

# Multiple deletion diagnostics in beta regression models

Li-Chu Chien

Received: 9 October 2011 / Accepted: 8 September 2012 / Published online: 23 October 2012  
© Springer-Verlag Berlin Heidelberg 2012

**Abstract** We consider the problem of identifying multiple outliers in a general class of beta regression models proposed by Ferrari and Cribari-Neto (J Appl Stat 31:799–815, 2004). The currently available single-case deletion diagnostic measures, e.g., the standardized weighted residual (SWR), the Cook-like distance (LD), etc., often fail to identify multiple outlying observations, because they suffer from the well-known problems of masking and swamping effects. In this article, we develop group deletion diagnostic measures, such as generalized SWR, generalized LD, generalized DFFITS and generalized DFBETAS, and suggest a simple procedure for identifying multiple outliers using these. The performance of the proposed methods is investigated through simulation studies and two practical examples.

**Keywords** Beta regression · Multiple outliers · Generalized SWR · Generalized LD · Generalized DFFITS · Generalized DFBETAS

## 1 Introduction

Many fields of studies involve data in the form of percentages, rates or proportions that are measured continuously in the open interval  $(0, 1)$ . For example, one may be interested in modeling the proportion of income spent on food as a function of the level of income and the number of persons in the household. The beta distribution is a flexible and useful tool for modeling data on the standard unit interval  $(0, 1)$ , since the beta density can display quite different shapes depending on the values of the parameters that index the distribution; see, for example, [Kieschnick and McCullough \(2003\)](#).

---

L.-C. Chien (✉)  
Institute of Statistics, National Chiao Tung University, Hsinchu, Taiwan  
e-mail: lcchien@nctu.edu.tw

Ferrari and Cribari-Neto (2004) proposed a class of beta regression models which is the class of models derived from generalized linear models (GLMs), except the response variable is not from a linear exponential family distribution. Under generalized beta linear models (GBLMs), they developed maximum likelihood inference including parameter and interval estimation, and also hypothesis tests. They provided complete inference tools for the new class of models. The tools are freely available in the R package *betareg*. See Cribari-Neto and Zeileis (2010) for details. Hence, the execution of the beta regression techniques in practical problems is convenient.

Identifying observations that may affect the results of a regression analysis is a fundamental step in regression model building processes. In general, observations that lie well outside the majority of the data are termed outliers, in the sense that outliers come from a different probability distribution or from a different deterministic model than the mass of the data. The existence of outliers always distorts the outcome and accuracy of regression results, and hence, outliers must be detected in the regression analysis processes.

In GBLMs, to identify the outlying observations that depart from the postulated model of the bulk of the data, some authors have provided guidelines for diagnostic analysis. For example, Ferrari and Cribari-Neto (2004) provided some diagnostic measures to identify atypical observations and to detect model misspecification. Espinheira et al. (2008a) proposed two new beta residuals and numerically compared their behavior to those originally suggested by Ferrari and Cribari-Neto (2004). The results indicate a preference for one of the new residuals, more specifically the residual that accounts for the different leverages of the observations. On the other hand, Espinheira et al. (2008b) developed the Cook-like distance (LD) to measure the effects of influential observations on regression parameter estimates of GBLMs.

These currently available outlier measures and influence diagnostics seem to be available only when the data merely contain a single outlier. However, if the data contain more than one outlying observation, these existing methods may become ineffective, due to the problems of masking and swamping effects. Hence, multiple outlier detection methods that are free from these problems are proposed in this article.

This article unfolds as follows. Section 2 contains a concise review of GLBMs proposed by Ferrari and Cribari-Neto (2004). In the next section, we briefly introduce some of the current diagnostic tools, e.g., the standardized weighted residual (SWR), LD, etc. We also define the GBLM versions of the influence measures DFFITS and DFBETAS in this section. In the section after, we introduce SWR, LD, DFFITS and DFBETAS based on group deletion techniques and suggest an easy procedure for detecting multiple outliers using these group deletion diagnostics. Sections 5 and 6 illustrate applications of these newly proposed deletion diagnostic methods in simulated and real data examples, respectively. Finally, in Sect. 7, some conclusions about these proposed diagnostic measures are set out.

## 2 Beta regression model

The probability function for a single response variable  $Y$  in a GBLM (Ferrari and Cribari-Neto 2004) is

$$f(y; \mu, \phi) = \frac{\Gamma(\phi)}{\Gamma(\mu\phi)\Gamma((1-\mu)\phi)} y^{\mu\phi-1} (1-y)^{(1-\mu)\phi-1}, \quad 0 < y < 1 \quad (1)$$

where  $0 < \mu < 1$ ,  $\phi > 0$  and  $\Gamma(\cdot)$  is the gamma function. The mean and variance of  $Y$  are  $E(Y) = \mu$  and  $Var(Y) = V(\mu)/(1 + \phi)$ , respectively, where  $V(\mu) = \mu(1 - \mu)$  is a variance function. Note that  $\phi$  can be viewed as a precision parameter, in the sense that, for a fixed mean  $\mu$ , the variance of  $Y$  decreases as  $\phi$  increases. Here our concern is with the mean parameter  $\mu$ , and  $\phi$  may be viewed as a nuisance parameter.

Now consider observations  $y_1, \dots, y_n$  which are regarded as realizations of independent random variables  $Y_1, \dots, Y_n$  and each  $y_i, i = 1, \dots, n$ , follows the density (1) with mean  $\mu_i$  and unknown precision  $\phi$ . The mean of  $Y_i$  involves explanatory variables through a link function  $g(\cdot)$ , so that  $g(\mu_i) = \eta_i$  where  $\eta_i = \mathbf{x}_i^T \boldsymbol{\beta} = \beta_1 x_{i,1} + \beta_2 x_{i,2} + \dots + \beta_p x_{i,p}$ . Here  $\eta_i$  is a linear predictor,  $\mathbf{x}_i = (x_{i,1}, \dots, x_{i,p})^T$  is a  $p$ -vector of explanatory variables,  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^T$  is a  $p$ -vector of unknown regression parameters, and  $g(\cdot)$  is a strictly monotonic and twice differentiable link function that forms a mapping from the interval  $(0, 1)$  to  $\mathbb{R}$ . The parameters  $\boldsymbol{\beta}$  and  $\phi$  can be estimated by maximum likelihood methods that can easily be implemented through the R package *betareg*. For details, see [Cribari-Neto and Zeileis \(2010\)](#).

### 3 Measures of influence based on single-case deletion diagnostics

In this section, we succinctly review the currently available single-case deletion diagnostic methods. We also define the GBLM versions of the influence measures DFFITS and DFBETAS according to the classical versions of DFFITS and DFBETAS established under the normal regression settings and discussed in Subsection 2.1 of [Belsley et al. \(1980\)](#).

#### 3.1 Assessing the influence of case deletion on residuals

This subsection is focused on the single-case outlier identification tools which use the residuals to detect the atypical observations and check model adequacy. We review the weighted residual (WR) and SWR proposed by [Espinheira et al. \(2008a\)](#).

*WR and SWR.* [Espinheira et al. \(2008a\)](#) defined WR for point  $i$  as

$$r_i^* = \frac{y_i^* - \hat{\mu}_i^*}{\sqrt{\hat{\phi} \hat{v}_i}} \quad (2)$$

where  $y_i^* = \log(y_i/(1 - y_i))$ ,  $\hat{\mu}_i^* = \psi(\hat{\mu}_i \hat{\phi}) - \psi((1 - \hat{\mu}_i) \hat{\phi})$ ,  $\hat{v}_i = \{\psi'(\hat{\mu}_i \hat{\phi}) + \psi'((1 - \hat{\mu}_i) \hat{\phi})\}$  and  $\psi'(\cdot)$  is the trigamma function. Here the symbol “ $\hat{\cdot}$ ” is used for quantities that are evaluated at parameter values of the maximum likelihood solution based on all observations. Then the standardized version, SWR, for point  $i$  is defined by

$$r_i^{ww} = \frac{r_i^*}{\sqrt{(1 - \hat{h}_i)/\hat{\phi}}} = \frac{y_i^* - \hat{\mu}_i^*}{\sqrt{\hat{v}_i (1 - \hat{h}_i)}} \quad (3)$$

where  $\widehat{h}_i$  is the  $i$ th diagonal element of  $\widehat{H} = \widehat{W}^{1/2} X(X^T \widehat{W} X)^{-1} X^T \widehat{W}^{1/2}$  and  $\widehat{W}^{1/2}$  is a symmetric square root of  $\widehat{W}$ . Here  $X = (x_1, \dots, x_n)^T$  and  $\widehat{W} = \widehat{\phi} \widehat{G} \widehat{V} \widehat{G}$  with diagonal matrices  $\widehat{G} = \text{Diag}(1/g'(\widehat{\mu}_1), \dots, 1/g'(\widehat{\mu}_n))$  and  $\widehat{V} = \text{Diag}(\widehat{v}_1, \dots, \widehat{v}_n)$ . An excessively large or small value of  $WR_i$  or  $SWR_i$  indicates point  $i$  having an unusual residual.

### 3.2 Assessing the influence of case deletion on regression parameter estimates

This subsection is concerned with the single-case influence identification tools which measure the influence of deleting one observation on the regression parameter estimates or on the fitted values. We review LD proposed by [Espinheira et al. \(2008b\)](#) and propose the GBLM versions of the influence measures DFFITS and DFBETAS.

LD. [Espinheira et al. \(2008b\)](#) defined LD for point  $i$  as

$$LD_i = (\widehat{\beta} - \widehat{\beta}^{(-i)})^T (\widehat{\phi} X^T \widehat{W} X) (\widehat{\beta} - \widehat{\beta}^{(-i)})$$

where  $\widehat{\beta}^{(-i)}$  is the maximum likelihood estimator (MLE) of  $\beta$  without the  $i$ th observation. From the approximate relation

$$\widehat{\beta}^{(-i)} \approx \widehat{\beta} - \frac{(X^T \widehat{W} X)^{-1} x_i \widehat{w}_i^{1/2}}{1 - \widehat{h}_i} r_i^* \tag{4}$$

where  $\widehat{w}_i^{1/2}$  is the  $i$ th diagonal entry of  $\widehat{W}^{1/2}$ , it follows that

$$LD_i \approx \left( \frac{r_i^*}{1 - \widehat{h}_i} \right)^2 \widehat{\phi} \widehat{h}_i = (r_i^{ww})^2 \frac{\widehat{h}_i}{1 - \widehat{h}_i}. \tag{5}$$

A large  $LD_i$  indicates that point  $i$  has a large impact on  $\widehat{\beta}$ .

DFFITS. We define DFFITS for point  $i$  as

$$DFFITS_i = \frac{\widehat{w}_i^{1/2} x_i^T (\widehat{\beta} - \widehat{\beta}^{(-i)})}{\sqrt{\widehat{h}_i / \widehat{\phi}^{(-i)}}}$$

which, in view of (4), is expressed by

$$DFFITS_i \approx \frac{r_i^*}{\sqrt{1 - \widehat{h}_i}} \frac{\sqrt{\widehat{h}_i \widehat{\phi}^{(-i)}}}{\sqrt{1 - \widehat{h}_i}} = r_i^{ww} \sqrt{\frac{\widehat{\phi}^{(-i)}}{\widehat{\phi}}} \sqrt{\frac{\widehat{h}_i}{1 - \widehat{h}_i}} \tag{6}$$

where  $\widehat{\phi}^{(-i)}$  is the MLE of  $\phi$  with observation  $i$  omitted. Using (6), we obtain, from (5), that

$$LD_i = (\text{DFFITS}_i)^2 \frac{\widehat{\phi}}{\widehat{\phi}^{(-i)}}.$$

An extreme  $\text{DFFITS}_i$  implies that observation  $i$  is influential on the weighted fit,  $\widehat{w}_i^{1/2} \mathbf{x}_i^T \widehat{\beta}$ .

*DFBETAS*. We define *DFBETAS* for point  $i$  for  $\beta_j$  as

$$\text{DFBETAS}_{ij} = \frac{\widehat{\beta}_j - \widehat{\beta}_j^{(-i)}}{\sqrt{(X^T \widehat{W} X)_j^{-1} / \widehat{\phi}^{(-i)}}}$$

where  $(X^T \widehat{W} X)_j^{-1}$  is the  $j$ th diagonal element of the inverse of  $X^T \widehat{W} X$ ,  $(X^T \widehat{W} X)^{-1}$ . Let  $c_{ji}$  be the  $j$ th component of the  $p$ -vector  $(X^T \widehat{W} X)^{-1} \mathbf{x}_i \widehat{w}_i^{1/2}$ . We then have

$$\text{DFBETAS}_{ij} \approx r_i^* \frac{\sqrt{\widehat{\phi}^{(-i)}}}{1 - \widehat{h}_i} \frac{c_{ji}}{\sqrt{\sum_{i=1}^n c_{ji}^2}} = r_i^{ww} \sqrt{\frac{\widehat{\phi}^{(-i)}}{\widehat{\phi}(1 - \widehat{h}_i)}} \frac{c_{ji}}{\sqrt{\sum_{i=1}^n c_{ji}^2}}. \tag{7}$$

A large absolute value of  $\text{DFBETAS}_{ij}$  shows that observation  $i$  is influential on  $\widehat{\beta}_j$ .

