# A Note on Cook's Distance in the Multivariate Linear Model

Whasoo Bae[a], Hyunmi Hwang[b], Choongrak Kim[1,b]

[a]Department of Data Science/Institute of Statistical Information, Inje University
[b]Department of Statistics, Pusan National University

## Abstract

We propose a version of Cook's distance (called local distance) in the multivariate linear model. The proposed version is a matrix, while the existing version of Cook's distance (called global distance) is a scalar. The existing Cook's distance is the trace of the proposed Cook's distance. In addition, we argue that the proposed Cook's distance has a more natural extension of the Cook's distance in the univariate linear model than the existing Cook's distance. An illustrative example based on a real data set is given.

Keywords: Global distance, influential observations, local distance.

## 1. Introduction

Most research on regression diagnostics are done for statistical models with one-dimensional response (univariate models); however, studies on multivariate regression diagnostics are relatively limited.

For the influence measures in the multivariate linear model, Caroni (1987) investigated a Studentized residual and suggested a version of Cook's distance (Cook, 1977) based on a confidence ellipsoid analogue. Altunkaynak and Ekni (2002) suggested a useful algorithm to compute the Cook's distance suggested by Caroni (1987). In addition, Diaz-Garcia *et al.* (2003) extended the concept of local influence (Cook, 1986) and likelihood displacement (Cook *et al.*, 1988) in the univariate linear model to the multivariate linear model. However, Tang and Fung (1997) considered case-deletion diagnostics for test statistics and Fung (1999) studied outlier diagnostics in several multivariate samples. For the diagnostics in the repeated measures or the longitudinal data, Preisser and Qaqish (1996) proposed deletion diagnostics for generalized estimating equations, Banerjee and Frees (1997) suggested influence diagnostics for linear longitudinal models, and Lindsey and Lindsey (2000) suggested some diagnostic tools for random effects in the repeated measures growth curve model.

In this paper, we suggest a new version of Cook's distance in the multivariate model. First, we mention the misleading aspect of the existing Cook's distance. The existing Cook's distance of the $i^{th}$ observation on the estimator of regression coefficient matrix is a scalar; however, the proposed Cook's distance in this thesis is a matrix that simultaneously reveals the influence of multiple outputs. We show that the existing version is sum of diagonal elements of the suggested version; therefore, the suggested version contains more diagnostic information than the existing version. In addition, this feature is demonstrated through numerical studies.

## 2. Multivariate Linear Model

Consider a multivariate linear model

$$Y = XB + U, \tag{2.1}$$

where $Y$ is $n \times q$ response matrix, $X$ is $n \times p$ design matrix, $B$ is unknown $p \times q$ regression coefficient matrix, and $U$ is $n \times q$ error matrix. Specifically let $y_i = (y_{i1}, y_{i2}, \ldots, y_{iq})'$, $x_i = (1, x_{i1}, x_{i2}, \ldots, x_{i,p-1})'$, $u_i = (u_{i1}, u_{i2}, \ldots, u_{iq})'$, $i = 1, 2, \ldots, n$ be the $i^{th}$ row of $Y, X, U$, respectively, and let $\beta_i = (\beta_{i1}, \beta_{i2}, \ldots, \beta_{iq})'$, $i = 1, 2, \ldots, p$ be the $i^{th}$ row of $B$. It is usually assumed that $y_1, y_2, \ldots, y_n$ are independent, and $\mathrm{E}[u_i] = 0$ and $\mathrm{Cov}(u_i) = \Sigma$ for all $i = 1, 2, \ldots, n$, where $\Sigma$ is a $q \times q$ variance-covariance matrix of the random vector $u_i$, and has $\sigma_{ij}$ as its $ij^{th}$ component. The goal is to suggest a new version of Cook's distance and comparing with the existing ones, so that we restrict our attention to the assumption of $\mathrm{Cov}(u_i) = \Sigma$ for all $i = 1, 2, \ldots, n$. For more general cases, see Preisser and Qaqish (1996), Banerjee and Frees (1997), and Lindsey and Lindsey (2000) among others.

If $X'X$ is non-singular, then the least squares estimator of the regression coefficient matrix $B$ is given by

$$\hat{B} = (X'X)^{-1}X'Y.$$

The fitted matrix can be expressed as $\hat{Y} = HY$, where $H = X(X'X)^{-1}X'$ is the hat matrix with $h_{ij} = x_i'(X'X)^{-1}x_j$ as the $ij^{th}$ component of $H$. Using this notation, the residual matrix is defined as $E = Y - \hat{Y}$, and let $E = (e_1, \ldots, e_n)'$, where $e_i'$, the $i^{th}$ row of $E$, is the $i^{th}$ residual vector. As an unbiased estimator of $\Sigma$, $\hat{\Sigma} = E'E/(n-p)$ is often used.

