

Small-sample performance of the robust score test and its modifications in generalized estimating equations

Xu Guo^{1,3}, Wei Pan^{1,*†}, John E. Connett¹, Peter J. Hannan² and Simone A. French²

¹*Division of Biostatistics, School of Public Health, University of Minnesota, A460 Mayo Building (MMC 303), Minneapolis, MN 55455, U.S.A.*

²*Division of Epidemiology, School of Public Health, University of Minnesota, A460 Mayo Building (MMC 303), Minneapolis, MN 55455, U.S.A.*

³*Clinical Genomics R&D, Affymetrix, Inc. Santa Clara, CA 95051, U.S.A.*

SUMMARY

The sandwich variance estimator of generalized estimating equations (GEE) may not perform well when the number of independent clusters is small. This could jeopardize the validity of the robust Wald test by causing inflated type I error and lower coverage probability of the corresponding confidence interval than the nominal level. Here, we investigate the small-sample performance of the robust score test for correlated data and propose several modifications to improve the performance. In a simulation study, we compare the robust score test to the robust Wald test for correlated Bernoulli and Poisson data, respectively. It is confirmed that the robust Wald test is too liberal whereas the robust score test is too conservative for small samples. To explain this puzzling operating difference between the two tests, we consider their applications to two special cases, one-sample and two-sample comparisons, thus motivating some modifications to the robust score test. A modification based on a simple adjustment to the usual robust score statistic by a factor of $J/(J-1)$ (where J is the number of clusters) reduces the conservativeness of the generalized score test. Simulation studies mimicking group-randomized clinical trials with binary and count responses indicated that it may improve the small-sample performance over that of the generalized score and Wald tests with test size closer to the nominal level. Finally, we demonstrate the utility of our proposal by applying it to a group-randomized clinical trial, trying alternative cafeteria options in schools (TACOS). Copyright © 2005 John Wiley & Sons, Ltd.

KEY WORDS: bias-correction; correlated data; GEE; robust score test; robust Wald test; sandwich variance estimator

*Correspondence to: Wei Pan, Division of Biostatistics, School of Public Health, University of Minnesota, A460 Mayo Building (MMC 303), Minneapolis, MN 55455, U.S.A.

†E-mail: weip@biostat.umn.edu

Contract/grant sponsor: NIH; contract/grant numbers: R01-HL65462, 5210-HL59275, N01-HR-96140, R18-HL61305
Contract/grant sponsor: Minnesota Medical Foundation

1. INTRODUCTION

Correlated data often arise from biomedical research due to repeated measurements on the same individual or clustered sampling. For example, in the trying alternative cafeteria options in schools (TACOS) study, a group-randomized trial (GRT), 20 secondary schools were randomly allocated to control and intervention [1, 2]. Student survey data were collected from a random sample of students from each school to evaluate the effect of the intervention. A within-school correlation for the student responses occurs when the outcomes within the same schools are more similar to each other than the outcomes from different schools. Appropriate statistical analyses need to take account of the within-school correlation. Two most popular and general statistical methods that can account for within-group (i.e. within-school in the TACOS study) correlations in GRT data are mixed-effects models and generalized estimating equations (GEE) [3, 4]. A challenge in analysing GRT data is how to handle a small number of groups; see two recent reviews and references therein [5, 6].

GEE [7, 8] has been widely used to model correlated data and make proper statistical inference. When the number of independent clusters is sufficiently large, GEE method has some desirable properties. The regression coefficient estimates from GEE are consistent and asymptotically normal. Their covariance is consistently estimated by the robust variance estimator in a sandwich form (thus also called as the sandwich estimator), which is robust to the misspecification of the covariance of the correlated responses. Using a chi-squared reference distribution, the generalized or robust Wald test and the generalized or robust score test can be applied to test a hypothesis about a subset of regression coefficients [9].