### 4 Measures of influence based on group deletion diagnostics

In this section, we introduce the group deletion diagnostic measures and suggest a diagnostic procedure for detecting multiple outlying observations using these. We assume that  $d$  observations among a set of  $n$  observations are unusual observations and omitted before the fitting of the model. Let  $[R]$  index a set of the  $(n - d)$  observations that are remaining in the analysis after deleting a set of the  $d$  unusual observations indexed by  $D$ .

Let  $\mathbf{y}^* = (y_1^*, \dots, y_n^*)^T$ ,  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_n)^T$  and  $\boldsymbol{\mu}^* = (\mu_1^*, \dots, \mu_n^*)^T$  be the  $n \times 1$  vectors. Without loss of generality, we assume that observations to be deleted are the last  $d$  components of  $\mathbf{y}^*$  and  $X$ , so that  $\mathbf{y}^{*T} = (\mathbf{y}_{[R]}^{*T}, \mathbf{y}_D^{*T})$  and  $X^T = (X_{[R]}^T, X_D^T)$ . Let  $\widehat{\beta}_{[R]}$ ,  $\widehat{\phi}_{[R]}$ ,  $\widehat{\boldsymbol{\mu}}_{[R]}$ ,  $\widehat{\boldsymbol{\mu}}_{[R]}^*$ ,  $\widehat{V}_{[R]}$ ,  $\widehat{G}_{[R]}$  and  $\widehat{W}_{[R]}$  be the MLEs, respectively, of  $\beta$ ,  $\phi$ ,  $\boldsymbol{\mu}$ ,  $\boldsymbol{\mu}^*$ ,  $V$ ,  $G$  and  $W$  without the  $d$  observations in the deletion set  $D$ . Let  $\widehat{W}_{[R]}^{-1}$  be the inverse of  $\widehat{W}_{[R]}$ . Then write

$$\mathbf{z}_{[R]} = X \widehat{\beta}_{[R]} + \widehat{W}_{[R]}^{-1} \widehat{G}_{[R]} (\mathbf{y}^* - \widehat{\boldsymbol{\mu}}_{[R]}^*) \tag{8}$$

and partition  $\widehat{W}_{[R]}$  as

$$\widehat{W}_{[R]} = \widehat{\phi}_{[R]} \widehat{G}_{[R]} \widehat{V}_{[R]} \widehat{G}_{[R]} = \begin{bmatrix} \widehat{W}_{[R][R]} & \mathbf{0} \\ \mathbf{0}^T & \widehat{W}_{D[R]} \end{bmatrix}$$

with diagonal matrices  $\widehat{W}_{[R][R]}$  and  $\widehat{W}_{D[R]}$  being related to observations from the  $[R]$  and  $D$  sets, respectively, in which  $\theta$  is a zero matrix of order  $(n - d) \times d$ . When a group of observations  $D$  is omitted, we define the leverage measure for case  $i$  as  $\widehat{h}_{i[R]} = \widehat{w}_{i[R]}^{1/2} \mathbf{x}_i^T (X_{[R]}^T \widehat{W}_{[R][R]} X_{[R]})^{-1} \mathbf{x}_i \widehat{w}_{i[R]}^{1/2}$  where  $\widehat{w}_{i[R]}^{1/2}$  is the  $i$ th element of  $\widehat{W}_{[R]}^{1/2}$  that is a symmetric square root of  $\widehat{W}_{[R]}$ .

#### 4.1 Identifying multiple outliers using group-deleted versions of WR and SWR

The generalized WR (GWR) and the generalized SWR (GSWR), based on group deletion techniques, are proposed by the single-case outlier detection measures WR and SWR.

*GWR.* For point  $i$  from the  $[R]$  set, WR is defined, in view of (2), as

$$r_{i[R]}^* = \frac{y_i^* - \widehat{\mu}_{i[R]}^*}{\sqrt{\widehat{\phi}_{[R]} \widehat{v}_{i[R]}}} \tag{9}$$

where  $\widehat{\mu}_{i[R]}^* = \psi(\widehat{\mu}_{i[R]} \widehat{\phi}_{[R]}) - \psi((1 - \widehat{\mu}_{i[R]}) \widehat{\phi}_{[R]})$  and  $\widehat{v}_{i[R]}$  is the  $i$ th diagonal unit of  $\widehat{V}_{[R]}$ . Using (8),  $r_{i[R]}^*$  is expressed in the alternative form by

$$r_{i[R]}^* = \widehat{w}_{i[R]}^{1/2} \left( z_{i[R]} - \mathbf{x}_i^T \widehat{\beta}_{[R]} \right) \tag{10}$$

which is called the internal WR (IWR), because the point  $i$  is from the  $[R]$  set.

In view of (10), WR for point  $i$  from the  $D$  set is defined as  $r_{i[R+i]}^* = \widehat{w}_{i[R+i]}^{1/2} (z_{i[R+i]} - \mathbf{x}_i^T \widehat{\beta}_{[R+i]})$  where notations “ $\widehat{\phantom{x}}$ ” and “[ $R + i$ ]” are used for quantities that are evaluated at parameter values of the maximum likelihood solution based on the remaining set  $[R + i]$ , which consists of  $(n - d + 1)$  observations with the point  $i$  from the  $D$  set and the others from the  $[R]$  set. Using

$$\widehat{\beta}_{[R+i]} \approx \widehat{\beta}_{[R]} + \frac{\left( X_{[R]}^T \widehat{W}_{[R][R]} X_{[R]} \right)^{-1} \mathbf{x}_i \widehat{w}_{i[R]}^{1/2}}{1 + \widehat{h}_{i[R]}} r_{i[R]}^*, \tag{11}$$

$r_{i[R+i]}^*$  is equivalently expressed in the form

$$r_{i[R+i]}^* = \frac{r_{i[R]}^*}{1 + \widehat{h}_{i[R]}} \tag{12}$$

which is called the external WR (EWR), because the point  $i$  is from the  $D$  set. Then  $GWR_i$  for any data points is defined by (9) for  $i$  in  $[R]$  and by (12) for  $i$  in  $D$ .

*GSWR.* For point  $i$  from the  $[R]$  set, the internal SWR (ISWR) is defined, using (3), as

$$I r_{i[R]}^{sw} = \frac{r_{i[R]}^*}{\sqrt{(1 - \widehat{h}_{i[R]}) / \widehat{\phi}_{[R]}}} . \tag{13}$$

On the other hand, for point  $i$  from the  $D$  set, the external SWR (ESWR) is defined, using (13), as  $Er_{i[R+i]}^{ww} = r_{i[R+i]}^* / \sqrt{(1 - \widehat{h}_{i[R+i]}) / \widehat{\phi}_{[R+i]}}$ . By writing

$$\widehat{h}_{i[R+i]} = \widehat{w}_{i[R]}^{1/2} \mathbf{x}_i^T \left( X_{[R]}^T \widehat{W}_{[R][R]} X_{[R]} + \mathbf{x}_i \widehat{w}_{i[R]} \mathbf{x}_i^T \right)^{-1} \mathbf{x}_i \widehat{w}_{i[R]}^{1/2} = \frac{\widehat{h}_{i[R]}}{1 + \widehat{h}_{i[R]}} \tag{14}$$

and using (12), ESWR is seen to be equal to

$$Er_{i[R+i]}^{ww} = \frac{r_{i[R]}^*}{\sqrt{(1 + \widehat{h}_{i[R]}) / \widehat{\phi}_{[R+i]}}} \tag{15}$$

Then  $GSWR_i$  for any data point is defined by (13) for  $i$  in  $[R]$  and by (15) for  $i$  in  $D$ .

In fact, group deletion residuals for both the  $[R]$  set and the  $D$  set, which are measured based on a similar scale, have received a great deal of attention in the literature in recent years. For example, similar residuals were derived by Hadi and Simonoff (1993) for detecting multiple outliers in linear models and, later, a similar diagnostic technique was sketched by Atkinson (1994). In addition, this type of residual was introduced, under linear regression and logistic regression, by Imon (2005) and Imon and Hadi (2008).

#### 4.2 Identifying multiple outliers using group-deleted versions of LD, DFFITS and DFBETAS

The generalized LD (GLD), the generalized DFFITS (GDFFITS) and the generalized DFBETAS (GDFBETAS), based on group deletion techniques, are proposed by the single-case influence detection measures LD, DFFITS and DFBETAS.

*GLD.* For point  $i$  from the  $[R]$  set, the internal LD is computed by

$$ILD_i = \left( \widehat{\beta}_{[R]} - \widehat{\beta}_{[R]}^{(-i)} \right)^T \left( \widehat{\phi}_{[R]} X_{[R]}^T \widehat{W}_{[R][R]} X_{[R]} \right) \left( \widehat{\beta}_{[R]} - \widehat{\beta}_{[R]}^{(-i)} \right)$$

whereas, for point  $i$  from the  $D$  set, the external LD is computed by

$$ELD_i = \left( \widehat{\beta}_{[R+i]} - \widehat{\beta}_{[R]} \right)^T \left( \widehat{\phi}_{[R+i]} X_{[R+i]}^T \widehat{W}_{[R+i][R]} X_{[R+i]} \right) \left( \widehat{\beta}_{[R+i]} - \widehat{\beta}_{[R]} \right)$$

where  $\widehat{W}_{[R+i][R]}$  is an  $(n-d+1) \times (n-d+1)$  diagonal matrix, in which the first  $(n-d)$  diagonal units are relative to observations from the  $[R]$  set and the last diagonal unit is relative to the point  $i$  from the  $D$  set.  $GLD_i$  is derived by combining  $ILD_i$  and  $ELD_i$ . Through (11) and using  $X_{[R+i]}^T \widehat{W}_{[R+i][R]} X_{[R+i]} = X_{[R]}^T \widehat{W}_{[R][R]} X_{[R]} + \mathbf{x}_i \widehat{w}_{i[R]} \mathbf{x}_i^T$  and

$$\widehat{\beta}_{[R]} \approx \widehat{\beta}_{[R]}^{(-i)} + \frac{\left( X_{[R]}^T \widehat{W}_{[R][R]} X_{[R]} \right)^{-1} \mathbf{x}_i \widehat{w}_{i[R]}^{1/2}}{1 - \widehat{h}_{i[R]}} r_{i[R]}^* \tag{16}$$

GLD<sub>*i*</sub> is expressed in terms of  $Ir_{i[R]}^{ww}$  and  $\widehat{h}_{i[R]}$  as

$$GLD_i = \begin{cases} (Ir_{i[R]}^{ww})^2 \frac{\widehat{h}_{i[R]}}{1-\widehat{h}_{i[R]}} & \text{if } i \in [R] \\ (Ir_{i[R]}^{ww})^2 \frac{(1-\widehat{h}_{i[R]})\widehat{h}_{i[R]} \widehat{\phi}_{[R+i]}}{1+\widehat{h}_{i[R]} \widehat{\phi}_{[R]}} & \text{if } i \notin [R]. \end{cases}$$

*GDFFITs*. The internal DFFITS (IDFFITS), for point *i* from the [R] set, is defined by

$$IDFFITS_i = \frac{\widehat{w}_{i[R]}^{1/2} \mathbf{x}_i^T (\widehat{\beta}_{[R]} - \widehat{\beta}_{[R]}^{(-i)})}{\sqrt{\widehat{h}_{i[R]} / \widehat{\phi}_{[R]}^{(-i)}}}$$

whereas the external DFFITS (EDFFITS), for point *i* from the **D** set, is defined by

$$EDFFITS_i = \frac{\widehat{w}_{i[R]}^{1/2} \mathbf{x}_i^T (\widehat{\beta}_{[R+i]} - \widehat{\beta}_{[R]})}{\sqrt{\widehat{h}_{i[R+i]} / \widehat{\phi}_{[R]}}}.$$

