

Influence in Fitting an Equicorrelation Model¹⁾

Myung Geun Kim²⁾ and Kang-Mo Jung³⁾

Abstract

The influence in fitting an equicorrelation model is investigated using the influence function. The influence functions for the model parameters are derived and its sample versions are used for investigating the influence of observations on the estimators of the parameters. Some relationships among the sample versions are found. We will derive a measure for identifying observations that have a large influence on the test of fitting the equicorrelation model using the influence function method. An example is given for illustration.

Keywords : Equicorrelation, influence functions, test.

1. Introduction

The equicorrelation model, also called the intraclass model has a pattern of equal variances and equal covariances in the covariance matrix. For this model the variables are correlated and every pair of variables has the same correlation coefficient. An acceptance of the equicorrelation model is sufficient for the validity of standard ANOVA approach to repeated measures design (Rencher, 1995). Also, the common correlation coefficient is used for measuring the agreement between quantitative measures in epidemiological studies and Giraudeau et al. (1996) studied the single case deletion effect on the common correlation coefficient estimate. There are few or no works on diagnostics in the equicorrelation model.

The influence function (Hampel, 1974) for a parameter at a point measures the effect of an infinitesimal contamination at that point on the estimator of the parameter. Hence the influence function can serve as a diagnostic method of detecting influential observations in estimating the model parameters and in performing a test (Jung and Kim, 1999).

1) This work was in partial supported by grants from Institute of Applied Science and Technology at Seowon University.

2) Professor, Dept of Applied Statistics, Seowon University, Chongju 361-742.
E-mail : mgkim@seowon.ac.kr

3) Assistant Professor, Dept. of Informatics & Statistics, Kunsan National University, Kunsan 573-701.
E-mail : kmjung@kunsan.ac.kr.

In this work we will study the influence in the equicorrelation model using the influence function. In Section 2 the likelihood equations for estimating the model parameters are reviewed. In Section 3 we will derive influence functions for σ^2 and ρ by defining appropriate functionals, and consider three sample versions of these influence functions that will be used for investigating the influence of observations on $\hat{\sigma}^2$ and $\hat{\rho}$. Some relationships among the sample versions are found. In Section 4 we will derive a measure for identifying observations that have a large influence on the test of fitting the equicorrelation model using the influence function method. It will be seen that even a single observation can greatly influence the result of this test. In Section 5 an illustrative example is given.

2. Preliminaries

Let $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ be a random sample from a p -variate normal distribution

$$f(\mathbf{x}) = |2\pi \Sigma|^{-1/2} \exp\{-(1/2)(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu)\},$$

where the covariance matrix Σ is of the form

$$\Sigma = \sigma^2 \{(1 - \rho)I_p + \rho \mathbf{1}_p \mathbf{1}_p^T\},$$

where $\sigma > 0$ and $0 < \rho < 1$. We write $\mathbf{1}_p$ as the p by 1 vector with all elements equal to 1.

Let $\bar{\mathbf{x}} = (1/n) \sum_{i=1}^n \mathbf{x}_i$ and $\mathbf{S} = (1/n) \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$. In what follows, we use the hat notation to denote its maximum likelihood estimator of a parameter. Then the maximum likelihood estimators of μ , σ^2 and ρ are given by

$$\begin{aligned} \hat{\mu} &= \bar{\mathbf{x}}, \\ \hat{\sigma}^2 &= \frac{1}{p} \text{tr}(\mathbf{S}), \\ \hat{\rho} &= \frac{2}{p(p-1) \hat{\sigma}^2} \sum_{i < j} s_{ij}, \end{aligned}$$

where s_{ij} is the (i, j) th element of \mathbf{S} . More details can be found in Rencher (1995).

3. Influence function

In this section we will derive influence functions for σ^2 and ρ and consider three sample versions of these influence functions that will be used for investigating the influence of observations on $\hat{\sigma}^2$ and $\hat{\rho}$. Some relationships among the sample versions are found.