However, the sandwich estimator does not perform well for small samples [10], which has been confirmed by empirical studies [11–13]. This has jeopardized the validity of the robust Wald test by causing inflated type I errors relative to specified nominal levels, and thus lowering coverage probabilities of the corresponding confidence intervals. To be more explicit, a small sample refers to a small number of independent clusters instead of a small number of observations for each cluster in this paper. Recently, the small-sample inference for correlated data has become an active research area and drawn much attention; Feng and Braun [14] gave an updated review. For small-sample inference in GEE, most efforts focus on improving the small-sample performance of the sandwich estimator to control the size of the resulting robust Wald test [15]. In general, there are two ways [16, 17]. One is to correct the bias of the sandwich estimator, and the other is to take account of its variability. Since the sandwich estimator is downward biased in estimating the covariance of the regression coefficient estimates and the bias is large for small samples, Mancl and DeRouen [18] proposed to use a bias-corrected sandwich estimator, rather than the usual sandwich estimator, to improve the performance of the robust Wald test. Another approach was developed by Pan and Wall [19]. Since the high variability of the sandwich estimator directly affects the size of the robust Wald test and the coverage probability of resulting confidence interval [20], Pan and Wall [19] took account of the variability of the sandwich estimator and constructed an approximate *t*-test for testing a single parameter and an *F*-test for multiple parameters, which were guaranteed to reduce the inflated size of the robust Wald test.

In addition to the robust Wald test, another important robust test is the generalized or robust score test [21]. Breslow [22] derived a robust score test for overdispersed independent data; in contrast to the liberal behaviour of the robust Wald test, the robust score test was shown to be conservative with a small sample size. Although the robust score test has long appeared

in some applications for correlated data [23–25], its performance for small sample correlated data has not been fully investigated, which may have caused its lack of use in some popular statistical packages that implement GEE, such as S-Plus. As to be shown later, our simulation studies confirm that such operating characteristics of the robust Wald test and score test will persist in GEE for correlated data with small sample sizes.

In Section 2, we briefly review the GEE method. We introduce the robust Wald test and its bias-corrected modification [18]. We also introduce the robust score test for correlated data. In Section 3, we conduct simulations to evaluate the small-sample performance of the robust score test. In Section 4, we apply the two robust statistics to some special cases that help explain the operating characteristics of the robust score test and thus motivate the modifications proposed in Section 5. We also consider constructing a confidence interval by inverting the robust score statistic in Section 6. Section 7 provides some comparison with the F -test. We illustrate our method using the data from the TACOS study in Section 8 and end with a short discussion in Section 9.

2. GENERALIZED ESTIMATING EQUATIONS

Suppose we have a correlated data set with J clusters. For each cluster j ($j = 1, \dots, J$), there is a response vector, $Y_j = (y_{j1}, \dots, y_{jn_j})^T$, and an $n_j \times p$ covariate matrix X_j . y_{jk} 's ($k = 1, \dots, n_j$) are assumed correlated within each cluster and independent across clusters. The marginal expectation of the response, $E(y_{jk}) = \mu_{jk}$, is modelled by a regression equation $g(\mu_{jk}) = X_{jk}\beta$, where $\beta = (\beta_1, \dots, \beta_p)^T$ is a p -dimensional vector of unknown regression coefficients and $g(\cdot)$ is a known link function. The marginal variance is $\text{Var}(y_{jk}) = v(\mu_{jk})\phi$, where v is a known variance function and ϕ is a scale parameter. The within-cluster correlation matrix $R_{0j} = \text{corr}(Y_j)$ is generally unknown. We can consistently estimate β by specifying a common working correlation matrix $R_w(\alpha)$, which may depend on some parameter α , and solving the following generalized estimating equations:

$$U(\beta, \phi) = \sum_{j=1}^J D_j^T V_j^{-1} S_j = 0 \quad (1)$$

where $D_j = \partial \mu_j / \partial \beta^T$, $\mu_j = (\mu_{j1}, \dots, \mu_{jn_j})^T$, $V_j = \phi A_j^{1/2} R_w(\alpha) A_j^{1/2}$, $A_j = \text{diag}(v(\mu_{j1}), \dots, v(\mu_{jn_j}))$, and $S_j = Y_j - \mu_j$.