Together these IDFFITS<sub>*i*</sub> and EDFFITS<sub>*i*</sub> give GDFFITS<sub>*i*</sub> for point *i* in the entire data set. Using (16) in conjunction with (11) and (14), it turns out that GDFFITS<sub>*i*</sub> is displayed in terms of  $Ir_{i[R]}^{ww}$  and  $\widehat{h}_{i[R]}$  as

$$GDFFITS_i = \begin{cases} Ir_{i[R]}^{ww} \frac{\sqrt{\widehat{\phi}_{[R]}^{(-i)} \widehat{h}_{i[R]}}}{\sqrt{\widehat{\phi}_{[R]} (1-\widehat{h}_{i[R]})}} & \text{if } i \in [R] \\ Ir_{i[R]}^{ww} \frac{\sqrt{(1-\widehat{h}_{i[R]}) \widehat{h}_{i[R]}}}{\sqrt{1+\widehat{h}_{i[R]}}} & \text{if } i \notin [R]. \end{cases}$$

After some algebraic manipulation, it is observed

$$GLD_i = \begin{cases} (GDFFITS_i)^2 \frac{\widehat{\phi}_{[R]}}{\widehat{\phi}_{[R]}^{(-i)}} & \text{if } i \in [R] \\ (GDFFITS_i)^2 \frac{\widehat{\phi}_{[R+i]}}{\widehat{\phi}_{[R]}} & \text{if } i \notin [R]. \end{cases}$$

*GDFBETAS*. For point *i* from the [R] set, we define the internal DFBETAS (IDFBETAS), for  $\beta_j$ , as

$$IDFBETAS_{ij} = \frac{\widehat{\beta}_{j[R]} - \widehat{\beta}_{j[R]}^{(-i)}}{\sqrt{\left(\mathbf{X}_{[R]}^T \widehat{\mathbf{W}}_{[R][R]} \mathbf{X}_{[R]}\right)_j^{-1} / \widehat{\phi}_{[R]}^{(-i)}}}$$

For point *i* from the **D** set, we define the external DFBETAS (EDFBETAS), for  $\beta_j$ , as

$$EDFBETAS_{ij} = \frac{\widehat{\beta}_{j[R+i]} - \widehat{\beta}_{j[R]}}{\sqrt{\left(\mathbf{X}_{[R+i]}^T \widehat{\mathbf{W}}_{[R+i][R]} \mathbf{X}_{[R+i]}\right)_j^{-1} / \widehat{\phi}_{[R]}}}.$$



Collecting these IDFBETAS<sub>ij</sub> and EDFBETAS<sub>ij</sub>, for any point *i* in the data set, GDFBETAS, for β<sub>*j*</sub>, is exhibited in terms of *I*r<sub>*i*[**R**]</sub><sup>ww</sup> and  $\widehat{h}_{i[\mathbf{R}]}$  as

$$GDFBETAS_{ij} = \begin{cases} I r_{i[\mathbf{R}]}^{ww} \frac{\sqrt{\widehat{\phi}_{[\mathbf{R}]}}^{(-i)}}{\sqrt{\widehat{\phi}_{[\mathbf{R}]}(1-\widehat{h}_{i[\mathbf{R}]})}} \frac{c_{ji[\mathbf{R}]}}{\sqrt{\sum_{i \in \mathbf{R}} c_{ji[\mathbf{R}]}^2}} & \text{if } i \in [\mathbf{R}] \\ I r_{i[\mathbf{R}]}^{ww} \frac{\sqrt{1-\widehat{h}_{i[\mathbf{R}]}}}{1+\widehat{h}_{i[\mathbf{R}]}} \frac{c_{ji[\mathbf{R}]}}{\sqrt{\sum_{i \in \mathbf{R}} c_{ji[\mathbf{R}]}^2 - (c_{ji[\mathbf{R}]}^2 / (1+\widehat{h}_{i[\mathbf{R}]})}} & \text{if } i \notin [\mathbf{R}] \end{cases}$$

where *c*<sub>*ji*[**R**]</sub> is the *j*th component of the *p*-vector  $(\mathbf{X}_{[\mathbf{R}]}^T \widehat{\mathbf{W}}_{[\mathbf{R}]} \mathbf{X}_{[\mathbf{R}]})^{-1} \mathbf{x}_i \widehat{w}_i^{1/2}$ .

Similar diagnostic measures are derived by Imon (2005) who proposed the group deletion versions of LD and DFFITS for the identification of multiple influential observations in linear regression. Nurunnabi et al. (2010) developed the generalized DFFITS in logistic regression for the same purpose.

### 4.3 The diagnostic algorithm

In this subsection, attention is given to the introduction of a test procedure for the detection of multiple outliers through the group deletion diagnostic measures, GSWR<sub>*i*</sub>, GLD<sub>*i*</sub>, GDFFITs<sub>*i*</sub> and GDFBETAS<sub>*i*</sub>, respectively. The main idea of the proposed method is to first form a basic subset of one fourth of the data which is possibly free from potential outliers and then employ the group deletion diagnostics, based on the basic subset, in identifying outlying observations. The detailed diagnostic algorithm is illustrated with GSWR<sub>*i*</sub> as follows.

*Step 0* Fit the regression model to the full data and compute *r*<sub>*i*</sub><sup>ww</sup> for *i* = 1, . . . , *n*.

*Step 1* Arrange the *n* observations in ascending order according to the absolute values of *r*<sub>*i*</sub><sup>ww</sup>, |*r*<sub>*i*</sub><sup>ww</sup>|, *i* = 1, . . . , *n*. Then the first {*n*/4} observations form the original [**R**] set, where {*n*/4} is the integer part of *n*/4 and represents the initial size of the [**R**] set. If the initial subset is not of full rank, increase the initial subset by as many observations as needed for the initial subset to become full rank (the observations are added according to their ranked order).

*Step 2* Fit a regression model to the current [**R**] set and compute GSWR<sub>*i*</sub>.

*Step 3* Arrange observations in accordance with an increasing order of |GSWR<sub>*i*</sub>| and let GSWR<sub>(*s*+1)</sub> be the (*s* + 1)th order statistic of |GSWR<sub>*i*</sub>|, where *s* is the current size of the [**R**] set.

- (a) If GSWR<sub>(*s*+1)</sub> > 2*c*(*s*) with  $c(s) = \sqrt{\widehat{\phi}_{[\mathbf{R}]} / \widehat{\phi}} \sqrt{\sum_{i=1}^n |\text{GSWR}_i| / \sum_{i=1}^n |\text{SWR}_i|}$ , then declare all members satisfying |GSWR<sub>*i*</sub>| > 2*c*(*s*) as outliers, and stop.
- (b) Otherwise, the current [**R**] set is replaced by the first (*s* + 1) ordered observations. If *s* + 1 = *n*, then go to Step 4; otherwise return to Step 2.

*Step 4* GSWR<sub>*i*</sub> is recalculated with all data points involved into the [**R**] set. Then declare all observations satisfying |GSWR<sub>*i*</sub>| > 2 as outliers and stop.

Here we suggest to use the cut-off values of 2 and 2*c*(*s*) for |SWR<sub>*i*</sub>| and |GSWR<sub>*i*</sub>|, respectively. That is because, SWR<sub>*i*</sub> and GSWR<sub>*i*</sub> are standardized quantities. Hence there will be very few SWR<sub>*i*</sub>'s larger than 2 and GSWR<sub>*i*</sub>'s larger than 2*c*(*s*), in which

$c(s)$  plays the role as an adjustment, which adjusts the effect from the differences between the sample size,  $n$ , and the current size of the  $\mathbf{R}$  set,  $s$ . It should be explained that, in practice, the criteria of  $|\text{SWR}_i|$  and  $|\text{GSWR}_j|$  can be flexibly defined by other cut-off values, e.g., 2.5 and  $2.5c(s)$ , 3 and  $3c(s)$ , etc.

Note that if there are a few outliers involved in the initially basic subset of size  $\{n/4\}$ , these outliers are still revealed as outlying because these observations should be gradually excluded from the subsets with sizes  $\{n/4\} + 1, \{n/4\} + 2, \dots$ , respectively.

The diagnostic algorithms for  $\text{GLD}_i, \text{GDFFIT}_i$  and  $\text{GDFBETAS}_{ij}$  are similar, except the cut-off values, and hence they are omitted. Their respective suggested cut-off values are as follows.

Under the assumption that the weighted matrix  $\widehat{\mathbf{W}}^{1/2}\mathbf{X}$  is of full rank, the average of  $\widehat{h}_i$  is  $p/n$ , and then  $\text{LD}_i$  in (5) becomes

$$\text{LD}_i = (r_i^{ww})^2 \frac{p}{n - p}.$$

Hence we suggest to use the cut-off values of  $4p/n$  and  $4c^2(s)p/s$  for  $\text{LD}_i$  and  $\text{GLD}_i$ , respectively. Under the same assumption,  $\text{DFFIT}_i$  reduces from (6) to

$$\text{DFFIT}_i = r_i^{ww} \sqrt{\frac{\widehat{\phi}^{(-i)}}{\widehat{\phi}}} \sqrt{\frac{p}{n - p}}.$$

Hence we suggest to use the cut-off values of  $2\sqrt{p/n}$  and  $2c(s)\sqrt{p/s}$  for  $|\text{DFFIT}_i|$  and  $|\text{GDFFIT}_i|$ , respectively.

On the other hand, in the special case of location, i.e., in the special case with  $\widehat{\mathbf{W}}^{1/2}\mathbf{X}$  that is an  $n \times 1$  vector of ones,  $\text{DFBETAS}_{ij}$  reduces from (7) to

$$\text{DFBETAS}_i = \frac{y_i^* - \widehat{\mu}_i^*}{\sqrt{\widehat{v}_i}} \sqrt{\frac{\widehat{\phi}^{(-i)}}{\widehat{\phi}}} \frac{\sqrt{n}}{n - 1}$$

where the quantity  $(y_i^* - \widehat{\mu}_i^*)/\sqrt{\widehat{v}_i}$  is a standardized quantity of  $y_i^*$ , since the expectation and variance of the random variable  $Y_i^* = \log(Y_i/(1 - Y_i))$  are  $E(Y_i^*) = \mu_i^*$  and  $\text{Var}(Y_i^*) = v_i$ , respectively. Thus, the cut-off values of  $\text{DFBETAS}_{ij}$  and  $\text{GDFBETAS}_{ij}$  are suggested as  $2/\sqrt{n}$  and  $2c(s)/\sqrt{s}$ , respectively.

Here it should be pointed out that the proposed diagnostic method has a resemblance to the forward search (FS) algorithm that was introduced by [Hadi and Simonoff \(1993\)](#) and [Atkinson \(1994\)](#) for identifying multiple outliers in linear models and in multivariate data, respectively. The proposed method and the FS algorithm have two similar stages. In the first stage, it is attempted to form a basic subset that is presumably free from potential outliers. In the second stage, it uses an appropriate diagnostic measure such as the adjusted residual or Cook distance to examine the potential outliers to see how extreme they are related to the basic subset. The possible outliers are then declared as outliers if they are greatly inconsistent with the majority of the data. Such a FS diagnostic technique was also applied to GLMs. For details, see [Atkinson and Riani \(2000, Chapter 6\)](#). A recent survey of theoretical development in work on the FS diagnostic was given in [Atkinson et al. \(2010\)](#), who tried to get suitable envelopes

for the outlier tests as the size of the current  $[R]$  set grows. The modified FS diagnostic procedure has not been applied to beta regression models.

