



ELSEVIER

Available online at www.sciencedirect.com

SCIENCE @ DIRECT®

Computational Statistics & Data Analysis 48 (2005) 755–764

COMPUTATIONAL
STATISTICS
& DATA ANALYSIS

www.elsevier.com/locate/csda

Alternative computational formulae for generalized linear model diagnostics: identifying influential observations with SAS software

John S. Preisser^{a,*}, Daniel I. Garcia^b

^a*Department of Biostatistics, School of Public Health, University of North Carolina,
CB 7420, Chapel Hill, NC 27599, USA*

^b*AstraZeneca, DCC-II E3-961, 1800 Concord Pike, P.O. Box 15437, Wilmington,
DE 19850-5437, USA*

Received 19 December 2003; received in revised form 17 March 2004; accepted 19 March 2004

Abstract

In generalized linear models, regression diagnostics including leverage, DFBETA and Cook's distance are commonly used to assess the influence of observations on the fit of a model. We illustrate how familiarity with the construction of common regression diagnostics formulae can lead to useful alternative formulae when the computer software of interest provides numerical values for only some of the component statistics. In particular, SAS software version 8.2 offers these diagnostics for logistic regression through PROC LOGISTIC, however PROC GENMOD does not compute them, so that, aside from residuals, diagnostics are not directly available from SAS for many generalized linear models. This article describes how these diagnostics may be obtained indirectly with alternative computational formulae based upon observation statistics that are produced as output by PROC GENMOD. Data from the Guidelines for Urinary Incontinence Discussion and Evaluation study, a randomized controlled trial directed at assessing the impact of urinary incontinence guideline adoption by primary care providers on patient outcomes, is used to illustrate the alternative computations.

© 2004 Elsevier B.V. All rights reserved.

Keywords: Cook's Distance; Leverage; Logistic regression

1. Introduction

Any statistical analysis of data should consider the role that influential observations have on the results. For linear regression, exact formulae exist for the change in

* Corresponding author. Tel.: +1-919-968-4460; fax: +1-919-968-3804.

E-mail address: jpreisse@bios.unc.edu (J.S. Preisser).

regression coefficients when an observation is deleted from the data set (Belsley et al., 1980). In contrast, iterative computational routines are required to obtain a solution for most generalized linear models (GLIM) and so exact formulae for the change in regression coefficients are not available. However, simple computational formulae for single-case deletion diagnostics, such as leverage, Cook's distance and DFBETA, provide approximations to the desired quantities (Pregibon, 1981; Williams, 1987). Although these GLIM diagnostics were introduced about two decades ago, their use in practice appears to be more limited than their corresponding linear regression counterparts. One indication of this is that these diagnostics are not computed in the current version of SAS PROC GENMOD (version 8.2, SAS Institute Inc., 1999). This article introduces the use of alternative computational formulae for calculating GLIM diagnostics. The proposed formulae are implemented with SAS PROC IML using observation statistics that are output from PROC GENMOD. The paper provides one illustration of how familiarity with the construction of common regression diagnostics formulae can lead to useful alternative formulae when the computer software of interest provides numerical values for only some of the component statistics, but not for the diagnostics directly.

The illustration involves data (Preisser and Qaqish, 1999) from the Guidelines for Urinary Incontinence Discussion and Evaluation (GUIDE) study, a randomized controlled trial directed at assessing the impact of urinary incontinence (UI) guideline adoption by primary care providers on patient outcomes. UI, or loss of bladder control, is a disorder that affects over 13 million men and women in the United States (Dugan et al., 2001). Although our illustration of the use of PROC GENMOD to obtain the diagnostics is based upon special formulae for a generalized linear model with canonical link function, general procedures for computing diagnostics are also described.

2. Traditional computational formulae for GLIM diagnostics

A generalized linear model has the form $g(\mu_i) = \eta_i$ for $i = 1, \dots, n$ independent observations, where $\eta_i = x_i^T \beta$ is the linear predictor defined by a $p \times 1$ covariate vector x_i and a β parameter vector of interest. The link function, $g(\cdot)$, equates the expected value of the response $\mu_i = E(Y_i)$ with η_i . The mean, $\mu_i = b'(\theta_i)$, expressed as the first derivative of the function $b(\theta_i)$, derives from the random part of the model given by a distribution belonging to the general exponential family (McCullagh and Nelder, 1989)

$$f(Y_i; \theta_i, \phi) = \exp\{(Y_i \theta_i - b(\theta_i))/a(\phi) + c(Y_i, \phi)\}. \quad (1)$$