## 2.1. Existing version of Cook's distance $C_i^G$

In the multivariate linear model $E(Y) = XB$, Cook's distance of the $i^{th}$ observation on the estimator of the regression coefficient matrix $B$ is based on $\hat{B} - \hat{B}_{(i)}$, where $\hat{B}_{(i)}$ is the least squares estimator of $B$ based on $n - 1$ observations after deleting the $i^{th}$ observation $(x_i, y_i)$. Note that $\hat{B} - \hat{B}_{(i)}$ is a $p \times q$ matrix, it is not straightforward to normalize to a scalar. To overcome this situation, by using the vec operation the multivariate linear model in (2.1) can be reexpressed as

$$\mathrm{vec}(Y) = \left(I_q \otimes X\right)\mathrm{vec}(B) + \mathrm{vec}(U).$$

Using this notation, Caroni (1987) and Diaz-Garcia $et\ al.$ (2003) suggested a version of Cook's distance in the multivariate linear model as

$$C_i^G = \frac{1}{p}\left[\mathrm{vec}\left(\hat{B} - \hat{B}_{(i)}\right)'\left(\mathrm{Cov}\left(\mathrm{vec}\left(\hat{B}\right)\right)\right)^{-1}\mathrm{vec}\left(\hat{B} - \hat{B}_{(i)}\right)\right].$$

Here we call $C_i^G$ a global distance, because it does not distinguish $q$ multiple outputs $(y_{i1}, y_{i2}, \ldots, y_{iq})$, but it computes the effect of $q$ multiple outputs simultaneously. That is to say, $C_i^G$ renders one scalar value as influence of $q$ multiple outputs.

To express $C_i^G$ as a function of basic building blocks, we note that $\mathrm{Cov}(\mathrm{vec}(\hat{B})) = \Sigma \otimes (X'X)^{-1}$ and

$$\hat{B} - \hat{B}_{(i)} = \frac{(X'X)^{-1}x_i e_i'}{1 - h_{ii}}. \tag{2.2}$$

In addition, if we use

$$\left(\text{vec}(\boldsymbol{Z}')\right)' \left(\boldsymbol{A} \otimes \boldsymbol{B}'\right) \text{vec}(\boldsymbol{Z}) = \text{tr}(\boldsymbol{A}\boldsymbol{Z}'\boldsymbol{B}\boldsymbol{Z}) \tag{2.3}$$

then, we have

$$C_i^G = \frac{1}{p} \frac{h_{ii}}{1 - h_{ii}} \frac{\boldsymbol{e}_i' \hat{\boldsymbol{\Sigma}}^{-1} \boldsymbol{e}_i}{1 - h_{ii}}. \tag{2.4}$$

This equation also can be expressed as

$$C_i^G = \frac{1}{p} \frac{r_i^2 h_{ii}}{1 - h_{ii}},$$

where

$$r_i^2 = \frac{\boldsymbol{e}_i' \hat{\boldsymbol{\Sigma}}^{-1} \boldsymbol{e}_i}{1 - h_{ii}}$$

is a version of studentized residual in the multivariate model.

For the influence of a set of observations, let $K = \{i_1, i_2, \ldots, i_k\}$ be an index set containing $k$ sets. Then the influential set version of the Cook's distance considered above is

$$C_K^G = \frac{1}{p} \left\{ \text{vec}\left(\hat{\boldsymbol{B}} - \hat{\boldsymbol{B}}_{(K)}\right)' \left[\boldsymbol{\Sigma}^{-1} \otimes (\boldsymbol{X}'\boldsymbol{X})\right] \text{vec}\left(\hat{\boldsymbol{B}} - \hat{\boldsymbol{B}}_{(K)}\right) \right\}.$$