## 5 Simulations

In this section, simulation studies are carried out to investigate the finite sample performance of the proposed group deletion diagnostic measures. We consider the two beta models

$$\text{Model 1 : } \text{logit}(\mu_i) = \log(\mu_i/(1 - \mu_i)) = \beta_1 + \beta_2 x_{i,2} + \beta_3 x_{i,3}$$

and

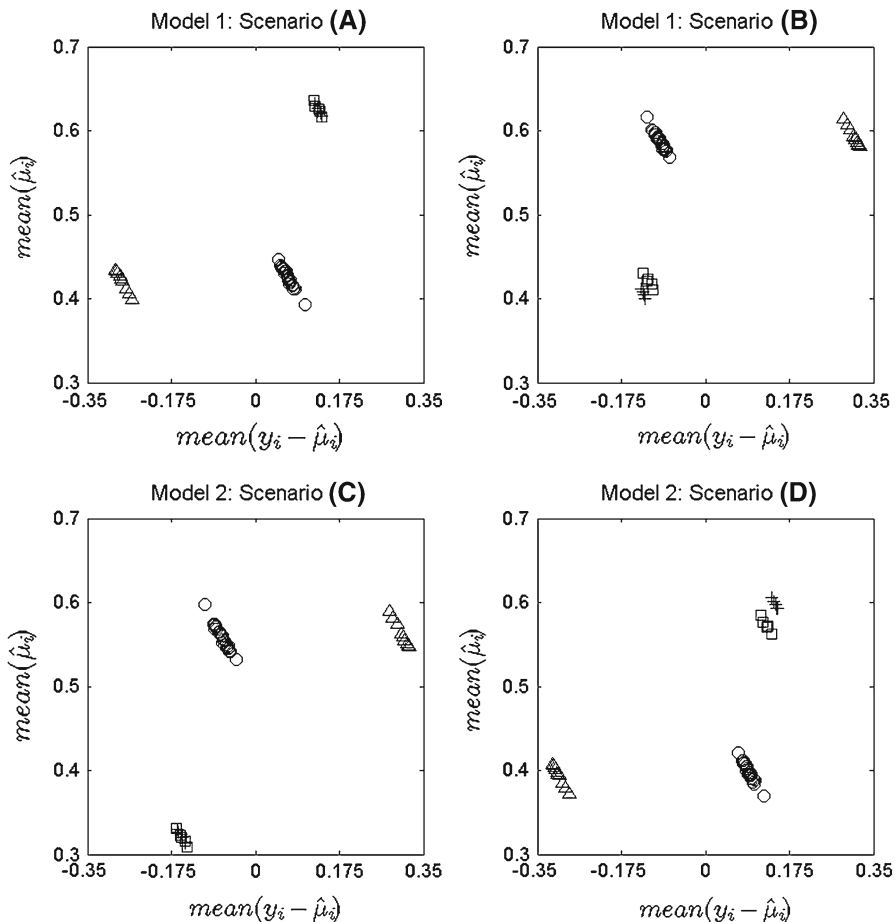
$$\text{Model 2 : } \Phi(\mu_i) = \beta_1 + \beta_2 x_{i,2} + \beta_3 x_{i,3}$$

with the logit and probit link functions, respectively. To show the different performance characters between  $\text{GSR}_i$ ,  $\text{GLD}_i$ ,  $\text{GDFITS}_i$  and  $\text{GDFBETAS}_i$ , numerical studies are performed for sample sizes  $n = 50$  and  $150$  with the first  $0.2n$  observations generated as outlying observations that have unusual residuals, the next  $0.1n$  observations generated as outlying data points that are influential on  $\hat{\beta}_2$ , and subsequently the next  $0.1n$  observations generated as outlying data points that have an influence on  $\hat{\beta}_3$ . To further contrast the differences between the performances based on the single-case and group deletion diagnostics, we consider two different outlier scenarios under each model. These scenarios are chosen because they are situations in which groups of outliers are influential for regression analyses but they are not easily identified by the single-case deletion diagnostics.

Let  $\text{logit}^{-1}(\cdot)$  and  $\Phi^{-1}(\cdot)$  be the inverses of  $\text{logit}(\cdot)$  and  $\Phi(\cdot)$ , respectively. Now consider Model 1 with  $n = 50$  and  $\phi = 30$ . Let  $\beta_1 = \beta_2 = \beta_3 = -0.01$  and  $\eta_i = -0.01 - 0.01x_{i,2} - 0.01x_{i,3}$ . Then the usual observations  $y_{21}, \dots, y_{50}$  are independently generated from the beta distribution,  $\text{Beta}(\mu_i\phi, (1 - \mu_i)\phi)$ , with  $\mu_i = \text{logit}^{-1}(\eta_i)$  in which  $x_{i,2}$  and  $x_{i,3}$  are independently from the uniform distribution  $\text{Uniform}(0.85, 1.15)$ . Under the outlier scenario (A), the outlying observations  $y_1, \dots, y_{10}$  that have extra small residuals are independently generated as the usual data points, except their means set by  $\mu_i = \text{logit}^{-1}(\eta_i) - 0.35$ . The outlying data points  $y_{11}, \dots, y_{15}$  that are influential on  $\hat{\beta}_2$  are independently from the beta distribution with  $\mu_i = \text{logit}^{-1}(-0.01 + 0.85x_{i,2} - 0.01x_{i,3})$  where  $x_{i,2}$  and  $x_{i,3}$  are independently from  $\text{Uniform}(1.3, 1.4)$  and  $\text{Uniform}(0.85, 1.15)$ . The outlying data points  $y_{16}, \dots, y_{20}$  that have an influence on  $\hat{\beta}_3$  are independently generated by the same processes, except their means set by  $\mu_i = \text{logit}^{-1}(-0.01 - 0.01x_{i,2} + 0.85x_{i,3})$  with  $x_{i,2}$  and  $x_{i,3}$  respectively from  $\text{Uniform}(0.85, 1.15)$  and  $\text{Uniform}(1.3, 1.4)$ . On the other hand, under the outlier scenario (B), the outliers  $y_1, \dots, y_{20}$  are generated as under the outlier scenario (A), except their respective corresponding means reset by  $\mu_i = \text{logit}^{-1}(\eta_i) + 0.41$ ,  $\mu_i = \text{logit}^{-1}(-0.01 - 0.6x_{i,2} - 0.01x_{i,3})$  and  $\mu_i = \text{logit}^{-1}(-0.01 - 0.01x_{i,2} - 0.7x_{i,3})$  for observation  $i, i = 1, \dots, 10, i = 11, \dots, 15$  and  $i = 16, \dots, 20$ , respectively. Data for  $n = 150$  are similarly generated by the same steps.

Next consider Model 2 with  $n = 50$  and  $\phi = 30$ . The usual observations  $y_{21}, \dots, y_{50}$  are independently from  $\text{Beta}(\mu_i\phi, (1 - \mu_i)\phi)$  with  $\mu_i = \Phi^{-1}(\eta_i)$

where  $x_{i,2}$  and  $x_{i,3}$  are independently from *Uniform* (0.85, 1.15). Under the outlier scenario (C), the outlying observations  $y_1, \dots, y_{10}$  that have extremely large residuals are independently generated as the usual data points, except their means set by  $\mu_i = \Phi^{-1}(\eta_i) + 0.38$ . The outlying data points  $y_{11}, \dots, y_{15}$  and  $y_{16}, \dots, y_{20}$  that respectively have an effect on  $\hat{\beta}_2$  and  $\hat{\beta}_3$  are independently from the beta distribution with  $\mu_i = \Phi^{-1}(-0.01 - 0.7x_{i,2} - 0.01x_{i,3})$  and  $\mu_i = \Phi^{-1}(-0.01 - 0.01x_{i,2} - 0.7x_{i,3})$ , respectively, where  $x_{i,2}$  and  $x_{i,3}$  are generated as in Model 1. Under the outlier scenario (D), the outliers  $y_1, \dots, y_{20}$  are generated as under the outlier scenario (C), except their respective corresponding means reset by  $\mu_i = \Phi^{-1}(\eta_i) - 0.4$ ,  $\mu_i = \Phi^{-1}(-0.01 + 0.4x_{i,2} - 0.01x_{i,3})$  and  $\mu_i = \Phi^{-1}(-0.01 - 0.01x_{i,2} + 0.5x_{i,3})$ , respectively. Data for  $n = 150$  are also generated in the same way.



**Fig. 1** Scatter plots of the averages of  $2000\hat{\mu}_i$  against  $2000(y_i - \hat{\mu}_i)$  for each outlier scenario for sample size  $n = 50$ . Triangle indexes the outlying observations that have unusual residuals; square indexes the outlying data points that are influential on  $\hat{\beta}_2$ ; plus indexes the outlying data points that have an influence on  $\hat{\beta}_3$ ; circle indexes the usual data points

Two thousand simulation runs are carried out for each case, with the  $x_{i,2}$ 's and  $x_{i,3}$ 's being regenerated after every 50 simulations. Figure 1 shows scatter plots of the averages of  $2000\widehat{\mu}_i$  against  $2000(y_i - \widehat{\mu}_i)$  for each outlier scenario for sample size  $n = 50$ . Because the results of sample sizes  $n = 50$  and  $150$  are similar, we omit the results for  $n = 150$  to reduce space. From Fig. 1, it is evident that the  $mean(y_i - \widehat{\mu}_i)$  values of the outlying observations that have unusual residuals are greater or less than that of the usual observations and the outlying data points that are influential on  $\widehat{\beta}_2$  or  $\widehat{\beta}_3$ . On the contrary, the  $mean(\widehat{\mu}_i)$  values of the outlying data points that have an influence on  $\widehat{\beta}_2$  or  $\widehat{\beta}_3$  are bigger or lower than that of the usual data points and the outlying observations that have extreme residuals. The indication reveals that the former have extremely unusual residuals but they don't have significant impacts on  $\widehat{\mu}_i$ , whereas the latter have influences on  $\widehat{\mu}_i$  but they don't have extraordinarily unusual residuals.

Tables 1, 2, 3 and 4 display the averages of the proportions of observation  $i$ ,  $i = 1, \dots, 0.2n$ , observation  $i$ ,  $i = 0.2n + 1, \dots, 0.3n$ , and observation  $i$ ,  $i = 0.3n + 1, \dots, 0.4n$ , claimed by the single-case and group deletion diagnostics as outlying observations in the two thousand simulated data sets. Also displayed are the averages of the proportions of observation  $i$ ,  $i = 1, \dots, n$ , identified as outliers in the 2000 replications, in which no observations are generated as outlying observations in the simulated data sets. In addition, Tables 1, 2, 3 and 4 also exhibit the results based on the different diagnostic criteria, in order to compare the results from the suggested cut-off values.

From Tables 1, 2, 3 and 4, it is clear that, in the cases where no observations are generated as outliers in the simulated data sets and the diagnostic criteria are based on the suggested cut-off values, the averages of the proportions of observation  $i$ ,  $i = 1, \dots, n$ , detected by the single-case and group deletion diagnostics as the unusual observations in the 2000 replications gradually approach 0.05 as the sample size  $n$  increases. The indications imply that, on the basis of the suggested cut-off values, the probability that the usual observations are misdiagnosed as unusual observations by the single-case and group deletion diagnostics is around 0.05. This, in turn, means that, in the cases where no observations are generated as outliers in the simulated data sets, the results from the single-case and group deletion diagnostics, based on the suggested cut-off values, are reasonable.

On the other hand, from Tables 1, 2, 3 and 4, it is clear that, in the cases with observations generated as outliers in the simulated data sets, the group deletion diagnostics are more effective than the single-case deletion diagnostics, in picking up the outliers from the data sets. As in the case in Model 2 with  $n = 150$  under the outlier scenario (C), when the cut-off values of  $DFBETAS_{i2}$  and  $GDFBETAS_{i2}$  are considered by  $2/\sqrt{n}$  and  $2c(s)/\sqrt{s}$ , respectively, the average of the proportions of observation  $i$ ,  $i = 31, \dots, 45$ , highlighted as outlying data points in the 2000 simulations by  $GDFBETAS_{i2}$  is 0.9438, bigger than the corresponding result 0.3661 from  $DFBETAS_{i2}$ . Similar results are also obtained by comparing differences between the results from  $GDFBETAS_{i3}$  and  $DFBETAS_{i3}$  or from  $GLD_i$  and  $LD_i$ , etc.

From Tables 1, 2, 3 and 4, it is also noted that  $SWR_i$  and  $GSWR_i$  identify outlying observations that have extra unusual residuals as atypical observations, whereas the outlying data points that have influences on  $\widehat{\beta}_2$  or  $\widehat{\beta}_3$  are not flagged by  $SWR_i$  and  $GSWR_i$  as unusual data points. This is because, the latter don't have enough extreme