Under mild regularity conditions, $\hat{\beta}$ is consistent and asymptotically normal [7]. $\text{Cov}(\hat{\beta})$ can be consistently estimated by the so-called sandwich estimator,

$$V_s = \left(\sum_{j=1}^J D_j^T V_j^{-1} D_j \right)^{-1} \left(\sum_{j=1}^J D_j^T V_j^{-1} \text{Cov}(Y_j) V_j^{-1} D_j \right) \left(\sum_{j=1}^J D_j^T V_j^{-1} D_j \right)^{-1} \quad (2)$$

where $\text{Cov}(Y_j) = E(S_j S_j^T)$.

2.1. Robust Wald test

To test $H_0: \beta_2 = \beta_2^0$ versus $H_1: \beta_2 \neq \beta_2^0$, we decompose the p -dimensional regression coefficient vector β as $(\beta_1^T, \beta_2^T)^T$, where β_1 and β_2 are vectors of dimensions of p_1 and p_2 respectively. Similarly, suppose that $\hat{V}_{s(2)}$ is the submatrix of \hat{V}_s for β_2 . Without loss of generality, we

consider testing $H_0: \beta_2 = \beta_2^0$ versus $H_1: \beta_2 \neq \beta_2^0$. The robust Wald statistic

$$W = (\hat{\beta}_2 - \beta_2^0)^T \hat{V}_{S(2)}^{-1} (\hat{\beta}_2 - \beta_2^0) \quad (3)$$

has an asymptotic chi-squared distribution $\chi_{p_2}^2$ under H_0 . When the number of clusters is large, valid statistical inference can be accomplished. However, the robust Wald test does not work well for small samples. The test sizes are inflated and the corresponding confidence intervals have lower coverage probabilities.

In V_S , $\hat{S}_j \hat{S}_j^T$ is used to estimate $\text{Cov}(Y_j)$. It is well known that \hat{S}_j is downward biased in estimating S_j and the bias is large when the number of clusters is small. Hence, the sandwich estimator of $\hat{\beta}$ is downward biased and can cause inflated test sizes for the robust Wald test. Since $E(\hat{S}_j \hat{S}_j^T) \doteq (I_j - H_{jj}) \text{Cov}(Y_j) (I_j - H_{jj}^T)$, Mancl and DeRouen [18] proposed to use $(I_j - H_{jj})^{-1} \hat{S}_j \hat{S}_j^T (I_j - H_{jj}^T)^{-1}$ to estimate $\text{Cov}(Y_j)$ to correct the downward bias of the sandwich estimator V_S , where I_j is an $n_j \times n_j$ identity matrix and $H_{jj} = D_j (\sum_{j=1}^J D_j^T V_j^{-1} D_j)^{-1} D_j^T V_j^{-1}$. The bias-corrected sandwich estimator is denoted as V_{BC} ,

$$\begin{aligned} V_{BC} &= \left(\sum_{j=1}^J D_j^T V_j^{-1} D_j \right)^{-1} \left(\sum_{j=1}^J D_j^T V_j^{-1} (I_j - H_{jj})^{-1} \hat{S}_j \hat{S}_j^T (I_j - H_{jj}^T)^{-1} V_j^{-1} D_j \right) \\ &\quad \times \left(\sum_{j=1}^J D_j^T V_j^{-1} D_j \right)^{-1} \end{aligned} \quad (4)$$

The bias-corrected robust Wald statistic is

$$W_{BC} = (\hat{\beta}_2 - \beta_2^0)^T \hat{V}_{BC(2)}^{-1} (\hat{\beta}_2 - \beta_2^0) \quad (5)$$