While θ_i is known as the canonical parameter, a random outcome, Y_i , with sampling distribution belonging to this family has variance that can be expressed as a function of the mean. In particular, $\text{Var}(Y_i) = b''(\theta_i)a(\phi) = v(\mu_i)a(\phi)$, where $v(\mu_i)$ is the variance function, and ϕ is the dispersion parameter. In GLIMs, β , is estimated by maximum likelihood using iterative weighted least squares (IWLS). Asymptotically, $\hat{\beta}$ is normally distributed with mean β and covariance matrix estimated by $\phi(X^T W X)^{-1}$, where W is

a diagonal matrix of $w_i = (\partial\mu_i/\partial\eta_i)^2/v(\mu_i)$. The hat matrix is

$$H = W^{1/2}X(X^TWX)^{-1}X^TW^{1/2}. \quad (2)$$

The property $\text{trace}(H) = p$ leads to various rules of thumb about what is a large leverage value, for instance those greater than twice the mean, $2p/n$. If $g(\cdot)$ is chosen such that $\eta_i = \theta_i$, $g(\cdot)$ is known as the canonical link function. In such cases, W has elements $w_i = [v(\mu_i)]^2/v(\mu_i) = v(\mu_i)$, and (2) simplifies to

$$H = V^{1/2}X(X^TVX)^{-1}X^TV^{1/2}, \quad (3)$$

where $V = \text{Diag}\{v(\mu_i)\}$. For the bernoulli distribution, the mean is equal to $\mu_i = \exp(\theta_i)/[1 + \exp(\theta_i)]$. So if $\eta_i = \log[\mu_i/(1 - \mu_i)]$, the canonical link function for a binomial distribution, then $w_i = v(\mu_i) = \mu_i(1 - \mu_i)$.

Residuals are commonly used to check the adequacy of a regression model. In GLIMs, one of the most commonly used residuals is the Pearson residual,

$$r_{pi} = e_i/v^{1/2}(\hat{\mu}_i), \quad (4)$$

where $e_i = y_i - \hat{\mu}_i$. The studentized form of the Pearson residual is

$$r'_{pi} = e_i/[v(\hat{\mu}_i)\hat{\phi}(1 - h_i)]^{1/2}, \quad (5)$$

where h_i is the i th diagonal element of H .

Case deletion diagnostics approximate the effect of deletion of a case, e.g., an observation, on statistics of interest. The DFBETA_{*i*} diagnostic is a one-step approximation of $\hat{\beta} - \hat{\beta}_{[i]}$, the difference between the estimated regression coefficients based upon the full data set and the estimates when one observation is deleted, with subscript $[i]$ denoting deletion of the i th observation (Williams, 1987). The diagnostic is given by

$$\text{DFBETA}_i = (X^TWX)^{-1}x_i^Tw_i^{1/2}(1 - h_i)^{-1/2}r'_{pi}\phi^{1/2}. \quad (6)$$

While the exact value of $\hat{\beta} - \hat{\beta}_{[i]}$ could be determined by two applications of the IWLS algorithm, this approach becomes computationally intensive since $n + 1$ applications would be required to determine the full set of n statistics (with advances in computing, approaches that require refitting the model, like the diagnostics of Fay (2002), have become feasible in many cases). Instead, expression (6) is used to estimate $\hat{\beta} - \hat{\beta}_{[i]}$ without any further applications, or even without any further iterations of the algorithm. It is called a one-step approximation because it is equivalent to the procedure that upon convergence of the algorithm for the full data set to obtain $\hat{\beta}$, omits the i th observation and applies one more iteration to determine $\hat{\beta}_{[i]}$. However, since all the components of formula (6) are available at convergence of the algorithm applied to the full data set, no more computational effort is required to obtain the full set of n case-deletion diagnostics. For a canonical link function

$$\text{DFBETA}_i = (X^TVX)^{-1}x_i^Tv_i^{1/2}(1 - h_i)^{-1/2}r'_{pi}\phi^{1/2}, \quad (7)$$

where $v_i = v(\mu_i)$. Expression (6) or (7) can be standardized as

$$\text{DFBETAS}_i = \text{DFBETA}_i/\sigma(\hat{\beta}_i) \quad (8)$$

where $\sigma(\hat{\beta}_i)$ is the square root of the appropriate diagonal element of $\phi(X^T W X)^{-1}$. Cook's Distance (D_i) is a scalar measure used to assess the overall fit of a model. For GLIMs, the measure it estimates is given by McCullagh and Nelder (1989, p. 407) as