**Table 1** Model 1:  $\text{logit}(\mu_i) = -0.01 - 0.01x_{i,2} - 0.01x_{i,3}, i = 1, \dots, 50$

Diagnostic	Cut-off value	Data without outlying observations	Data with outlying observations <sup>a</sup>				
			Observations 1–50	Observations 1–10	Observations 11–15	Observations 16–20	Observations 21–50
Scenario (A)							
SWR <sub><i>i</i></sub>	2	0.0518	<b>0.3205</b>	0.0116	0.0109	0.0008	
GSWR <sub><i>i</i></sub>	2 <i>c</i> ( <i>s</i> )	0.0424	<b>0.3434</b>	0.0037	0.0037	0.0002	
LD <sub><i>i</i></sub>	4 <i>p</i> / <i>n</i>	0.0658	<b>0.1027</b>	<b>0.1601</b>	<b>0.1562</b>	0.0033	
GLD <sub><i>i</i></sub>	4 <i>c</i> <sup>2</sup> ( <i>s</i> ) <i>p</i> / <i>s</i>	0.0488	0.0431	<b>0.2786</b>	<b>0.2854</b>	0.0005	
DFFITS <sub><i>i</i></sub>	2√ <i>p</i> / <i>n</i>	0.0714	<b>0.1319</b>	<b>0.1616</b>	<b>0.1587</b>	0.0004	
GDFITS <sub><i>i</i></sub>	2 <i>c</i> ( <i>s</i> )√ <i>p</i> / <i>s</i>	0.0623	<b>0.1033</b>	<b>0.4049</b>	<b>0.3972</b>	0.0011	
DFBETAS <sub><i>i</i>2</sub>	2/√ <i>n</i>	0.0810	0.0898	<b>0.3769</b>	0.0134	0.0093	
GDFBETAS <sub><i>i</i>2</sub>	2 <i>c</i> ( <i>s</i> )/√ <i>s</i>	0.0787	<b>0.1202</b>	<b>0.8435</b>	0.0061	0.0019	
DFBETAS <sub><i>i</i>3</sub>	2/√ <i>n</i>	0.0816	0.0855	0.0205	<b>0.3829</b>	0.0126	
GDFBETAS <sub><i>i</i>3</sub>	2 <i>c</i> ( <i>s</i> )/√ <i>s</i>	0.0800	<b>0.1100</b>	0.0136	<b>0.8257</b>	0.0024	
SWR <sub><i>i</i></sub>	2.5	0.0142	<b>0.1235</b>	0.0015	0.0008	0.0000	
GSWR <sub><i>i</i></sub>	2.5 <i>c</i> ( <i>s</i> )	0.0131	<b>0.1228</b>	0.0012	0.0006	0.0000	
LD <sub><i>i</i></sub>	6.25 <i>p</i> / <i>n</i>	0.0259	0.0203	0.0728	0.0639	0.0002	
GLD <sub><i>i</i></sub>	6.25 <i>c</i> <sup>2</sup> ( <i>s</i> ) <i>p</i> / <i>s</i>	0.0225	0.0142	<b>0.1096</b>	<b>0.1054</b>	0.0001	
DFFITS <sub><i>i</i></sub>	2.5√ <i>p</i> / <i>n</i>	0.0307	0.0337	0.0767	0.0672	0.0006	
GDFITS <sub><i>i</i></sub>	2.5 <i>c</i> ( <i>s</i> )√ <i>p</i> / <i>s</i>	0.0288	0.0309	<b>0.2733</b>	<b>0.2791</b>	0.0001	
DFBETAS <sub><i>i</i>2</sub>	2.5/√ <i>n</i>	0.0425	0.0344	<b>0.2570</b>	0.0044	0.0020	
GDFBETAS <sub><i>i</i>2</sub>	2.5 <i>c</i> ( <i>s</i> )/√ <i>s</i>	0.0423	0.0458	<b>0.7400</b>	0.0018	0.0003	
DFBETAS <sub><i>i</i>3</sub>	2.5/√ <i>n</i>	0.0446	0.0319	0.0065	<b>0.2606</b>	0.0035	
GDFBETAS <sub><i>i</i>3</sub>	2.5 <i>c</i> ( <i>s</i> )/√ <i>s</i>	0.0443	0.0447	0.0039	<b>0.7312</b>	0.0009	
Scenario (B)							
SWR <sub><i>i</i></sub>	2	0.0518	<b>0.3861</b>	0.0020	0.0016	0.0002	
GSWR <sub><i>i</i></sub>	2 <i>c</i> ( <i>s</i> )	0.0424	<b>0.4485</b>	0.0004	0.0004	0.0000	
LD <sub><i>i</i></sub>	4 <i>p</i> / <i>n</i>	0.0658	<b>0.1426</b>	0.0742	0.0819	0.0014	
GLD <sub><i>i</i></sub>	4 <i>c</i> <sup>2</sup> ( <i>s</i> ) <i>p</i> / <i>s</i>	0.0488	<b>0.1004</b>	<b>0.1055</b>	<b>0.1483</b>	0.0005	
DFFITS <sub><i>i</i></sub>	2√ <i>p</i> / <i>n</i>	0.0714	<b>0.1727</b>	0.0755	0.0845	0.0017	
GDFITS <sub><i>i</i></sub>	2 <i>c</i> ( <i>s</i> )√ <i>p</i> / <i>s</i>	0.0623	<b>0.1876</b>	<b>0.1569</b>	<b>0.2074</b>	0.0004	
DFBETAS <sub><i>i</i>2</sub>	2/√ <i>n</i>	0.0810	<b>0.1213</b>	<b>0.2585</b>	0.0045	0.0056	
GDFBETAS <sub><i>i</i>2</sub>	2 <i>c</i> ( <i>s</i> )/√ <i>s</i>	0.0787	<b>0.1602</b>	<b>0.6769</b>	0.0019	0.0013	
DFBETAS <sub><i>i</i>3</sub>	2/√ <i>n</i>	0.0816	<b>0.1131</b>	0.0069	<b>0.2886</b>	0.0081	
GDFBETAS <sub><i>i</i>3</sub>	2 <i>c</i> ( <i>s</i> )/√ <i>s</i>	0.0800	<b>0.1470</b>	0.0045	<b>0.7397</b>	0.0014	
SWR <sub><i>i</i></sub>	2.5	0.0142	<b>0.1609</b>	0.0001	0.0000	0.0000	
GSWR <sub><i>i</i></sub>	2.5 <i>c</i> ( <i>s</i> )	0.0131	<b>0.1726</b>	0.0001	0.0000	0.0000	
LD <sub><i>i</i></sub>	6.25 <i>p</i> / <i>n</i>	0.0259	0.0351	0.0254	0.0243	0.0002	
GLD <sub><i>i</i></sub>	6.25 <i>c</i> <sup>2</sup> ( <i>s</i> ) <i>p</i> / <i>s</i>	0.0225	0.0317	0.0322	0.0350	0.0002	

**Table 1** continued

Diagnostic	Cut-off value	Data without outlying observations	Data with outlying observations <sup>a</sup>				
			Observations 1–50	Observations 1–10	Observations 11–15	Observations 16–20	Observations 21–50
DFFITS <sub><i>i</i></sub>	$2.5\sqrt{p/n}$	0.0307	0.0519	0.0266	0.0260	0.0003	
GDFFITs <sub><i>i</i></sub>	$2.5c(s)\sqrt{p/s}$	0.0288	0.0531	0.0744	<b><u>0.1054</u></b>	0.0002	
DFBETAS <sub><i>i</i>2</sub>	$2.5/\sqrt{n}$	0.0425	0.0537	<b>0.1464</b>	0.0008	0.0009	
GDFBETAS <sub><i>i</i>2</sub>	$2.5c(s)/\sqrt{s}$	0.0423	0.0683	<b><u>0.5066</u></b>	0.0003	0.0002	
DFBETAS <sub><i>i</i>3</sub>	$2.5/\sqrt{n}$	0.0446	0.0461	0.0023	<b>0.1681</b>	0.0019	
GDFBETAS <sub><i>i</i>3</sub>	$2.5c(s)/\sqrt{s}$	0.0443	0.0625	0.0011	<b><u>0.5812</u></b>	0.0006	

The diagnostic results based on single-case and group deletion diagnostics are larger than 0.1, which are respectively highlighted as the significance of bold and underline values

<sup>a</sup> observations 1–10 have unusual residuals; observations 11–15 and 16–20 respectively have influences on  $\hat{\beta}_2$  and  $\hat{\beta}_3$ ; observations 21–50 are usual data points

residuals of  $y_i - \hat{\mu}_i$ , as the results shown in Fig. 1, in which the  $mean(y_i - \hat{\mu}_i)$  values of the latter are closer to that of the usual observations, in comparison with the former. On the other hand, it is seen that due to the latter generated as having impacts on regression parameter estimates, they are more easily identified by DFBETAS<sub>*i*</sub> and GDFBETAS<sub>*i*</sub> as unusual data points.

In addition, it is observed from Tables 1, 2, 3 and 4 that, in general, the averages of the proportions of the outlying data points that have influences on  $\hat{\beta}_2$  or  $\hat{\beta}_3$  detected as unusual data points are higher than that of the outlying observations that have unusual residuals. In other words, it is observed that the latter are identified harder than the former. This is because, the latter are not outlying enough. As compared the results from Model 2 under the outlier scenarios (C) and (D), it is shown that the latter become easier to be identified when their residuals become more excessive.

Obviously, the group deletion diagnostics are able to provide the more valid inferences in the regression diagnostic analyses, under the data sets with multiple outliers.

## 6 Examples

Two practical applications are presented in this section to illustrate the usefulness of the proposed group deletion diagnostic measures.

### 6.1 Example 1: Reading accuracy data

The first application uses the data analyzed by Espinheira et al. (2008a) from 44 children (19 children with dyslexia and 25 controls) recruited from primary schools in the Australian Capital Territory. The ages of the children ranged from eight years five months to twelve years three months. The response ( $y$ ) gives the scores on a test of reading accuracy, and the explanatory variables represent dyslexia versus non-dyslexia

**Table 2** Model 1:  $\text{logit}(\mu_i) = -0.01 - 0.01x_{i,2} - 0.01x_{i,3}, i = 1, \dots, 150$

Diagnostic	Cut-off value	Data without outlying observations	Data with outlying observations <sup>a</sup>				
			Observations 1–150	Observations 1–30	Observations 31–45	Observations 46–60	Observations 61–150
Scenario (A)							
SWR <sub><i>i</i></sub>	2	0.0481	<b>0.3096</b>	0.0088	0.0087	0.0006	
GSWR <sub><i>i</i></sub>	2 <i>c</i> ( <i>s</i> )	0.0382	<b>0.3290</b>	0.0021	0.0018	0.0001	
LD <sub><i>i</i></sub>	4 <i>p</i> / <i>n</i>	0.0547	<b>0.1040</b>	<b>0.1173</b>	<b>0.1198</b>	0.0013	
GLD <sub><i>i</i></sub>	4 <i>c</i> <sup>2</sup> ( <i>s</i> ) <i>p</i> / <i>s</i>	0.0377	0.0512	<b>0.1680</b>	<b>0.1618</b>	0.0002	
DFFITS <sub><i>i</i></sub>	2√ <i>p</i> / <i>n</i>	0.0566	<b>0.1120</b>	<b>0.1183</b>	<b>0.1206</b>	0.0015	
GDFITS <sub><i>i</i></sub>	2 <i>c</i> ( <i>s</i> )√ <i>p</i> / <i>s</i>	0.0470	0.0765	<b>0.3752</b>	<b>0.3575</b>	0.0003	
DFBETAS <sub><i>i</i>2</sub>	2/√ <i>n</i>	0.0697	0.0858	<b>0.3299</b>	0.0096	0.0058	
GDFBETAS <sub><i>i</i>2</sub>	2 <i>c</i> ( <i>s</i> )/√ <i>s</i>	0.0628	0.0978	<b>0.8951</b>	0.0028	0.0003	
DFBETAS <sub><i>i</i>3</sub>	2/√ <i>n</i>	0.0688	0.0805	0.0103	<b>0.3301</b>	0.0062	
GDFBETAS <sub><i>i</i>3</sub>	2 <i>c</i> ( <i>s</i> )/√ <i>s</i>	0.0622	0.0902	0.0021	<b>0.8973</b>	0.0003	
SWR <sub><i>i</i></sub>	2.5	0.0136	<b>0.1226</b>	0.0009	0.0006	0.0000	
GSWR <sub><i>i</i></sub>	2.5 <i>c</i> ( <i>s</i> )	0.0121	<b>0.1235</b>	0.0005	0.0004	0.0000	
LD <sub><i>i</i></sub>	6.25 <i>p</i> / <i>n</i>	0.0199	0.0228	0.0460	0.0442	0.0001	
GLD <sub><i>i</i></sub>	6.25 <i>c</i> <sup>2</sup> ( <i>s</i> ) <i>p</i> / <i>s</i>	0.0161	0.0169	0.0500	0.0471	0.0000	
DFFITS <sub><i>i</i></sub>	2.5√ <i>p</i> / <i>n</i>	0.0213	0.0269	0.0470	0.0455	0.0001	
GDFITS <sub><i>i</i></sub>	2.5 <i>c</i> ( <i>s</i> )√ <i>p</i> / <i>s</i>	0.0191	0.0215	<b>0.1666</b>	<b>0.1497</b>	0.0001	
DFBETAS <sub><i>i</i>2</sub>	2.5/√ <i>n</i>	0.0345	0.0307	<b>0.2064</b>	0.0023	0.0008	
GDFBETAS <sub><i>i</i>2</sub>	2.5 <i>c</i> ( <i>s</i> )/√ <i>s</i>	0.0320	0.0342	<b>0.7709</b>	0.0007	0.0001	
DFBETAS <sub><i>i</i>3</sub>	2.5/√ <i>n</i>	0.0343	0.0290	0.0025	<b>0.2067</b>	0.0008	
GDFBETAS <sub><i>i</i>3</sub>	2.5 <i>c</i> ( <i>s</i> )/√ <i>s</i>	0.0319	0.0317	0.0005	<b>0.7838</b>	0.0001	
Scenario (B)							
SWR <sub><i>i</i></sub>	2	0.0481	<b>0.3715</b>	0.0009	0.0012	0.0001	
GSWR <sub><i>i</i></sub>	2 <i>c</i> ( <i>s</i> )	0.0382	<b>0.4288</b>	0.0000	0.0001	0.0000	
LD <sub><i>i</i></sub>	4 <i>p</i> / <i>n</i>	0.0547	<b>0.1429</b>	0.0421	0.0503	0.0004	
GLD <sub><i>i</i></sub>	4 <i>c</i> <sup>2</sup> ( <i>s</i> ) <i>p</i> / <i>s</i>	0.0377	<b>0.1127</b>	0.0315	0.0039	0.0001	
DFFITS <sub><i>i</i></sub>	2√ <i>p</i> / <i>n</i>	0.0566	<b>0.1519</b>	0.0426	0.0507	0.0005	
GDFITS <sub><i>i</i></sub>	2 <i>c</i> ( <i>s</i> )√ <i>p</i> / <i>s</i>	0.0470	<b>0.1605</b>	0.0591	<b>0.1038</b>	0.0001	
DFBETAS <sub><i>i</i>2</sub>	2/√ <i>n</i>	0.0697	<b>0.1176</b>	<b>0.1987</b>	0.0031	0.0028	
GDFBETAS <sub><i>i</i>2</sub>	2 <i>c</i> ( <i>s</i> )/√ <i>s</i>	0.0628	<b>0.1383</b>	<b>0.6972</b>	0.0008	0.0002	
DFBETAS <sub><i>i</i>3</sub>	2/√ <i>n</i>	0.0688	<b>0.1060</b>	0.0024	<b>0.2222</b>	0.0033	
GDFBETAS <sub><i>i</i>3</sub>	2 <i>c</i> ( <i>s</i> )/√ <i>s</i>	0.0622	<b>0.1241</b>	0.0006	<b>0.7654</b>	0.0004	
SWR <sub><i>i</i></sub>	2.5	0.0136	<b>0.1588</b>	0.0000	0.0001	0.0000	
GSWR <sub><i>i</i></sub>	2.5 <i>c</i> ( <i>s</i> )	0.0121	<b>0.1724</b>	0.0000	0.0000	0.0000	
LD <sub><i>i</i></sub>	6.25 <i>p</i> / <i>n</i>	0.0199	0.0369	0.0100	0.0128	0.0000	
GLD <sub><i>i</i></sub>	6.25 <i>c</i> <sup>2</sup> ( <i>s</i> ) <i>p</i> / <i>s</i>	0.0161	0.0338	0.0092	0.0114	0.0001	