To express $C_K^G$ as a function of basic building blocks we first note that

$$\hat{\boldsymbol{B}} - \hat{\boldsymbol{B}}_{(K)} = (\boldsymbol{X}'\boldsymbol{X})^{-1} \boldsymbol{X}_K' (\boldsymbol{I} - \boldsymbol{H}_K)^{-1} \boldsymbol{e}_K,$$

where $\boldsymbol{X}_K = (\boldsymbol{x}_{i_1}, \boldsymbol{x}_{i_2}, \ldots, \boldsymbol{x}_{i_k})$ is $k \times p$ matrix, $\boldsymbol{e}_K = (\boldsymbol{e}_{i_1}, \boldsymbol{e}_{i_2}, \ldots, \boldsymbol{e}_{i_k})$ is $k \times q$ matrix, and $\boldsymbol{H}_K = \boldsymbol{X}_K(\boldsymbol{X}'\boldsymbol{X})^{-1}\boldsymbol{X}_K'$. Then, by (2.3),

$$\begin{aligned}
C_K^G &= \frac{1}{p} \left\{ \text{vec}\left(\hat{\boldsymbol{B}} - \hat{\boldsymbol{B}}_{(K)}\right)' \left[\boldsymbol{\Sigma}^{-1} \otimes (\boldsymbol{X}'\boldsymbol{X})\right] \text{vec}\left(\hat{\boldsymbol{B}} - \hat{\boldsymbol{B}}_{(K)}\right) \right\} \\
&= \frac{1}{p} \text{tr}\left(\left(\hat{\boldsymbol{B}} - \hat{\boldsymbol{B}}_{(K)}\right)' \boldsymbol{X}'\boldsymbol{X}\left(\hat{\boldsymbol{B}} - \hat{\boldsymbol{B}}_{(K)}\right)\boldsymbol{\Sigma}^{-1}\right) \\
&= \frac{1}{p} \text{tr}\left(\boldsymbol{e}_K' (\boldsymbol{I} - \boldsymbol{H}_K)^{-1} \boldsymbol{H}_K (\boldsymbol{I} - \boldsymbol{H}_K)^{-1} \boldsymbol{e}_K \hat{\boldsymbol{\Sigma}}^{-1}\right).
\end{aligned}$$

## 2.2. A proposed version of Cook's distance $C_i^L$

Recall that Cook's distance for the $i^{th}$ observation in the univariate linear model can be written as

$$\begin{aligned}
D_i &= \frac{1}{p} \left(\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}_{(i)}\right)' \text{Cov}\left(\hat{\boldsymbol{\beta}}\right)^{-1} \left(\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}_{(i)}\right) \\
&= \frac{1}{p} \left(\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}_{(i)}\right)' \sigma^{-2} (\boldsymbol{X}'\boldsymbol{X})\left(\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}_{(i)}\right) \\
&= \frac{1}{p} \frac{\left(\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}_{(i)}\right)'}{\sigma} (\boldsymbol{X}'\boldsymbol{X}) \frac{\left(\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}_{(i)}\right)}{\sigma},
\end{aligned}$$

which motivates Cook's distance in the multivariate linear model as

$$C_i^L = \frac{1}{p}\hat{\Sigma}^{-\frac{1}{2}}\left(\hat{B} - \hat{B}_{(i)}\right)'(X'X)\left(\hat{B} - \hat{B}_{(i)}\right)\hat{\Sigma}^{-\frac{1}{2}}.$$

If we use $\hat{B} - \hat{B}_{(i)} = (X'X)^{-1}x_ie_i'/(1 - h_{ii})$, then we can express $C_i^L$ as basic building blocks, *i.e.*,

$$C_i^L = \frac{1}{p}\hat{\Sigma}^{-\frac{1}{2}}\frac{e_ix_i'(X'X)^{-1}}{1 - h_{ii}}(X'X)\frac{(X'X)^{-1}x_ie_i'}{1 - h_{ii}}\hat{\Sigma}^{-\frac{1}{2}}$$

$$= \frac{1}{p}\frac{h_{ii}}{(1 - h_{ii})^2}\hat{\Sigma}^{-\frac{1}{2}}e_ie_i'\hat{\Sigma}^{-\frac{1}{2}},$$

where $e_i = y_i - x_i'\hat{B}$ and $h_{ii} = x_i'(X'X)^{-1}x_i$. Note that $C_i^L$ is not scalar but a $q \times q$ matrix. In addition, it is clear that $\text{tr}(C_i^L) = C_i^G$.

## 2.3. Remarks on local influence and global influence

There is a serious disadvantage in $C_i^G$. Note that some of the diagonal elements of the local distance $C_i^L$ could be small or large, and therefore, $C_i^G$ may not reveal the actual influence of the $i^{th}$ observation. For example, $C_i^G = 2.0$ is due to either $C_i^L(1,1) = 1.0$, $C_i^L(2,2) = 1.0$ or $C_i^L(1,1) = 1.9$, $C_i^L(2,2) = 0.1$. Hence, $C_i^L$ is more informative measure than $C_i^G$ to represent actual influence of the $i^{th}$ observation.