$$(\hat{\beta} - \hat{\beta}_{[i]})^T (X^T V X) (\hat{\beta} - \hat{\beta}_{[i]}) / p \hat{\phi} \quad (9)$$

which, upon substituting DFBETA_{*i*} for $\hat{\beta} - \hat{\beta}_{[i]}$, gives $D_i = p^{-1} h_i (1 - h_i)^{-1} (r'_{pi})^2$ as reported by Williams (1987).

3. Alternative computational formulae for GLIM diagnostics

Alternative expressions for h_i and DFBETA_{*i*} facilitate their computation in SAS. Given values of r_{pi} and r'_{pi} (e.g., provided by SAS PROC GENMOD), leverage, h_i , can be determined as follows:

$$\begin{aligned} 1 - (r_{pi}/r'_{pi})^2 / \phi &= 1 - \frac{1}{\phi} \left\{ \frac{e_i/[v_i]^{1/2}}{e_i/[v_i \hat{\phi} (1 - h_i)]^{1/2}} \right\}^2 \\ &= 1 - \{[1 - h_i]^{1/2}\}^2 = h_i. \end{aligned} \quad (10)$$

Given h_i and r'_{pi} the standard computational formula for D_i given at the end of Section 2 may be applied. For any GLIM including those with non-canonical link functions, expression (6) may be calculated after first determining w_i as a function of μ_i . Since the predicted means can be obtained from SAS PROC GENMOD, Eq. (6) can be evaluated directly after substituting $\hat{\phi}^{-1} \widehat{\text{cov}}(\hat{\beta})$ for $(X^T W X)^{-1}$. For binary responses, we note that PROC LOGISTIC allows probit and complementary log–log non-canonical link functions. For a model with canonical link, a formula for DFBETA_{*i*} free of h_i and v_i (apart from $\widehat{\text{cov}}(\hat{\beta})$) is available. Substituting the relations for the covariance of $\hat{\beta}$, $v_i = e_i^2 / r_{pi}^2$ determined from (4), and $(1 - h_i)^{1/2} = r_{pi} / (r'_{pi} \sqrt{\phi})$ from (10) into expression (7) yields

$$\begin{aligned} \hat{\beta} - \hat{\beta}_{[i]} &\approx (X^T V X)^{-1} x_i^T \begin{bmatrix} e_i \\ r_{pi} \end{bmatrix} \begin{bmatrix} \hat{\phi}^{1/2} r'_{pi} \\ r_{pi} \end{bmatrix} r'_{pi} \hat{\phi}^{1/2} \\ &= \hat{\phi}^{-1} \widehat{\text{cov}}(\hat{\beta}) x_i^T \begin{bmatrix} e_i (r'_{pi})^2 \\ r_{pi}^2 \end{bmatrix} \hat{\phi} \\ &= \widehat{\text{cov}}(\hat{\beta}) x_i^T \begin{bmatrix} e_i (r'_{pi})^2 \\ r_{pi}^2 \end{bmatrix}. \end{aligned} \quad (11)$$

4. Implementation of alternative formulae in SAS

This section describes the computation of leverage, Cook statistic and DFBETA for GLIMs with canonical link function through PROC GENMOD. Since PROC GENMOD

Table 1
Descriptive analysis for GUIDE data^a

Covariate		Bothered	
		Yes (<i>n</i> = 54)	No (<i>n</i> = 83)
GENDER	Female	43	74
	Male	11	9
AGE	Mean	80.8	81.5
	Standard deviation	3.6	4.0
	Range	(76,90)	(76,91)
WEEKACC	Mean	38.5	12.5
	Standard deviation	29.2	13.2
	Range	(1,117)	(1,65)
SEVERE	1 (least)	3	18
	2	30	56
	3	15	8
	4 (most)	6	1
TOILET	Mean	7.4	5.7
	Standard deviation	3.1	2.7
	Range	(3,20)	(2,20)

^aReproduced from Preisser and Qaqish (1999).

does not directly provide these diagnostics as output, PROC IML is used to perform the matrix arithmetic of expressions (10), thereby obtaining D_i as well, and (11). The code provided below models the binary data from GUIDE using logistic regression. This data has been considered by Preisser and Qaqish (1999) and is available for download at www.bios.unc.edu/~jpreisse/data.htm