**Table 2** continued

Diagnostic	Cut-off value	Data without outlying observations	Data with outlying observations <sup>a</sup>			
		Observations 1–150	Observations 1–30	Observations 31–45	Observations 46–60	Observations 61–150
DFFITS <sub><i>i</i></sub>	$2.5\sqrt{p/n}$	0.0213	0.0424	0.0104	0.0131	0.0000
GDFFITs <sub><i>i</i></sub>	$2.5c(s)\sqrt{p/s}$	0.0191	0.0428	0.0133	0.0247	0.0000
DFBETAS <sub><i>i</i>2</sub>	$2.5/\sqrt{n}$	0.0345	0.0487	0.0955	0.0005	0.0002
GDFBETAS <sub><i>i</i>2</sub>	$2.5c(s)/\sqrt{s}$	0.0320	0.0541	<b>0.4334</b>	0.0004	0.0000
DFBETAS <sub><i>i</i>3</sub>	$2.5/\sqrt{n}$	0.0343	0.0430	0.0003	<b>0.1122</b>	0.0003
GDFBETAS <sub><i>i</i>3</sub>	$2.5c(s)/\sqrt{s}$	0.0319	0.0471	0.0001	<b>0.5277</b>	0.0001

The diagnostic results based on single-case and group deletion diagnostics are larger than 0.1, which are respectively highlighted as the significance of bold and underline values

<sup>a</sup> observations 1–30 have unusual residuals; observations 31–45 and 46–60 respectively have influences on  $\hat{\beta}_2$  and  $\hat{\beta}_3$ ; observations 61–150 are usual data points

status ( $x_2$ ), non-verbal IQ scores converted to z-scores ( $x_3$ ) and an interaction variable ( $x_4$ ). The variable  $x_2$  is coded as 1 when the child is dyslexic and otherwise it is coded as  $-1$ . The non-dyslexic readers’ mean accuracy score is 0.900 whereas the mean for readers who have dyslexia is 0.606. The overall mean score is 0.773. Following the suggestions in [Espinheira et al. \(2008a\)](#), we analyze the data set using the GBLM with the logit link function as

$$\text{logit}(\mu_i) = \beta_1 + \beta_2x_{i,2} + \beta_3x_{i,3} + \beta_4x_{i,4}, \quad i = 1, \dots, 44.$$

Table 5 displays diagnostic results based on the single-case and group deletion diagnostics, respectively. Looking at Table 5, some interesting findings are as follows. First, from  $\text{GSWR}_i$  and  $\text{SWR}_i$ , it is observed that  $\text{GSWR}_i$  suggests observations 8, 9, 15 and 22 as the outlying observations that have unusual residuals, whereas, among them, only observation 8 is highlighted by  $\text{SWR}_i$  as an unusual observation on the residual. Then comparing the results from  $\text{LD}_i$  and  $\text{GLD}_i$ , it is noted that  $\text{LD}_i$  detects observations 6 and 8 as the outlying data points that have an impact on the regression parameter estimates,  $\hat{\beta}_1, \hat{\beta}_2, \hat{\beta}_3$  and  $\hat{\beta}_4$ , while  $\text{GLD}_i$  identifies observations 8 and 15 as the atypical observations that are influential on  $\hat{\beta}_1, \hat{\beta}_2, \hat{\beta}_3$  and  $\hat{\beta}_4$ . The joint deletion of observations 6 and 8 shows that the relative changes in  $\hat{\beta}_1, \hat{\beta}_2, \hat{\beta}_3$  and  $\hat{\beta}_4$  are  $-1.773, -2.761, 32.77$  and  $24.36\%$ , whereas the relative changes in  $\hat{\beta}_1, \hat{\beta}_2, \hat{\beta}_3$  and  $\hat{\beta}_4$  are  $-11.32, -16.03, 117.5$  and  $86.82\%$  with observations 8 and 15 deleted. The indication implies that observations 8 and 15 jointly have more powerful influences on the regression parameter estimates, in contrast with observations 6 and 8.

Similarly, comparing the results from  $\text{DFFITS}_i$  and  $\text{GDFFITs}_i$ , it is also observed that  $\text{GDFFITs}_i$  detects observations 5, 8, 9, 15 and 22 as the atypical points on the weighted fits. Among these atypical points, observations 5, 9, 15 and 22 are obscured by  $\text{DFFITS}_i$ . Jointly removing observations 5, 8, 9, 15 and 22 from the data results in the apparently relative variations,  $-28.10, -39.86, 269.1$  and  $199.0\%$  in  $\hat{\beta}_1, \hat{\beta}_2, \hat{\beta}_3$  and  $\hat{\beta}_4$ .

**Table 3** Model 2:  $\Phi(\mu_i) = -0.01 - 0.01x_{i,2} - 0.01x_{i,3}, i = 1, \dots, 50$

Diagnostic	Cut-off value	Data without outlying observations	Data with outlying observations <sup>a</sup>				
			Observations 1–50	Observations 1–10	Observations 11–15	Observations 16–20	Observations 21–50
Scenario (C)							
SWR <sub><i>i</i></sub>	2	0.0517	<b>0.3082</b>	0.0193	0.0182	0.0003	
GSWR <sub><i>i</i></sub>	2 <i>c</i> ( <i>s</i> )	0.0423	<b>0.3242</b>	0.0087	0.0082	0.0001	
LD <sub><i>i</i></sub>	4 <i>p</i> / <i>n</i>	0.0657	0.0828	<b>0.1917</b>	<b>0.1907</b>	0.0030	
GLD <sub><i>i</i></sub>	4 <i>c</i> <sup>2</sup> ( <i>s</i> ) <i>p</i> / <i>s</i>	0.0489	0.0231	<b>0.3460</b>	<b>0.3614</b>	0.0003	
DFFITS <sub><i>i</i></sub>	2√ <i>p</i> / <i>n</i>	0.0714	<b>0.1488</b>	<b>0.1099</b>	<b>0.1040</b>	0.0037	
GDFFITs <sub><i>i</i></sub>	2 <i>c</i> ( <i>s</i> )√ <i>p</i> / <i>s</i>	0.0624	<b>0.1440</b>	<b>0.2729</b>	<b>0.2739</b>	0.0010	
DFBETAS <sub><i>i</i>2</sub>	2/√ <i>n</i>	0.0425	0.0761	<b>0.4132</b>	0.0189	0.0091	
GDFBETAS <sub><i>i</i>2</sub>	2 <i>c</i> ( <i>s</i> )/√ <i>s</i>	0.0424	<b>0.1017</b>	<b>0.8840</b>	0.0079	0.0015	
DFBETAS <sub><i>i</i>3</sub>	2/√ <i>n</i>	0.0815	0.0723	0.0295	<b>0.4245</b>	0.0122	
GDFBETAS <sub><i>i</i>3</sub>	2 <i>c</i> ( <i>s</i> )/√ <i>s</i>	0.0799	0.0981	0.0173	<b>0.8841</b>	0.0016	
SWR <sub><i>i</i></sub>	2.5	0.0142	<b>0.1084</b>	0.0028	0.0017	0.0000	
GSWR <sub><i>i</i></sub>	2.5 <i>c</i> ( <i>s</i> )	0.0132	<b>0.1051</b>	0.0025	0.0014	0.0000	
LD <sub><i>i</i></sub>	6.25 <i>p</i> / <i>n</i>	0.0258	0.0136	0.0908	0.0856	0.0005	
GLD <sub><i>i</i></sub>	6.25 <i>c</i> <sup>2</sup> ( <i>s</i> ) <i>p</i> / <i>s</i>	0.0225	0.0080	<b>0.1467</b>	<b>0.1434</b>	0.0002	
DFFITS <sub><i>i</i></sub>	2.5√ <i>p</i> / <i>n</i>	0.0307	0.0403	0.0451	0.0384	0.0006	
GDFFITs <sub><i>i</i></sub>	2.5 <i>c</i> ( <i>s</i> )√ <i>p</i> / <i>s</i>	0.0287	0.0393	<b>0.1636</b>	<b>0.1530</b>	0.0003	
DFBETAS <sub><i>i</i>2</sub>	2.5/√ <i>n</i>	0.0811	0.0260	<b>0.2940</b>	0.0064	0.0019	
GDFBETAS <sub><i>i</i>2</sub>	2.5 <i>c</i> ( <i>s</i> )/√ <i>s</i>	0.0787	0.0358	<b>0.8121</b>	0.0013	0.0002	
DFBETAS <sub><i>i</i>3</sub>	2.5/√ <i>n</i>	0.0446	0.0248	0.0109	<b>0.3007</b>	0.0035	
GDFBETAS <sub><i>i</i>3</sub>	2.5 <i>c</i> ( <i>s</i> )/√ <i>s</i>	0.0442	0.0369	0.0049	<b>0.8163</b>	0.0005	
Scenario (D)							
SWR <sub><i>i</i></sub>	2	0.0517	<b>0.3891</b>	0.0024	0.0025	0.0001	
GSWR <sub><i>i</i></sub>	2 <i>c</i> ( <i>s</i> )	0.0423	<b>0.4513</b>	0.0002	0.0002	0.0000	
LD <sub><i>i</i></sub>	4 <i>p</i> / <i>n</i>	0.0657	<b>0.1448</b>	0.0843	<b>0.1009</b>	0.0018	
GLD <sub><i>i</i></sub>	4 <i>c</i> <sup>2</sup> ( <i>s</i> ) <i>p</i> / <i>s</i>	0.0489	0.0923	<b>0.1229</b>	<b>0.2048</b>	0.0004	
DFFITS <sub><i>i</i></sub>	2√ <i>p</i> / <i>n</i>	0.0714	<b>0.1933</b>	0.0553	0.0601	0.0020	
GDFFITs <sub><i>i</i></sub>	2 <i>c</i> ( <i>s</i> )√ <i>p</i> / <i>s</i>	0.0624	<b>0.2295</b>	0.0746	<b>0.1109</b>	0.0005	
DFBETAS <sub><i>i</i>2</sub>	2/√ <i>n</i>	0.0425	<b>0.1192</b>	<b>0.2803</b>	0.0071	0.0056	
GDFBETAS <sub><i>i</i>2</sub>	2 <i>c</i> ( <i>s</i> )/√ <i>s</i>	0.0424	<b>0.1537</b>	<b>0.7037</b>	0.0036	0.0014	
DFBETAS <sub><i>i</i>3</sub>	2/√ <i>n</i>	0.0815	<b>0.1069</b>	0.0083	<b>0.3224</b>	0.0087	
GDFBETAS <sub><i>i</i>3</sub>	2 <i>c</i> ( <i>s</i> )/√ <i>s</i>	0.0799	<b>0.1431</b>	0.0045	<b>0.7903</b>	0.0013	
SWR <sub><i>i</i></sub>	2.5	0.0142	<b>0.1704</b>	0.0001	0.0001	0.0000	
GSWR <sub><i>i</i></sub>	2.5 <i>c</i> ( <i>s</i> )	0.0132	<b>0.1820</b>	0.0001	0.0000	0.0000	
LD <sub><i>i</i></sub>	6.25 <i>p</i> / <i>n</i>	0.0258	0.0368	0.0290	0.0314	0.0001	
GLD <sub><i>i</i></sub>	6.25 <i>c</i> <sup>2</sup> ( <i>s</i> ) <i>p</i> / <i>s</i>	0.0225	0.0317	0.0386	0.0535	0.0001	