Let F be a distribution function defined on the p -dimensional Euclidean space and

$\theta = \theta(F)$ be a parameter of interest which is a functional of F . The mean vector and covariance matrix for the distribution F are written as $\mu = \mu(F)$ and $\Sigma = \Sigma(F)$, respectively. For $0 \leq \varepsilon \leq 1$, the perturbation of F at \mathbf{x} is defined by $F_\varepsilon = (1 - \varepsilon)F + \varepsilon \delta_{\mathbf{x}}$, where $\delta_{\mathbf{x}}$ denotes the distribution having unit mass at \mathbf{x} . The perturbation of θ at \mathbf{x} is $\theta(F_\varepsilon)$. The influence function for θ at \mathbf{x} (Hampel, 1974) is defined by

$$\lim_{\varepsilon \rightarrow 0} \frac{\theta(F_\varepsilon) - \theta}{\varepsilon}. \quad (1)$$

The influence function for a parameter at \mathbf{x} measures the effect of an infinitesimal contamination at \mathbf{x} on the estimator of the parameter.

Consider the functionals $\sigma^2 = \sigma^2(F)$ and $\rho = \rho(F)$ defined by

$$\begin{aligned} \sigma^2(F) &= \frac{1}{p} \text{tr}(\Sigma(F)), \\ \rho(F) &= \frac{2}{p(p-1)\sigma^2(F)} \sum_{i < j} \sigma_{ij}(F), \end{aligned}$$

where $\sigma_{ij} = \sigma_{ij}(F)$ is the (i, j) th element of Σ . Let \hat{F} be the empirical distribution function based on the sample of size n . Then it is clear that $\sigma^2(\hat{F}) = \hat{\sigma}^2$ and $\rho(\hat{F}) = \hat{\rho}$ because $\mu(\hat{F}) = \bar{\mathbf{x}}$ and $\Sigma(\hat{F}) = \mathbf{S}$. Thus the functionals $\sigma^2(F)$ and $\rho(F)$ are well defined. Since the perturbation of Σ is

$$\Sigma(F_\varepsilon) = \Sigma + \varepsilon \{(\mathbf{x} - \mu)(\mathbf{x} - \mu)^T - \Sigma\} + O(\varepsilon^2),$$

the influence function for σ^2 is easily computed as

$$IF(\sigma^2; \mathbf{x}) = \frac{1}{p} \{(\mathbf{x} - \mu)^T(\mathbf{x} - \mu) - \text{tr}(\Sigma)\}. \quad (2)$$

Note that the influence function for σ_{ij} , $IF(\sigma_{ij}; \mathbf{x})$ is just the (i, j) th element of $(\mathbf{x} - \mu)(\mathbf{x} - \mu)^T - \Sigma$. By the Taylor series expansion, we have

$$\frac{1}{\sigma^2(F_\varepsilon)} = \frac{1}{\sigma^2} \left\{ 1 - \varepsilon \frac{IF(\sigma^2; \mathbf{x})}{\sigma^2} \right\} + O(\varepsilon^2).$$

Thus it is easily shown that the influence function for ρ is given by

$$\begin{aligned} IF(\rho; \mathbf{x}) &= \frac{2}{p(p-1)\sigma^2} \sum_{i < j} \left\{ IF(\sigma_{ij}; \mathbf{x}) - IF(\sigma^2; \mathbf{x}) \frac{\sigma_{ij}}{\sigma^2} \right\} \\ &= \frac{2}{p(p-1)\sigma^2} \sum_{i < j} (x_i - \mu_i)(x_j - \mu_j) - \frac{\rho}{p\sigma^2} (\mathbf{x} - \mu)^T(\mathbf{x} - \mu), \end{aligned} \quad (3)$$

where x_i and μ_i are the i th components of \mathbf{x} and μ , respectively.