2.2. Robust score test

To test $H_0: \beta_2 = \beta_2^0$ versus $H_1: \beta_2 \neq \beta_2^0$, corresponding to the decomposition of regression coefficients β to $(\beta_1^T, \beta_2^T)^T$, we can decompose the generalized estimating equation (1) as $U(\beta_1, \beta_2, \phi) = (U_{(1)}^T(\beta_1, \beta_2, \phi), U_{(2)}^T(\beta_1, \beta_2, \phi))^T$. Under H_0 , we can solve $U_{(1)}(\beta_1, \beta_2^0, \phi) = 0$ to obtain an estimate of β_1 , denoted as $\tilde{\beta}_1$. A_{11} , A_{12} , A_{21} , A_{22} , B_{11} , B_{12} , B_{21} and B_{22} are the corresponding decomposed submatrices of A and B , where $A = \lim_J J^{-1} A_J$ is the normed limit of

$$A_J = -E \frac{\partial U}{\partial \beta^T} = \sum_{j=1}^J D_j^T V_j^{-1} D_j \quad (6)$$

$B = \lim_J J^{-1} B_J$ is the normed limit of

$$B_J = E U U^T = \sum_{j=1}^J D_j^T V_j^{-1} \text{Cov}(Y_j) V_j^{-1} D_j \quad (7)$$

Using the Taylor expansion, one can establish that $(1/\sqrt{J})U_{(2)}(\tilde{\beta}_1, \beta_2^0, \phi)$ follows an asymptotic multivariate normal distribution with zero mean and covariance matrix

$$\text{Cov}[(1/\sqrt{J})U_{(2)}(\tilde{\beta}_1, \beta_2^0, \phi)] \doteq C \text{Cov}[(1/\sqrt{J})U(\beta_1, \beta_2^0, \phi)] C^T \quad (8)$$

where $C = (-A_{21}A_{11}^{-1}, I)$ and I is a $p_2 \times p_2$ identity matrix. In practice, $\text{Cov}[(1/\sqrt{J})U(\beta_1, \beta_2^0, \phi)]$ is typically estimated by $\sum_{j=1}^J \tilde{U}_j \tilde{U}_j^T / J$, where \tilde{U}_j is $U_j = D_j^T V_j^{-1} S_j$ evaluated at $(\tilde{\beta}_1, \beta_2^0)$. The robust score statistic for testing H_0 is

$$S = \left[\tilde{C} \sum_{j=1}^J \tilde{U}_j \right]^T \left[\tilde{C} \left(\sum_{j=1}^J \tilde{U}_j \tilde{U}_j^T \right) \tilde{C}^T \right]^{-1} \left[\tilde{C} \sum_{j=1}^J \tilde{U}_j \right] \quad (9)$$

where \tilde{C} is evaluated at $(\tilde{\beta}_1, \beta_2^0)$. It follows an asymptotic chi-squared distribution $\chi_{p_2}^2$ under H_0 .

3. SIMULATION STUDIES

Simulations are conducted to study the small-sample performance of the robust score test for correlated Bernoulli and Poisson data, respectively. Correlated responses can be easily generated from a generalized normal random-effects model, which can be well approximated by a corresponding generalized marginal model [26, 27]. For each data configuration, 1000 simulations are generated. For each simulation, the corresponding generalized marginal regression model is fitted using either the independence or the compound symmetry working correlation matrix in GEE. To facilitate comparisons, we include both the robust score test and the robust Wald test for testing null hypotheses involving a single parameter or multiple parameters. For simplicity, we only consider cases with an equal cluster size; later, we apply the robust score test to the TACOS data with unequal cluster sizes.