**Table 3** continued

Diagnostic	Cut-off value	Data without outlying observations	Data with outlying observations <sup>a</sup>			
		Observations 1–50	Observations 1–10	Observations 11–15	Observations 16–20	Observations 21–50
DFFITS <sub><i>i</i></sub>	$2.5\sqrt{p/n}$	0.0307	0.0619	0.0168	0.0155	0.0001
GDFFITs <sub><i>i</i></sub>	$2.5c(s)\sqrt{p/s}$	0.0287	0.0661	0.0323	0.0457	0.0001
DFBETAS <sub><i>i</i>2</sub>	$2.5/\sqrt{n}$	0.0811	0.0533	<b>0.1638</b>	0.0015	0.0008
GDFBETAS <sub><i>i</i>2</sub>	$2.5c(s)/\sqrt{s}$	0.0787	0.0663	<b>0.5440</b>	0.0005	0.0001
DFBETAS <sub><i>i</i>3</sub>	$2.5/\sqrt{n}$	0.0446	0.0425	0.0027	<b>0.1996</b>	0.0023
GDFBETAS <sub><i>i</i>3</sub>	$2.5c(s)/\sqrt{s}$	0.0442	0.0578	0.0010	<b>0.6485</b>	0.0005

The diagnostic results based on single-case and group deletion diagnostics are larger than 0.1, which are respectively highlighted as the significance of bold and underline values

<sup>a</sup> observations 1–10 have unusual residuals; observations 11–15 and 16–20 respectively have influences on  $\hat{\beta}_2$  and  $\hat{\beta}_3$ ; observations 21–50 are usual data points

On the other hand, DFBETAS<sub>*i*3</sub> and DFBETAS<sub>*i*4</sub> detect observations 6, 8 and 15 as having a large impact on  $\hat{\beta}_3$  and  $\hat{\beta}_4$ , while GDFBETAS<sub>*i*3</sub> and GDFBETAS<sub>*i*4</sub> spot observations 5, 8, 9, 15 and 22 as influential on the results for  $\hat{\beta}_3$  and  $\hat{\beta}_4$ . The exclusion of observations 5, 8, 9, 15 and 22 causes an apparent change 269.1 and 199.0% in  $\hat{\beta}_3$  and  $\hat{\beta}_4$ , whereas the relative change in  $\hat{\beta}_3$  and  $\hat{\beta}_4$  are 86.31 and 63.88% with observations 6, 8 and 15 omitted together. This indicates that observations 5, 8, 9, 15 and 22 have a relatively larger effect on  $\hat{\beta}_3$  and  $\hat{\beta}_4$ , compared with observations 6, 8 and 15. In addition, it should be explained that, in this example, we don't consider the results of GDFBETAS<sub>*i*2</sub> and DFBETAS<sub>*i*2</sub>. This is because,  $x_2$  is a binary variable whose value equals  $-1$  or  $1$ . Thus we are not very interested that which observation is influential on  $\hat{\beta}_2$ .

Evidently, in reading accuracy data, the group deletion diagnostics detect some multiple outlying observations that are influential on the results of the regression analysis, while these unusual observations are ignored by the single-case deletion diagnostics.

### 6.2 Example 2: Stress, depression, and anxiety

The second example uses data collected from a sample of 166 nonclinical women in Townsville, Queensland, Australia. The data were analyzed by [Smithson and Verkuilen \(2006\)](#) using beta regression. The response variable ( $y$ ) represents the scores on a test of anxiety symptoms, and the explanatory variable ( $x_2$ ) gives the corresponding scores on the test of stress symptoms. We follow [Smithson and Verkuilen \(2006\)](#) in considering the model

$$\text{logit}(\mu_i) = \beta_1 + \beta_2 x_{i,2}, \quad i = 1, \dots, 166.$$

Table 6 presents diagnostic results based on the single-case and group deletion diagnostics, respectively. Careful inspection of Table 6 indicates some interesting findings as follows.

**Table 4** Model 2:  $\Phi(\mu_i) = -0.01 - 0.01x_{i,2} - 0.01x_{i,3}, i = 1, \dots, 150$

Diagnostic	Cut-off value	Data without outlying observations	Data with outlying observations <sup>a</sup>				
			Observations 1–150	Observations 1–30	Observations 31–45	Observations 46–60	Observations 61–150
<b>Scenario (C)</b>							
SWR <sub><i>i</i></sub>	2	0.0481	<b>0.2960</b>	0.0147	0.0146	0.0002	
GSWR <sub><i>i</i></sub>	2 <i>c</i> ( <i>s</i> )	0.0381	<b>0.3111</b>	0.0038	0.0049	0.0000	
LD <sub><i>i</i></sub>	4 <i>p</i> / <i>n</i>	0.0547	0.0843	<b>0.1448</b>	<b>0.1459</b>	0.0011	
GLD <sub><i>i</i></sub>	4 <i>c</i> <sup>2</sup> ( <i>s</i> ) <i>p</i> / <i>s</i>	0.0377	0.0277	<b>0.2434</b>	<b>0.2172</b>	0.0000	
DFFITS <sub><i>i</i></sub>	2√ <i>p</i> / <i>n</i>	0.0566	<b>0.1267</b>	0.0716	0.0723	0.0012	
GDFFITs <sub><i>i</i></sub>	2 <i>c</i> ( <i>s</i> )√ <i>p</i> / <i>s</i>	0.0470	<b>0.1153</b>	<b>0.1806</b>	<b>0.1648</b>	0.0004	
DFBETAS <sub><i>i</i>2</sub>	2/√ <i>n</i>	0.0700	0.0731	<b>0.3661</b>	0.0160	0.0050	
GDFBETAS <sub><i>i</i>2</sub>	2 <i>c</i> ( <i>s</i> )/√ <i>s</i>	0.0629	0.0807	<b>0.9438</b>	0.0023	0.0001	
DFBETAS <sub><i>i</i>3</sub>	2/√ <i>n</i>	0.0689	0.0677	0.0159	<b>0.3636</b>	0.0054	
GDFBETAS <sub><i>i</i>3</sub>	2 <i>c</i> ( <i>s</i> )/√ <i>s</i>	0.0623	0.0744	0.0027	<b>0.9466</b>	0.0001	
SWR <sub><i>i</i></sub>	2.5	0.0135	<b>0.1088</b>	0.0019	0.0022	0.0000	
GSWR <sub><i>i</i></sub>	2.5 <i>c</i> ( <i>s</i> )	0.0120	<b>0.1085</b>	0.0011	0.0010	0.0000	
LD <sub><i>i</i></sub>	6.25 <i>p</i> / <i>n</i>	0.0200	0.0154	0.0606	0.0608	0.0000	
GLD <sub><i>i</i></sub>	6.25 <i>c</i> <sup>2</sup> ( <i>s</i> ) <i>p</i> / <i>s</i>	0.0161	0.0098	0.0717	0.0675	0.0000	
DFFITS <sub><i>i</i></sub>	2.5√ <i>p</i> / <i>n</i>	0.0213	0.0310	0.0222	0.0212	0.0001	
GDFFITs <sub><i>i</i></sub>	2.5 <i>c</i> ( <i>s</i> )√ <i>p</i> / <i>s</i>	0.0191	0.0297	0.0458	0.0473	0.0000	
DFBETAS <sub><i>i</i>2</sub>	2.5/√ <i>n</i>	0.0345	0.0230	<b>0.2409</b>	0.0041	0.0005	
GDFBETAS <sub><i>i</i>2</sub>	2.5 <i>c</i> ( <i>s</i> )/√ <i>s</i>	0.0321	0.0252	<b>0.8718</b>	0.0008	0.0001	
DFBETAS <sub><i>i</i>3</sub>	2.5/√ <i>n</i>	0.0343	0.0220	0.0043	<b>0.2384</b>	0.0005	
GDFBETAS <sub><i>i</i>3</sub>	2.5 <i>c</i> ( <i>s</i> )/√ <i>s</i>	0.0319	0.0231	0.0005	<b>0.8763</b>	0.0000	
<b>Scenario (D)</b>							
SWR <sub><i>i</i></sub>	2	0.0481	<b>0.3723</b>	0.0010	0.0022	0.0001	
GSWR <sub><i>i</i></sub>	2 <i>c</i> ( <i>s</i> )	0.0381	<b>0.4271</b>	0.0000	0.0002	0.0000	
LD <sub><i>i</i></sub>	4 <i>p</i> / <i>n</i>	0.0547	<b>0.1450</b>	0.0487	0.0657	0.0006	
GLD <sub><i>i</i></sub>	4 <i>c</i> <sup>2</sup> ( <i>s</i> ) <i>p</i> / <i>s</i>	0.0377	<b>0.1092</b>	0.0355	0.0650	0.0001	
DFFITS <sub><i>i</i></sub>	2√ <i>p</i> / <i>n</i>	0.0566	<b>0.1698</b>	0.0259	0.0316	0.0007	
GDFFITs <sub><i>i</i></sub>	2 <i>c</i> ( <i>s</i> )√ <i>p</i> / <i>s</i>	0.0470	<b>0.1940</b>	0.0208	0.0353	0.0002	
DFBETAS <sub><i>i</i>2</sub>	2/√ <i>n</i>	0.0700	<b>0.1156</b>	<b>0.2196</b>	0.0044	0.0029	
GDFBETAS <sub><i>i</i>2</sub>	2 <i>c</i> ( <i>s</i> )/√ <i>s</i>	0.0629	<b>0.1370</b>	<b>0.7267</b>	0.0001	0.0002	
DFBETAS <sub><i>i</i>3</sub>	2/√ <i>n</i>	0.0689	<b>0.1013</b>	0.0030	<b>0.2653</b>	0.0036	
GDFBETAS <sub><i>i</i>3</sub>	2 <i>c</i> ( <i>s</i> )/√ <i>s</i>	0.0623	<b>0.1196</b>	0.0006	<b>0.8432</b>	0.0002	
SWR <sub><i>i</i></sub>	2.5	0.0135	<b>0.1678</b>	0.0000	0.0001	0.0000	
GSWR <sub><i>i</i></sub>	2.5 <i>c</i> ( <i>s</i> )	0.0120	<b>0.1826</b>	0.0000	0.0001	0.0000	
LD <sub><i>i</i></sub>	6.25 <i>p</i> / <i>n</i>	0.0200	0.0398	0.0118	0.0167	0.0000	
GLD <sub><i>i</i></sub>	6.25 <i>c</i> <sup>2</sup> ( <i>s</i> ) <i>p</i> / <i>s</i>	0.0161	0.0355	0.0110	0.0154	0.0000	

**Table 4** continued

Diagnostic	Cut-off value	Data without outlying observations	Data with outlying observations <sup>a</sup>			
		Observations 1–150	Observations 1–30	Observations 31–45	Observations 46–60	Observations 61–150
DFFITS <sub><i>i</i></sub>	$2.5\sqrt{p/n}$	0.0213	0.0514	0.0050	0.0006	0.0000
GDFFITs <sub><i>i</i></sub>	$2.5c(s)\sqrt{p/s}$	0.0191	0.0539	0.0044	0.0078	0.0000
DFBETAS <sub><i>i</i>2</sub>	$2.5/\sqrt{n}$	0.0345	0.0479	<b>0.1091</b>	0.0008	0.0003
GDFBETAS <sub><i>i</i>2</sub>	$2.5c(s)/\sqrt{s}$	0.0321	0.0530	<b>0.4839</b>	0.0004	0.0000
DFBETAS <sub><i>i</i>3</sub>	$2.5/\sqrt{n}$	0.0343	0.0408	0.0004	<b>0.1401</b>	0.0003
GDFBETAS <sub><i>i</i>3</sub>	$2.5c(s)/\sqrt{s}$	0.0319	0.0447	0.0002	<b>0.6547</b>	0.0000