Next we will consider three sample versions as in Critchley (1985): the empirical influence function (EIF), the sample influence function (SIF) and the deleted empirical influence function

(DIF). A large absolute value of each sample version indicates that the corresponding observation is influential.

3.1 Empirical influence function

The EIF is obtained by substituting the empirical distribution function \hat{F} for F in (2) and (3). It is easily seen from (2) that the EIF for $\hat{\sigma}^2$ at \mathbf{x}_u becomes

$$EIF(\hat{\sigma}^2; \mathbf{x}_u) = \frac{1}{p} (\mathbf{x}_u - \bar{\mathbf{x}})^T (\mathbf{x}_u - \bar{\mathbf{x}}) - \hat{\sigma}^2. \quad (4)$$

Let \bar{x}_i be the i th element of $\bar{\mathbf{x}}$. Using the relation that $2 \sum_{i < j} (x_i - \bar{x}_i)(x_j - \bar{x}_j)$

$= \{1_p^T (\mathbf{x} - \bar{\mathbf{x}})\}^2 - (\mathbf{x} - \bar{\mathbf{x}})^T (\mathbf{x} - \bar{\mathbf{x}})$, from (3) it is straightforward to show

$$EIF(\hat{\rho}; \mathbf{x}_u) = \frac{1}{p(p-1)\hat{\sigma}^2} [\{1_p^T (\mathbf{x}_u - \bar{\mathbf{x}})\}^2 - \{1 + (p-1)\hat{\rho}\}(\mathbf{x}_u - \bar{\mathbf{x}})^T (\mathbf{x}_u - \bar{\mathbf{x}})]. \quad (5)$$

3.2 Sample influence function

The SIF can be obtained by setting $F = \hat{F}$ and taking $\varepsilon = -1/(n-1)$ in the definition of the influence function (1) instead of taking a limit. Then the SIF for a parameter θ at \mathbf{x}_u can be rewritten as $(n-1)\{\theta(\hat{F}) - \theta(\hat{F}_{-u})\}$, where $\hat{F}_{-u} = (1 + (n-1)^{-1})\hat{F} - (n-1)^{-1}\delta_{\mathbf{x}_u}$ is the deleted version of \hat{F} with the u th observation \mathbf{x}_u deleted.

Let $\mathbf{S}_{-u} = \Sigma(\hat{F}_{-u})$. Since

$$\mathbf{S}_{-u} = \frac{n}{n-1} \mathbf{S} - \frac{n}{(n-1)^2} (\mathbf{x}_u - \bar{\mathbf{x}})(\mathbf{x}_u - \bar{\mathbf{x}})^T,$$

the SIF for $\hat{\sigma}^2$ at \mathbf{x}_u becomes

$$\begin{aligned} SIF(\hat{\sigma}^2; \mathbf{x}_u) &= \frac{n}{(n-1)p} (\mathbf{x}_u - \bar{\mathbf{x}})^T (\mathbf{x}_u - \bar{\mathbf{x}}) - \hat{\sigma}^2 \\ &= EIF(\hat{\sigma}^2; \mathbf{x}_u) + \frac{1}{(n-1)p} (\mathbf{x}_u - \bar{\mathbf{x}})^T (\mathbf{x}_u - \bar{\mathbf{x}}). \end{aligned} \quad (6)$$

The estimator of ρ without observation \mathbf{x}_u is computed as

$$\begin{aligned} \hat{\rho}_{-u} &= \frac{2}{p(p-1)\hat{\sigma}_{-u}^2} \sum_{i < j} s_{ij, -u} \\ &= \frac{2}{p(p-1)} \cdot \left[\frac{n-1}{n\hat{\sigma}^2} \left\{ 1 + \frac{\frac{1}{(n-1)p} (\mathbf{x}_u - \bar{\mathbf{x}})^T (\mathbf{x}_u - \bar{\mathbf{x}})}{\hat{\sigma}^2 - \frac{1}{(n-1)p} (\mathbf{x}_u - \bar{\mathbf{x}})^T (\mathbf{x}_u - \bar{\mathbf{x}})} \right\} \right] \\ &\quad \cdot \sum_{i < j} \left\{ \frac{n}{n-1} s_{ij} - \frac{n}{(n-1)^2} (x_{ui} - \bar{x}_i)(x_{uj} - \bar{x}_j) \right\}, \end{aligned}$$