The goal of the analysis of the GUIDE data is to identify factors that are predictive of the response to the question: “Do you consider this accidental loss of urine a problem that interferes with your day to day activities or bothers you in other ways?” Let the answer to this question be a binary variable BOTHERED, where BOTHERED = 1 if the patient answers yes and BOTHERED = 0 if the patients answers no. The covariates of interest that will be used to model the variable of interest BOTHERED are the following: WEEKACC, AGE, GENDER, SEVERE, and TOILET. Patients were asked how many leaking accidents they experience in an average week (WEEKACC). When the patient experiences ‘loss of urine’ the variable SEVERE takes the following values: SEVERE=1 if there is only some moisture, SEVERE=2 if the patient wet the underwear, SEVERE=3 if the urine trickled down the thigh, and SEVERE=4 if the patient wet the floor. The variable TOILET measures the number of times during the day they usually go to the toilet to urinate. The variable AGE = (age in yr–76)/10 is a standardized measure of the patient’s age. After omitting 2 patients with missing covariates, 137 elderly patients (age > 76) from 38 medical practices are analyzed, and 54 responded they were bothered by their UI. Table 1 shows the descriptive analysis for the GUIDE study.

The SAS PROC LOGISTIC code that produces the regression diagnostics, DFBETA, leverage and Cook's statistic, for the logistic regression model is

```
proc logistic data=UI descending;
    model bothered = female age dayacc severe toilet;
    output out=logcook c=cstat dfbetas=all H=H;
run;
```

(12)

Cook's distance D_i is obtained by dividing $cstat$ by $p\phi$, where $p = 6$ and $\phi = 1$ for the GUIDE data. The $dfbetas$ option provides the standardized version given in expression (8), and option H gives the leverage values. Equivalently, the diagnostics can be obtained using output from PROC GENMOD:

```
ods output covB=covB ParameterEstimates=ParameterEstimates;
proc genmod data=UI;
    model bothered = female age dayacc severe toilet/dist=bin link=logit r covB;
    output out=gencook resraw=resraw reschi=reschi stdreschi=stdreschi;
run;
```

(13)

where $resraw$, $reschi$ and $stdreschi$ refer to e_i , r_{pi} and r'_{pi} , respectively. The appendix gives PROC IML code to calculate $DFBETAS_i$, leverage and Cook's statistic from the three types of residuals and the covariance matrix of the regression parameter estimates provided in the output data set $covB$. For logistic regression $\phi = 1$.

5. Illustration with GUIDE study data

Regression diagnostics are produced for the logistic regression model specified by the SAS code given by (12) and (13). Table 2 shows the parameter estimates for the complete data set of 137 observations and for a data set obtained by deleting observations from three patients. A few observations appear to have a large impact on the results, particularly, for TOILET.

For the data set based upon all 137 observations, PROC GENMOD and PROC LOGISTIC gave identical results to three decimal places for all parameter estimates and standard errors. By default PROC GENMOD uses the observed fisher information in the fitting algorithm, whereas PROC LOGISTIC uses the expected fisher information by default (Fisher scoring). So the identical results are expected because, for canonical link models, the observed equals the expected information, and so Newton–Raphson and Fisher scoring are equivalent methods (McCullagh and Nelder, 1989).

The figure shows the leverage (h_i) and Cook's Distance values for the 137 observations in GUIDE. Two patients appear to have high leverage values relative to the remaining patients, and two patients appear to have large Cook's Distance. Patient 44 and Patient 122, have a measure of leverage equal to 0.27496 and 0.16069, respectively. Patient 44 is a 77 year old female who reported a high frequency of weekly

Table 2

Parameter estimates and standard errors for urinary incontinence for all 137 observations and for the model that has patients 8, 44 and 122 removed from the data set

	All observations ($n = 137$)		3 observations removed	
	Estimate	Standard error	Estimate	Standard error
INTERCEPT	−3.2929	1.1083	−3.4115	1.1905
FEMALE	−0.6723	0.6116	−1.0386	0.6614
AGE	−0.6405	0.5847	−1.0293	0.6285
DAYACC	0.4154	0.0958	0.4903	0.1112
SEVERE	0.8285	0.3645	0.6341	0.4128
TOILET	0.1108	0.0855	0.2625	0.1071

Results were obtained from SAS PROC LOGISTIC. Using all 137 observations, GENMOD gave estimate for INTERCEPT of −3.2930 and for AGE of −0.6406; all other estimates and standard errors were identical to four decimal places.