To see further relations between $C_i^L$ and $C_i^G$, note that $e_i'\hat{\Sigma}^{-1}e_i$ in $C_i^G$ is replaced by $\hat{\Sigma}^{-1/2}e_ie_i'\hat{\Sigma}^{-1/2}$ in $C_i^L$. To compare these two terms we assume, for simplicity, that $q = 2$ and correlation between $y_{i1}$ and $y_{i2}$ is close to zero. Then, $e_i'\hat{\Sigma}^{-1}e_i \simeq e_{i1}^2/\hat{\sigma}_1^2 + e_{i2}^2/\hat{\sigma}_2^2$, but

$$\hat{\Sigma}^{-\frac{1}{2}}e_ie_i'\hat{\Sigma}^{-\frac{1}{2}} \simeq \begin{bmatrix} \dfrac{e_{i1}^2}{\hat{\sigma}_1^2} & \dfrac{e_{i1}e_{i2}}{\hat{\sigma}_1\hat{\sigma}_2} \\[3mm] \dfrac{e_{i1}e_{i2}}{\hat{\sigma}_1\hat{\sigma}_2} & \dfrac{e_{i2}^2}{\hat{\sigma}_2^2} \end{bmatrix}.$$

Therefore, if we assume weak correlations between responses, we may say that the $j^{th}$ diagonal element $C_i^L(jj)$ represents the influence of $y_{ij}$ and $x_i$ on $\hat{B}$. Of course, if $q$ responses are highly correlated, then interpretations on each component of $C_i^L$ will be very complicated. Hence, $C_i^L$ can be called the local distance representing the influence of each response separately for the $i^{th}$ observation. On the other hand, $C_i^G$, sum of diagonal elements of $C_i^L$, can be called the global distance representing the influence of all the $q$ responses simultaneously for the $i^{th}$ observation.

We can easily extend $C_i^L$ to $C_K^L$, local Cook's distance for the set of observations in $K = (i_1, i_2, \ldots, i_k)$, *i.e.*,

$$C_K^L = \frac{1}{p}\hat{\Sigma}^{-\frac{1}{2}}\left(\hat{B} - \hat{B}_{(K)}\right)'(X'X)\left(\hat{B} - \hat{B}_{(K)}\right)\hat{\Sigma}^{-\frac{1}{2}}. \tag{2.5}$$

To express $C_K^L$ as a function of basic building blocks, we use

$$\hat{B} - \hat{B}_{(K)} = (X'X)^{-1}X_K'(I - H_K)^{-1}e_K.$$

Then the local Cook's distance for the set of observations in $K$ can be reexpressed as

$$C_K^L = \frac{1}{p}\hat{\Sigma}^{-\frac{1}{2}}e_K'(I - H_K)^{-1}H_K(I - H_K)^{-1}e_K\hat{\Sigma}^{-\frac{1}{2}}. \tag{2.6}$$

Table 1: Two versions of Cook's distance $C_K^G$ and $C_K^L$ when $K = \{i\}$ in the automobile tire data.

| $K$ | $C_K^L$ | | $C_K^G$ |
|---|---|---|---|
| 1 | 0.317 | 0.093 | 0.344 |
| | 0.093 | 0.027 | |
| 2 | 0.247 | −0.330 | 0.688 |
| | −0.330 | 0.441 | |
| 3 | 0.303 | −0.073 | 0.321 |
| | −0.073 | 0.018 | |
| 4 | 0.067 | 0.066 | 0.131 |
| | 0.066 | 0.064 | |
| 5 | 0.024 | −0.081 | 0.298 |
| | −0.081 | 0.274 | |
| 6 | 0.132 | −0.066 | 0.165 |
| | −0.066 | 0.033 | |
| 7 | 0.010 | 0.109 | 1.196 |
| | 0.109 | 1.186 | |
| 8 | 0.010 | −0.025 | 0.073 |
| | −0.025 | 0.063 | |
| 9 | 0.704 | −0.785 | 1.580 |
| | −0.785 | 0.876 | |
| 10 | 0.038 | 0.037 | 0.073 |
| | 0.037 | 0.035 | |

Table 2: Five largest Cook's distances for $C_K^G$ and $C_K^L$ when $K = \{i, j\}$ in the automobile tire data.