where $\hat{\sigma}_{-u}^2 = \sigma^2(\hat{F}_{-u})$, x_{ui} is the i th element of \mathbf{x}_u and $s_{ij,-u}$ is the (i, j) th element of \mathbf{S}_{-u} . Thus the SIF for $\hat{\rho}$ at \mathbf{x}_u is easily obtained as

$$\begin{aligned} SIF(\hat{\rho}; \mathbf{x}_u) &= \left[p(p-1) \left\{ \hat{\sigma}_{-u}^2 - \frac{1}{(n-1)p} (\mathbf{x}_u - \bar{\mathbf{x}})^T (\mathbf{x}_u - \bar{\mathbf{x}}) \right\} \right]^{-1} \\ &\quad \cdot \left[\{1_p^T (\mathbf{x}_u - \bar{\mathbf{x}})\}^2 - (1 + (p-1)\hat{\rho})(\mathbf{x}_u - \bar{\mathbf{x}})^T (\mathbf{x}_u - \bar{\mathbf{x}}) \right] \\ &= \frac{\hat{\sigma}_{-u}^2}{\hat{\sigma}_{-u}^2 - \frac{1}{(n-1)p} (\mathbf{x}_u - \bar{\mathbf{x}})^T (\mathbf{x}_u - \bar{\mathbf{x}})} EIF(\hat{\rho}; \mathbf{x}_u). \end{aligned} \quad (7)$$

3.3 Deleted empirical influence function

The DIF is obtained by replacing F with \hat{F}_{-u} in (2) and (3) and it measures the effect of deleting the u th observation on the estimator. The mean vector for \hat{F}_{-u} is given by

$$\mu(\hat{F}_{-u}) = \bar{\mathbf{x}} - \frac{1}{n-1} (\mathbf{x}_u - \bar{\mathbf{x}})$$

and the covariance matrix for \hat{F}_{-u} is $\Sigma(\hat{F}_{-u}) = \mathbf{S}_{-u}$ computed in the previous subsection.

Replacing F with \hat{F}_{-u} in (2) yields the DIF for $\hat{\sigma}^2$ at \mathbf{x}_u as

$$\begin{aligned} DIF(\hat{\sigma}^2; \mathbf{x}_u) &= \frac{n(n+1)}{(n-1)^2 p} (\mathbf{x}_u - \bar{\mathbf{x}})^T (\mathbf{x}_u - \bar{\mathbf{x}}) - \frac{n}{n-1} \hat{\sigma}^2 \\ &= \frac{n}{n-1} EIF(\hat{\sigma}^2; \mathbf{x}_u) + \frac{2n}{(n-1)^2 p} (\mathbf{x}_u - \bar{\mathbf{x}})^T (\mathbf{x}_u - \bar{\mathbf{x}}). \end{aligned} \quad (8)$$

Similarly the DIF for $\hat{\rho}$ at \mathbf{x}_u is computed as

$$\begin{aligned} DIF(\hat{\rho}; \mathbf{x}_u) &= \frac{n^2}{p(n-1)^2 \hat{\sigma}_{-u}^2} \left[\frac{1}{p-1} \{1_p^T (\mathbf{x}_u - \bar{\mathbf{x}})\}^2 \right. \\ &\quad \left. - \left(\frac{1}{p-1} + \hat{\rho}_{-u} \right) (\mathbf{x}_u - \bar{\mathbf{x}})^T (\mathbf{x}_u - \bar{\mathbf{x}}) \right] \\ &= \frac{n}{n-1} \left[\frac{\hat{\sigma}_{-u}^2}{\hat{\sigma}_{-u}^2 - \frac{1}{(n-1)p} (\mathbf{x}_u - \bar{\mathbf{x}})^T (\mathbf{x}_u - \bar{\mathbf{x}})} \right]^2 EIF(\hat{\rho}; \mathbf{x}_u). \end{aligned} \quad (9)$$