Table 3

Selected values of leverage and Cook's distance

Patient ID	Leverage		Cook's Statistic	
	LOGISTIC	GENMOD	LOGISTIC	GENMOD
8	0.08174	0.08173	0.18487	0.18488
44	0.27497	0.27496	0.25942	0.25944
122	0.16068	0.16069	0.02704	0.02704

leaking accidents (21) and toileting (20 times a day), but also reported that she was not bothered by her UI. Since the patient reported such a high frequency of toileting and leaking we would expect her to be bothered by her UI. Correspondingly, patient 44 reports the largest value of Cook's Distance in the data set. Similarly, Patient 8, a 77 year old male, with a large value of D_i reported a relatively high frequency of weekly leaking accidents (65) and toileting (10 times a day), while reporting not being bothered by his UI. The data from patient 8 illustrates that extreme values of x_i do not necessarily imply a large leverage in logistic regression; see Hosmer and Lemeshow (1989) for further discussion of this phenomenon. In fact, the leverage for patient 8 is less than $2p/n$. Conversely, the value of leverage for patient 122 was relatively large, while its Cook's distance was not. While reporting the most severe accidents (SEVERE = 4) without being bothered by her UI, this 79 year old female reported a low frequency of weakly accidents (1) and toileting (4 times per day).

Table 3 summarizes leverage and Cook's Distance for the observations labelled in Fig. 1. Alternative computational formula in the appendix were used to obtain the results labelled GENMOD. These values differed from those obtained directly from PROC LOGISTIC only in the fifth decimal place, thus illustrating the equivalency of the alternative computational formulae to the standard formulae.

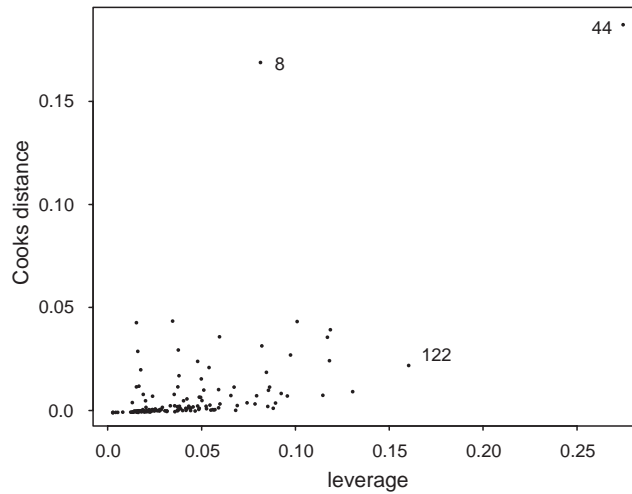


Fig. 1. GUIDE: logistic regression diagnostics.

Table 4
Selected values of DFBETAS

Variable	Patient 8		Patient 44		Patient 122	
	LOGISTIC	GENMOD	LOGISTIC	GENMOD	LOGISTIC	GENMOD
Female	0.51162	0.51163	0.00668	0.00668	−0.02953	−0.02953
Age	0.37937	0.37938	0.14742	0.14742	0.02065	0.02065
Dayacc	−0.62834	−0.62836	0.01141	0.01140	0.13116	0.13116
Severe	0.48629	0.48629	0.22566	0.22566	−0.34967	−0.34968
Toilet	−0.36007	−0.36007	−1.19649	−1.19653	0.13671	0.13671

Table 4 reports standardized $DFBETAS_i$ for the observations labelled in the figure. Again, equivalency of the two approaches is illustrated. The largest effect of a single-case deletion upon an individual beta coefficient pertains to the effect of patient 44 on the value for TOILET; as noted earlier this patient reported the maximum frequency of toileting.

6. Discussion

Alternative computational formulae for deletion diagnostics from generalized linear models are useful in many instances when SAS software is employed. While their use was illustrated with logistic regression for which SAS PROC LOGISTIC provides direct results, the formulae apply to other GLIMs such as poisson regression for which direct results are currently unavailable in SAS PROC GENMOD (version 8.2). In the case of non-canonical link GLIMs, the proposed diagnostic formulae apply to a broad assortment of generalized linear and quasi-likelihood models including those characterized

by a family of link functions (Nelder and Pregibon, 1987; Preisser, 2002). The lack of availability of the statistics, leverage, DFBETA, and Cook's distance, in particular, motivated the presentation. More broadly, we have illustrated how familiarity with the construction of common regression diagnostic formulae can lead to the determination of alternative formulae useful when the computer software of interest provides numerical values for only some of the component statistics. Finally, the broadest aim of this article was to encourage the use of regression diagnostics for generalized linear models.