The diagnostic results based on single-case and group deletion diagnostics are larger than 0.1, which are respectively highlighted as the significance of bold and underline values

<sup>a</sup> Observations 1–30 have unusual residuals; observations 31–45 and 46–60 respectively have influences on  $\hat{\beta}_2$  and  $\hat{\beta}_3$ ; observations 61–150 are usual data points

**Table 5** The single-case and group deletion diagnostics for example 1

Diagnostic	Cut-off value	Case <i>i</i> claimed as an outlier	Cut-off value	Case <i>i</i> claimed as an outlier	Cut-off value	Case <i>i</i> claimed as an outlier
SWR <sub><i>i</i></sub>	2	8	2.5	No observations <sup>a</sup>	3	No observations <sup>a</sup>
GSWR <sub><i>i</i></sub>	$2c(s)$	8, 9, 15, 22	$2.5c(s)$	8	$3c(s)$	No observations <sup>a</sup>
LD <sub><i>i</i></sub>	$4p/n$	6, 8	$6.25p/n$	8	$9p/n$	8
GLD <sub><i>i</i></sub>	$4c^2(s)p/s$	8, 15	$6.25c^2(s)p/s$	8	$9c^2(s)p/s$	8
DFFITS <sub><i>i</i></sub>	$2\sqrt{p/n}$	6, 8	$2.5\sqrt{p/n}$	8	$3\sqrt{p/n}$	8
GDFFITs <sub><i>i</i></sub>	$2c(s)\sqrt{p/s}$	5, 8, 9, 15, 22	$2.5c(s)\sqrt{p/s}$	8, 15, 22	$3c(s)\sqrt{p/s}$	8
DFBETAS <sub><i>i</i>3</sub>	$2/\sqrt{n}$	6, 8, 15	$2.5/\sqrt{n}$	6, 8	$3/\sqrt{n}$	8
GDFBETAS <sub><i>i</i>3</sub>	$2c(s)/\sqrt{s}$	5, 8, 9, 15, 22	$2.5c(s)/\sqrt{s}$	5, 8, 9, 15, 22	$3c(s)/\sqrt{s}$	8, 9, 15, 22
DFBETAS <sub><i>i</i>4</sub>	$2/\sqrt{n}$	6, 8, 15	$2.5/\sqrt{n}$	6, 8	$3/\sqrt{n}$	8
GDFBETAS <sub><i>i</i>4</sub>	$2c(s)/\sqrt{s}$	5, 8, 9, 15, 22	$2.5c(s)/\sqrt{s}$	5, 8, 9, 15, 22	$3c(s)/\sqrt{s}$	8, 9, 15, 22

<sup>a</sup> No observations are claimed as outliers

From the results of SWR<sub>*i*</sub> and GSWR<sub>*i*</sub>, it is shown that diagnostic results based on SWR<sub>*i*</sub> and GSWR<sub>*i*</sub> are similar. Observations 10, 89, 116 and 136 are flagged by SWR<sub>*i*</sub> and GSWR<sub>*i*</sub> as the outlying observations with extra large or small residuals.

Similarly, from the results of LD<sub>*i*</sub> and GLD<sub>*i*</sub>, it is shown that GLD<sub>*i*</sub> and LD<sub>*i*</sub> suggest observations 10, 55, 77, 89, 116, 125, 132, 151, 152 and 164 as the outlying data points that are influential on the regression parameter estimates  $\hat{\beta}_1$  and  $\hat{\beta}_2$ . It is also interesting to note that, on the basis of cut-off values  $9c^2(s)p/s$  and  $9p/n$ , GLD<sub>*i*</sub> suggests cases 89, 116, 151 and 152 as having an influence on  $\hat{\beta}_1$  and  $\hat{\beta}_2$ , while LD<sub>*i*</sub> identifies cases 55, 89, 116 and 152 as influential on the results for  $\hat{\beta}_1$  and  $\hat{\beta}_2$ . The joint deletion of observations 89, 116, 151 and 152 causes the relative changes in  $\hat{\beta}_1$  and  $\hat{\beta}_2$  as 5.21 and 13.80%, whereas the relative changes in  $\hat{\beta}_1$  and  $\hat{\beta}_2$  are 3.47 and 7.74% with

**Table 6** The single-case and group deletion diagnostics for example 2

Diagnostic	Cut-off value	Case <i>i</i> claimed as an outlier	Cut-off value	Case <i>i</i> claimed as an outlier	Cut-off value	Case <i>i</i> claimed as an outlier
SWR <sub><i>i</i></sub>	2	10, 89, 116, 136	2.5	10, 89, 116	3	89
GSWR <sub><i>i</i></sub>	2 <i>c(s)</i>	10, 89, 116, 136	2.5 <i>c(s)</i>	10, 89, 116	3 <i>c(s)</i>	89
LD <sub><i>i</i></sub>	4 <i>p/n</i>	10, 55, 77, 89, 116, 125, 132, 151, 152, 164	6.25 <i>p/n</i>	55, 77, 89, 116, 125, 132, 151, 152, 164	9 <i>p/n</i>	55, 89, 116, 152
GLD <sub><i>i</i></sub>	4 <i>c</i> <sup>2</sup> ( <i>s</i> ) <i>p/s</i>	10, 55, 77, 89, 116, 125, 132, 151, 152, 164	6.25 <i>c</i> <sup>2</sup> ( <i>s</i> ) <i>p/s</i>	55, 77, 89, 116, 125, 132, 152, 164	9 <i>c</i> <sup>2</sup> ( <i>s</i> ) <i>p/s</i>	89, 116, 151, 152
DFFITS <sub><i>i</i></sub>	2√ <i>p/n</i>	10, 55, 77, 89, 116, 125, 132, 151, 152, 164	2.5√ <i>p/n</i>	55, 77, 89, 116, 125, 132, 151, 152, 164	3√ <i>p/n</i>	55, 89, 116, 152
GDFFITs <sub><i>i</i></sub>	2 <i>c(s)</i> √ <i>p/s</i>	10, 51, 55, 77, 89, 116, 125, 132, 133, 151, 152, 164	2.5 <i>c(s)</i> √ <i>p/s</i>	10, 55, 77, 89, 116, 125, 132, 151, 152, 164	3 <i>c(s)</i> √ <i>p/s</i>	55, 77, 89, 116, 125, 132, 152, 164
DFBETAS <sub><i>i</i>2</sub>	2/√ <i>n</i>	55, 77, 89, 116, 125, 132, 151, 152, 164	2.5/√ <i>n</i>	55, 77, 89, 116, 125, 132, 151, 152, 164	3/√ <i>n</i>	55, 77, 89, 116, 125, 132, 151, 152, 164
GDFBETAS <sub><i>i</i>2</sub>	2 <i>c(s)</i> /√ <i>s</i>	10, 51, 55, 77, 89, 116, 117, 125, 132, 133, 151, 152, 164	2.5 <i>c(s)</i> /√ <i>s</i>	51, 55, 77, 89, 116, 125, 132, 151, 152, 164	3 <i>c(s)</i> /√ <i>s</i>	55, 77, 89, 116, 125, 132, 151, 152, 164

observations 55, 89, 116 and 152 eliminated. The indication shows that GLD<sub>*i*</sub> is more effective than LD<sub>*i*</sub>, in pointing out the outlying data points that have jointly bigger impacts on  $\hat{\beta}_1$  and  $\hat{\beta}_2$ .

From the results of DFFITS<sub>*i*</sub> and GDFFITs<sub>*i*</sub>, it is observed that DFFITS<sub>*i*</sub> misses observations 51 and 133 that are discovered by GDFFITs<sub>*i*</sub> as the atypical points on the weighted fits. The joint deletion of observations 10, 55, 77, 89, 116, 125, 132, 151, 152 and 164 leads to the relative variations in  $\hat{\beta}_1$  and  $\hat{\beta}_2$  as 2.60 and 2.00 %, whereas the joint deletion of observations 10, 51, 55, 77, 89, 116, 125, 132, 133, 151, 152 and 164 leads to the relative variations in  $\hat{\beta}_1$  and  $\hat{\beta}_2$  as 2.76 and -1.93 %. It is evident that observations 51 and 133 that are omitted by DFFITS<sub>*i*</sub> are influential, because deleting the atypical points with observations 51 and 133 results in a relative change of  $\hat{\beta}_2$  from 2.00 (positive) to -1.93 % (negative).

From the results of DFBETAS<sub>*i*2</sub> and GDFBETAS<sub>*i*2</sub>, it is shown that GDFBETAS<sub>*i*2</sub> detects observations 10, 51, 55, 77, 89, 116, 117, 125, 132, 133, 151, 152 and 164 as having a large impact on  $\hat{\beta}_2$ , while DFBETAS<sub>*i*2</sub> suggests observations 55, 77, 89, 116, 125, 132, 151, 152 and 164 as influential on the results for  $\hat{\beta}_2$ . When observations 10, 51, 55, 77, 89, 116, 117, 125, 132, 133, 151, 152 and 164 are simultaneously discarded,  $\hat{\beta}_2$  has a positive jump of 0.83 %, whereas the relative change in  $\hat{\beta}_2$  is a 1.51 % reduction, when observations 55, 77, 89, 116, 125, 132, 151, 152 and 164 are

jointly deleted. This indicates that the outlying data points, observations 10, 51, 117 and 133, are influential on  $\widehat{\beta}_2$ , but they are not found by  $\text{DFBETAS}_{i2}$ .

Clearly, in this illustration, the group deletion diagnostics suggest some multiple outlying observations that have joint effects on the regression outcome for further considerations, while these potential outlying data points are camouflaged by the single-case deletion diagnostics.

## 7 Conclusions

In this article, we provide a diagnostic way for the identification of multiple outliers in GBLMs. We suggest the group deletion diagnostic measures,  $\text{GSWR}_i$ ,  $\text{GLD}_i$ ,  $\text{GDFFIT}_i$  and  $\text{GDFBETAS}_i$ , respectively, and propose a test procedure for detecting multiple outlying observations using these. Simulation studies and analysis of two practical examples show that our proposed methods can assist the analyst in detecting multiple outlying observations in GBLMs.

Finally, we note that the GBLMs with the constant precision parameter were recently improved by Simas et al. (2010) who allowed a regression structure for the precision parameter. In future work we will extend the group deletion diagnostic techniques to beta regression models with regressors for the precision parameter.

**Acknowledgments** The author is deeply indebted to the associate editor and two referees for their helpful comments and suggestions that substantially improve this present version of the paper.

## References

- Atkinson AC (1994) Fast very robust methods for the detection of multiple outliers. *J Am Stat Assoc* 89:1329–1339
- Atkinson AC, Riani M (2000) Robust diagnostic regression analysis. Springer, New York
- Atkinson AC, Riani M, Cerioli A (2010) The forward search: theory and data analysis (with discussion). *J Korean Stat Soc* 39:117–134
- Belsley DA, Kuh E, Welsch RE (1980) Regression diagnostics: identifying influential data and sources of collinearity. Wiley, New York
- Cribari-Neto F, Zeileis A (2010) Beta regression in R. *J Stat Softw* 34:1–24
- Espinheira PL, Ferrari SLP, Cribari-Neto F (2008a) On beta regression residuals. *J Appl Stat* 35:407–419
- Espinheira PL, Ferrari SLP, Cribari-Neto F (2008b) Influence diagnostics in beta regression. *Comput Stat Data Anal* 52:4417–4431
- Ferrari SLP, Cribari-Neto F (2004) Beta regression for modeling rates and proportions. *J Appl Stat* 31:799–815
- Hadi AS, Simonoff JS (1993) Procedures for the identification of multiple outliers in linear models. *J Am Stat Assoc* 88:1264–1272
- Imon AHMR (2005) Identifying multiple influential observations in linear regression. *J Appl Stat* 32:929–946
- Imon AHMR, Hadi AS (2008) Identification of multiple outliers in logistic regression. *Commun Stat Theory Methods* 37:1697–1709
- Kieschnick R, McCullough BD (2003) Regression analysis of variates observed on (0, 1): percentage, proportions, fractions. *Stat Model* 3:193–213
- Nurunnabi AAM, Imon AHMR, Nasser M (2010) Identification of multiple influential observations in logistic regression. *J Appl Stat* 37:1605–1624
- Simas AB, Barreto-Souza W, Rocha AV (2010) Improved estimators for a general class of beta regression models. *Comput Stat Data Anal* 54:348–366
- Smithson M, Verkuilen J (2006) A better lemon squeezer? Maximum-likelihood regression with beta-distributed dependent variables. *Psychol Methods* 11:54–71