3.4 Comparison of three sample versions

First consider the sample versions of the influence function for ρ given in (5), (7) and (9). If the sample size n is sufficiently large, then three sample versions yield similar values. Furthermore, we get the inequality that

$$|EIF(\hat{\rho}; \mathbf{x})| \leq |SIF(\hat{\rho}; \mathbf{x})| \leq |DIF(\hat{\rho}; \mathbf{x})|.$$

This inequality implies that DIF is more sensitive to the influence of an observation than the others.

Next the same reasoning applies to the case $\hat{\sigma}^2$ by comparing (4), (6) and (8). The positiveness of $EIF(\hat{\sigma}^2; \mathbf{x})$ implies that

$$EIF(\hat{\sigma}^2; \mathbf{x}) \leq SIF(\hat{\sigma}^2; \mathbf{x}) \leq DIF(\hat{\sigma}^2; \mathbf{x}).$$

Thus SIF and DIF tend to have larger values than EIF whenever $EIF(\hat{\sigma}^2; \mathbf{x})$ has a large positive value. In this case DIF is most sensitive to the influence of an observation. When $(\mathbf{x} - \bar{\mathbf{x}})^T(\mathbf{x} - \bar{\mathbf{x}})$ is near 0 and the sample size n is sufficiently large, three sample versions yield similar values.

4. Influence on a test of fitting the equicorrelation model

The test of the hypothesis

$$H_0: \Sigma = \sigma^2\{(1-\rho)I_p + \rho \mathbf{1}_p \mathbf{1}_p^T\}$$

is usually performed by the likelihood ratio statistic given by

$$T = \frac{|\mathbf{S}|}{\hat{\sigma}^{2p}(1-\hat{\rho})^{p-1}(1+(p-1)\hat{\rho})}.$$

Then

$$T_* = -\left[n-1 - \frac{p(p+1)^2(2p-3)}{6(p-1)(p^2+p-4)}\right] \ln T \quad (10)$$

is approximately distributed as a chi-squared distribution with $[p(p+1)/2-2]$ degrees of freedom. The null hypothesis H_0 would be rejected for a significantly large value of T_* (Rencher, 1995, p. 277).

To investigate the influence of observations on the likelihood ratio statistic T , we consider a functional $\eta = \eta(F)$ defined by

$$\eta(F) = \frac{|\Sigma(F)|}{\sigma^{2p}(F)(1-\rho(F))^{p-1}\{1+(p-1)\rho(F)\}}.$$

Then we have $\eta(\hat{F}) = T$. The perturbation of $|\Sigma(F)|$ is easily computed as

$$|\Sigma(F_\varepsilon)| = |\Sigma|[1 + \{(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu) - p\}\varepsilon + O(\varepsilon^2)].$$

For a functional θ , we have $\theta(F_\varepsilon) = \theta + IF(\theta; \mathbf{x})\varepsilon + O(\varepsilon^2)$. From (2) and (3), we get

$$\begin{aligned}
\sigma^{-2p}(F_\varepsilon) &= \sigma^{-2p} \left[1 - \frac{p}{\sigma^2} IF(\sigma^2; \mathbf{x}) \varepsilon + O(\varepsilon^2) \right], \\
(1 - \rho(F_\varepsilon))^{-(p-1)} &= (1 - \rho)^{-(p-1)} \left[1 + \frac{p-1}{1-\rho} IF(\rho; \mathbf{x}) \varepsilon + O(\varepsilon^2) \right], \\
(1 + (p-1)\rho(F_\varepsilon))^{-1} &= (1 + (p-1)\rho)^{-1} \left[1 - \frac{p-1}{1+(p-1)\rho} IF(\rho; \mathbf{x}) \varepsilon + O(\varepsilon^2) \right].
\end{aligned}$$