Deletion diagnostics similar to those discussed here exist for clustered data problems analyzed with generalized estimating equations (Preisser and Qaqish, 1996; Hardin and Hilbe, 2003). In that situation, computational formulae for both observation- and cluster-deletion diagnostics are available. The GUIDE study discussed in the present article is from a survey of 137 individuals from 38 different practices. The clustering was ignored for the purpose of illustrating GLIM diagnostics. Preisser and Qaqish (1999), however, report on GEE diagnostics for the GUIDE data based upon applying GEE with an exchangeable correlation matrix. Regression diagnostics for GLIMs are equivalent to GEE diagnostics that use a model-based covariance estimator and a working independence correlation matrix.

Appendix.

The following PROC IML code calculates the diagnostics using output from (13):

```
proc iml;
use data_set;
read all var {inter female age dayacc severe toilet} into x;
read all var {bothered} into y;
read all var {resraw} into resraw;
read all var {reschi} into reschi;
read all var {stdreschi} into stdreschi;
use covB;
read all var {Prm1 Prm2 Prm3 Prm4 Prm5 Prm6 Prm7} into CovB;
use ParameterEstimates;
read all var {estimate} into estimate;
```

After data are read by IML, matrices are then formed for use in computing regression diagnostics:

```
p=ncol{x};
n=nrow{x};
scale=estimate[p+1,]##2;
rr=(reschi/stdreschi)##2;
invrr=1/rr;
stddev=sqrt(vecdiag(covB));
tx=t(x);
```

Finally, the observation diagnostics are calculated. These include a vector of the observation leverages, a vector of Cook's Distance for the observations, a $p \times n$ matrix of the DFBETA diagnostics, and their standardized version, DFBETAS. In DFBETA, the i th column contains the vector of diagnostics approximating $\hat{\beta} - \hat{\beta}_{[i]}$.

```
lev=J(n,1,1)- rr/scale;
D=(1/p)#(lev/(1- lev))#(stdreschi##2);
DFBETA=covB* TX* Diag(invRR#resraw);
DFBETAS=DFBETA#(stddev##- 1);
```

References

- Belsley, D.A., Kuh, E., Welsch, R.E., 1980. Regression Diagnostics. Wiley, New York.
- Dugan, E., Roberts, C.P., Cohen, S.J., Preisser, J.S., Davis, C.C., Bland, D.R., Albertson, E., 2001. Why older community-dwelling adults do not discuss urinary incontinence with their primary care physicians. *J. Amer. Geriatrics Soc.* 49, 462–465.
- Fay, M.P., 2002. Measuring a binary response's range of influence in logistic regression. *Amer. Statist.* 56, 5–9.
- Hardin, J.W., Hilbe, J.M., 2003. Generalized Estimating Equations. Chapman & Hall, London.
- Hosmer, D.W., Lemeshow, S., 1989. Applied Logistic Regression. Wiley, New York.
- McCullagh, P., Nelder, J.A., 1989. Generalized Linear Models, 2nd Edition. Chapman & Hall, London.
- Nelder, J.A., Pregibon, D., 1987. An extended quasi-likelihood function. *Biometrika* 74, 221–232.
- Pregibon, D., 1981. Logistic regression diagnostics. *Ann. Statist.* 9 (4), 705–724.
- Preisser, J.S., 2002. Quasi-likelihood analysis of patient satisfaction with medical care. *Health Services Outcomes Res. Methodol.* 3, 233–245.
- Preisser, J.S., Qaqish, B.F., 1996. Deletion diagnostics for generalized estimating equations. *Biometrika* 83, 551–562.
- Preisser, J.S., Qaqish, B.F., 1999. Robust regression for clustered data with application to binary responses. *Biometrics* 55, 574–579.
- SAS Institute Inc., 1999. SAS/STAT User's Guide, Version 8. SAS Institute Inc., Cary, NC, 3884pp.
- Williams, D.A., 1987. Generalized linear model diagnostics using the deviance and single-case deletions. *Appl. Statist.* 36, 181–191.