| $K$ | $C_K^L$ | | $C_K^G$ |
|---|---|---|---|
| 5, 9 | 13.1654 | 0.7721 | 13.3879 |
| | 0.7721 | 0.2225 | |
| 7, 9 | 4.9712 | 3.3713 | 8.0197 |
| | 3.3713 | 3.0485 | |
| 5, 7 | 1.4881 | 2.8318 | 7.1213 |
| | 2.8318 | 5.6332 | |
| 3, 7 | 1.9939 | 2.9972 | 6.9868 |
| | 2.9972 | 4.9929 | |
| 3, 9 | 0.2126 | −0.3611 | 6.9321 |
| | −0.3611 | 6.7195 | |

If we assume, for simplicity, that $q = 2$, $K = \{i, j\}$ and correlation between $y_{i1}$ and $y_{i2}$ is close to zero, then the first diagonal element of $C_K^L$ represents the influence of $y_{i1}$ and $y_{j1}$ on $\hat{B}$, and the second diagonal element of $C_K^L$ represents the influence of $y_{i2}$ and $y_{j2}$ on $\hat{B}$; however, $C_K^G$ shows only the sum of two influences of $\{y_{i1}, y_{j1}\}$ and $\{y_{i2}, y_{j2}\}$ on $\hat{B}$.

## 3. Example

As an illustrative example for two versions of Cook's distance $C_K^G$ and $C_K^L$, we take automobile tire data described in Green (1978). The data consist of 2 responses and 4 covariates, and the number of observations is 10, *i.e.*, $q = 2$, $p = 5$ and $n = 10$.

After fitting a multivariate linear model $Y = XB + U$, we evaluate two influence measures in Table 1. First, the global Cook's distance $C_i^G$ reveals that the $9^{th}$ observation is very influential, and the $7^{th}$ observation is also quite influential. If we see the local Cook's distance $C_i^L$, then the global influence of the $9^{th}$ observation 1.580 is sum 0.704 (local influence of $y_{9,1}$) and 0.876 (local influence of $y_{9,2}$); however, the global influence of the $7^{th}$ observation 1.196 is sum of 0.010 (local influence of $y_{7,1}$) and 1.186 (local influence of $y_{7,2}$). Therefore, the high influence of the $7^{th}$ observation is mainly due to

the $2^{nd}$ component of the response $y_{7,2}$. For the influence of two observations (Table 2), we see that $\{5, 9\}$ has the largest global Cook's distance mainly due to $y_{51}$ and $y_{91}$; in addition, the contribution of $y_{52}$ and $y_{92}$ is almost negligible.

## References

Altunkaynak, B. and Ekni, M. (2002). Detection of influential observation vectors for multivariate linear regression, *Journal of Mathematics and Statistics*, **31**, 139–151.

Banerjee, M. and Frees, E. W. (1997). Influence diagnostics for linear longitudinal models, *Journal of the American Statistical Association*, **92**, 999–1005.

Caroni, C. (1987). Residuals and influence in the multivariate linear model, *The Statistician*, **36**, 365–370.

Cook, R. D. (1977). Detection of influential observations in linear regression, *Technometrics*, **19**, 15–18.

Cook, R. D. (1986). Assessment of local influence, *Journal of the Royal Statistical Society, Series B*, **48**, 133–169.

Cook, R. D., Pena, D. and Weisberg, S. (1988). The likelihood displacement, *Communications in Statistics - Theory and Methods*, **17**, 623–640.

Diaz-Garcia, J. A., Rohas, M. G. and Leiva-Sanchez, V. (2003). Influence diagnostics for elliptical multivariate linear regression models, *Communications in Statistics - Theory and Methods*, **32**, 625–641.

Fung, W. K. (1999). Outlier diagnostics in several multivariate samples, *Journal of the Royal Statistical Society, Series D*, **48**, 73–84.

Green, P. E. (1978). *Analyzing Multivariate Data*, Hindsdale, Ill, The Dryden Press.

Lindsey, P. J. and Lindsey, J. K. (2000). Diagnostic tools for random effects in the repeated measures growth curve model, *Computational Statistics and Data Analysis*, **33**, 79–100.

Preisser, J. S. and Qaqish, B. F. (1996). Deletion diagnostics for generalized estimating equations, *Biometrika*, **83**, 551–562.

Tang, M. K. and Fung, W. K. (1997). Case-deletion diagnostics for test statistics in multivariate regression, *Australian and New Zealand Journal of Statistics*, **39**, 345–353.