Thus we have

$$\eta(F_\varepsilon) = \eta(F) + IF(\eta; \mathbf{x}) \varepsilon + O(\varepsilon^2),$$

where

$$\begin{aligned}
IF(\eta; \mathbf{x}) &= \eta(F) \left[(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) - p - \frac{p}{\sigma^2} IF(\sigma^2; \mathbf{x}) \right. \\
&\quad \left. + \frac{p(p-1)\rho}{(1-\rho)(1+(p-1)\rho)} IF(\rho; \mathbf{x}) \right].
\end{aligned}$$

Here only EIF is considered as a sample version of the influence function $IF(\eta; \mathbf{x})$ and the other two sample versions can be derived as in Section 3. In general three sample versions described in Section 3 give similar results. Let $\hat{\eta} = \eta(\hat{F})$. Then the EIF for $\hat{\eta}$ at \mathbf{x}_u is given by

$$\begin{aligned}
EIF(\hat{\eta}; \mathbf{x}_u) &= T \left[(\mathbf{x}_u - \bar{\mathbf{x}})^T \mathbf{S}^{-1} (\mathbf{x}_u - \bar{\mathbf{x}}) - \frac{1}{(1-\hat{\rho})\hat{\sigma}^2} (\mathbf{x}_u - \bar{\mathbf{x}})^T (\mathbf{x}_u - \bar{\mathbf{x}}) \right. \\
&\quad \left. + \frac{\hat{\rho}}{(1-\hat{\rho})(1+(p-1)\hat{\rho})\hat{\sigma}^2} \{1_p^T (\mathbf{x}_u - \bar{\mathbf{x}})\}^2 \right]. \quad (11)
\end{aligned}$$

5. A numerical example

In this section, we consider the cost data which consists of 36 measurements on the per-mile cost of three variables: fuel, repair and capital. This data set is taken from Johnson and Wichern (1992, p. 276). When the covariance matrix does not have any structure, Bacon-Shone and Fung (1987) analyzed the data set and concluded that observations 9 and 21 are possible outliers.

First we will check whether this data set follows the equicorrelation model using the test statistic T_* in (10). The hypothesis that the data set follows the equicorrelation model would be rejected for a significantly large value of T_* . For the cost data, the value of T_* is 9.11 and the p -value is 0.058. Thus we would not reject the assumption of equicorrelation model at any significance level less than 0.058.

Next we will investigate the influence of observations on $\hat{\sigma}^2$ and $\hat{\rho}$ using three sample versions of the influence functions for σ^2 and ρ . The results for σ^2 is similar to those for ρ and therefore the index plot of σ^2 is not provided here. In Figure 1, y -axis indicates

EIF, SIF, DIF for $\hat{\rho}$. From Figure 1, we may conclude that observations 9 and 21 are most influential and that three sample versions yield the same result. Furthermore, we can see that the inequalities among three sample versions of the influence function in Section 3.4 are satisfied.

