^{*1/2} \geq \mathbf{t}_{\alpha/2:(n-q-1)} \quad (21)$$

where $\mathbf{t}_{\alpha/2:(n-q-1)}$ is the upper $\alpha/2 - percentil$ of a \mathbf{t} distribution with $(n - q - 1)$ df .

In a similar way, when we deal with multiple observations:

Theorem 4. *Under the assumptions of Theorem 2, we have,*

$$\frac{\mathcal{D}_{m_I}^*}{n - q - k} \sim \mathcal{LH}_{s,m,h} \quad (22)$$

where $\mathcal{LH}_{s,m,h}$ denote the central distribution for the Lawley-Hotelling statistics with parameters $s = \min(p, k), m = (|p - k| - 1)/2$ and $h = (n - q - p - 1)/2$ and $\mathcal{D}_{m_I}^*$ is given by (17) substituting S_1 by S_{1_I} .

Proof: It follows immediately from Theorem 5.3.1 from Gupta A.K and Varga T. [pp. 182, 301, 1993] and Theorem 10.6.2, Corolary 10.6.3, in Muirhead [pp. 468-471 and p. 471, 1982].

■

5. AN APPLICATION

We illustrate the use of the exact test given in Section 4 under two scenarios : simple regression and multivariate multiple regression.

The first data set was presented by Cook and Weisberg (pp. 204-207, 1994). This is a small data set with observations on 21 children, giving their *AGE* in months at first spoken word, and a *SCORE*, which is a measure of the development of the child. A plot (*AGE*, *SCORE*) is given in Figure 1 a). It is clear that are three observations that have a distinguishable behavior: 18, 1 and 17. If we think on an ordinary least squares (OLS) linear fit, *SCORE* seems to decrease with *AGE*. Case 18 appears to be poorly fitted by the linear trend, relative to the other data. Cases 1 and 17, have relatively large values of *AGE*. Figure 1 b), shows the Q-Q plot of the residuals from the OLS linear fit. It is clear that only observation 18 seems to

be a candidate for an *outlier* as it is defined in Chatterjee (pp.94-95, 1988).

Figure 1. Original data for the adaptive score measure and a Q-Q plot for the residual on the simple regression of *SCORE* on *AGE*.

Figure 2 shows the identification and detection of influential and outlier points, using different techniques. This analysis emphasize the fact mentioned before: studentized residuals coincide with the modified distance for the univariate case. Moreover studentized residuals are the base for the John and Draper distance as it is discussed in Draper and Smith (pp. 169-175, 1981). Figure 2 a) shows how the Cook's distance detects observation 17th, which is a leverage point as it is described in Cook and Weisberg (1994). Taking into account Figure 1 b) and according with Figures 2 b) and 2 c), observation 18 was a candidate for an outlier and the tests are in favor to declare observation 18th as an outlier observation. Note that the critical value in Figure 2 c) is the approximate value of an \mathcal{F} -distribution multiplied by $s^2(i)$ for $i = 18$ in order to plot the original distance of Draper and John.

For the multivariate multiple regression example, we generated 19 observations from a model $y = X\beta + \epsilon$, with normal errors, $p = 2$, and $q = 4$.

Figure 2. Identification of influence and outlier points based on a) Cook's Distance, b) The Modified Cook's Distance according with (20), c) Draper & John Distance. $n = 21$, $q = 2$, $p = 1$ and $s^2(i)$, the residual variance without the i -th observation was use in all the cases according with Note 2.

We fitted the model and constructed E , the residual sum of square matrix, then applied the Mahalanobis squared distance as it is shown in Seber (pp. 152-153, 1984). Figure 3 a) shows the Mahalanobis squared distance and suggests observations 10 and 11 as possible outliers. We applied (22), with $k = 2$, $Dm = 4.77$, compared with a critical value of 3.015; the percentile is approximated by an \mathcal{F} distribution as suggested by Seber (pp. 38-39, 563-564, 1984). The test is in favor of considering observations 10 and 11 as outliers.

Figure 3 b) shows the same analysis taken one observation at a time. We use an \mathcal{F} test based on $(1 - \alpha/n)$ instead of $(1 - \alpha)$ to get a simultaneous test with a nominal level at least α . In this case we get the same conclusion as with the test based on k observations. We recommend the use of the test based on k observations as given in (22), and selecting the k point using a graphical method as the one given in Figure 3 a).

It is important to recall this test is valid under elliptical distributions and not only for the normal error case.

Figure 3. Identification of outliers based on the Mahalanobis Distance on the residual matrix and detection of outliers based on the Modified Cook's Distance as given in Theorem 3.

6. ACKNOWLEDGEMENT

We are grateful to J. Ramón Domínguez for helping us with the figures. An S-plus program for all the calculations is available under request to jrdguez@cimat.mx

REFERENCES

- [1] Anderson, T. W. (1984). *An Introduction to Multivariate Statistical Analysis*, John Wiley & Sons, New York.
- [2] Besley, D., Kuh, E., and Welsch, R. (1980). *Regression Diagnostics: Identifying Influential Data and Sources of Collinearity*, John Wiley & Sons, New York.
- [3] Cook, R. D. (1977). "Detection of influential observations in linear regression", *Technometrics* 19, 15-18.
- [4] Cook, R. D., and Weisberg, S. (1982). *Residual and Influence in Regression*, Chapman and Hall, London.
- [5] Cook, R. D., and Weisberg, S. (1994). *An Introduction to Regression Graphics*, John Wiley & Sons. New York.
- [6] Chatterjee, S., and Hadi, A. S. (1988). *Sensitivity Analysis in Linear Regression*, John Wiley & Sons. New York.
- [7] Díaz-García, J.A., Galea, M., and Leiva- Sánchez, V. (2001). "Influence diagnostics for elliptical regression linear models", Submitted for publication.
- [8] Draper, N., and Smith, H. (1981). *Applied Regression Analysis*, (2nd ed.), John Wiley & Sons, New York.
- [9] Fang, K. T., and Anderson T. W. (1990). *Statistical Inference in Elliptically Contoured and Related Distributions*, Allerton Press Inc., New York.
- [10] Fang, K. T., and Zhang, Y. T. (1990). *Generalized Multivariate Analysis*, Science Press, Beijing, Springer-Verlang.
- [11] Galea, M., Paula, G., and Bolfarine, H. (1997). "Local influence in elliptical linear regression models", *The Statistician* 46, 71-79.
- [12] Gupta, A. K., and Varga, T. (1993). *Elliptically Contoured Models in Statistics*, Kluwer Academic Publishers, Dordrecht.
- [13] Liu, S. Z. (2000). "On local influence for elliptical linear models", *Statistical Papers* 41, 211-224.
- [14] Muirhead, R. J. (1982). *Aspects of Multivariate Statistical Theory*, John Wiley & Sons, New York.
- [15] C. R. Rao,(1973). *Linear Statistical Inference and its Applications* (2nd ed.), John Wiley & Sons, New York.
- [16] Rao, C. R. and Mitra, S. K.(1971). *Generalized Inverse of Matrices and its Applications* (2nd ed.), John Wiley & Sons, New York.
- [17] Seber, G.A.F.(1984). *Multivariate Observations* , John Wiley & Sons, New York.