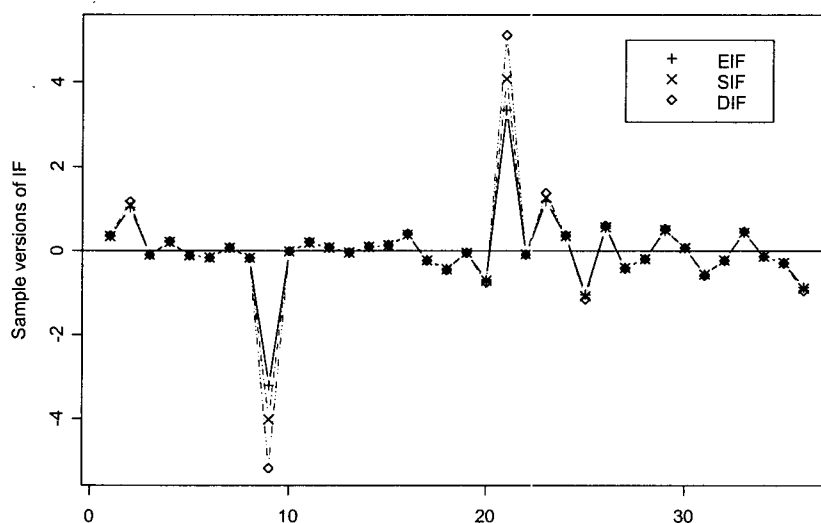


Figure 1 Sample versions of the influence function for ρ

The index plot of EIF in (11) for the test statistic T (or T_*), which is depicted in Figure 2, shows that only observation 9 is most influential on the test statistic. This result is supported by the case deletion of the test statistic T_* . The change in the value of the test statistic T_* due to a single case deletion is the maximum for the deletion of observation 9 and the next for the deletion of observation 21. After the deletion of observation 9, the p -value for the test statistic based on the remaining sample becomes 0.80 and thus the hypothesis that the data set follows the equicorrelation model would not be rejected at any reasonable significance level, whereas for the deletion of observation 21 the p -value is 0.09. It indicates that observation 9 has a large influence on the test but others do not.

From this example we can see that the influence function method yields useful information about the influence of observations on estimators or test statistics and that it can be a useful diagnostic tool.

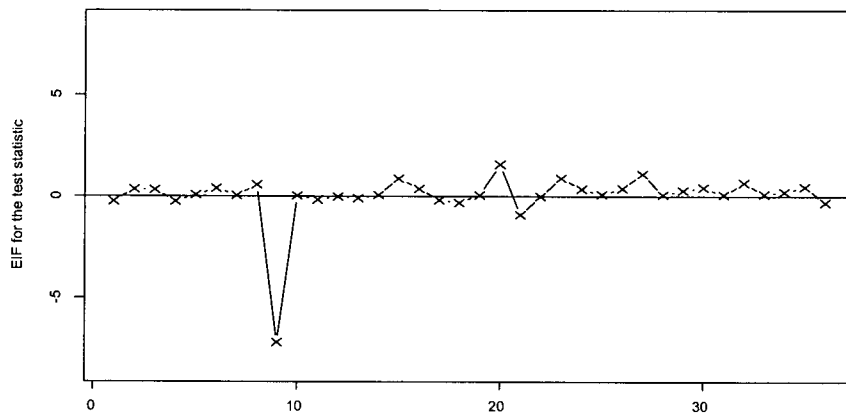


Figure 2 The empirical influence function for test statistic of equicorrelation model

References

- [1] Bacon-Shone, J. and Fung, W. K. (1987). A new graphical method for detecting single and multiple outliers in multivariate data, *Applied Statistics*, **36**, 153-162.
- [2] Critchley, F. (1985). Influence in principal component analysis, *Biometrika*, **72**, 627-636.
- [3] Giraudeau, B., Mallet, A. and Chastang, C. (1996). Case influence on the intraclass correlation coefficient estimate, *Biometrics*, **52**, 1492-1497.
- [4] Hampel, F. R. (1974). The influence curve and its role in robust estimation, *Journal of the American Statistical Association*, **69**, 383-393.
- [5] Johnson, A. J. and Wichern, D. W. (1992). *Applied Multivariate Statistical Analysis*, 3 ed., Prentice-Hall, Englewood Cliffs.
- [6] Jung, K.-M. and Kim, M. G. (1999). Influence analysis of the likelihood ratio test in multivariate Behrens-Fisher problem, *The Korean Communications in Statistics*, **6**, 939-947.
- [7] Rencher, A. C. (1995). *Methods of Multivariate Analysis*, Wiley, New York.