3.1. Correlated Bernoulli responses

Correlated responses, y_{jk} 's, are randomly generated from $\text{Bin}(1, \mu_{jk})$, where μ_{jk} are generated from the following normal random-effects logistic model:

$$\text{logit}(\mu_{jk} | b_j) = \beta_0 + x_{1jk}\beta_1 + x_{2jk}\beta_2 + x_{3jk}\beta_3 + b_j \quad (10)$$

where $\beta_0 = \beta_1 = \beta_2 = \beta_3 = 0$, $j = 1, \dots, J$ ($J = 10, 20, 30$), and $k = 1, \dots, 50$. The covariates x_{1jk} , x_{2jk} and x_{3jk} are all iid from a Bernoulli distribution $\text{Bin}(1, \frac{1}{2})$, b_j 's are iid from $N(0, \frac{1}{5})$, and they are independent of each other. The intra-cluster correlation (ICC) is about 0.05; see appendix for a derivation. Both individual and joint hypothesis tests with $H_0: \beta_1 = 0$ and $H_0: \beta_1 = \beta_2 = \beta_3 = 0$ are considered for simulated data.

A marginal logistic regression model $\text{logit}(\mu_{jk}) = \beta_0 + x_{1jk}\beta_1 + x_{2jk}\beta_2 + x_{3jk}\beta_3$ was fitted using GEE for each simulated data set. The size of a test is evaluated as the observed fraction of times when the null hypothesis is rejected while the null hypothesis is true. The sizes of the robust Wald, bias-corrected Wald, and robust score tests are shown in Table I. (Please ignore the columns S' and S'' for modified score tests, which will be discussed later.) When J is 10, the size of the robust Wald test is 0.096 for independence working correlation and 0.087 for compound symmetry working correlation, which are almost twice of the nominal level of 0.05. As J increases to 30, it gets closer to 0.05. For the robust score test, when J is 10, the size is 0.036 for independence working correlation and 0.034 for compound symmetry working correlation, which are smaller than 0.05. As J increases, it approaches to 0.05.

Table I. Empirical size for testing $H_0: \beta_1 = \beta_2 = \beta_3 = 0$ or $H_0: \beta_1 = 0$ at the nominal level of 0.05 in a mixed-effects logistic regression model for correlated Bernoulli responses.

Tests		Independence					Compound symmetry				
H_0	J	W	W_{BC}	S	S'	S''	W	W_{BC}	S	S'	S''
$\beta_1 = 0$	10	0.096	0.061	0.036	0.095	0.051	0.087	0.058	0.034	0.088	0.045
$\beta_1 = 0$	20	0.067	0.050	0.043	0.067	0.047	0.069	0.057	0.046	0.068	0.053
$\beta_1 = 0$	30	0.067	0.059	0.049	0.068	0.054	0.066	0.055	0.047	0.066	0.052
$\beta_1 = \beta_2 = \beta_3 = 0$	10	0.222	0.167	0.015	0.234	0.036	0.214	0.162	0.012	0.218	0.038
$\beta_1 = \beta_2 = \beta_3 = 0$	20	0.122	0.097	0.040	0.123	0.051	0.124	0.099	0.039	0.130	0.049
$\beta_1 = \beta_2 = \beta_3 = 0$	30	0.100	0.084	0.050	0.102	0.055	0.108	0.091	0.050	0.106	0.054

For joint hypothesis $H_0: \beta_1 = \beta_2 = \beta_3 = 0$, the size of the robust Wald test or its bias-corrected version is more inflated than that for individual hypothesis $H_0: \beta_1 = 0$ (Table I). For example, when J is 10, the size of the robust Wald test is 0.222 for independence working correlation and 0.214 for compound symmetry working correlation. Even when J increases to 30, the robust Wald test still has inflated size of 0.100 for independence working correlation and 0.108 for compound symmetry working correlation. The bias-corrected Wald test still has much inflated type I error rates. In contrast, the robust score test is more conservative for joint hypothesis tests. When J is 10, it has size of 0.015 for independence working correlation and 0.012 for compound symmetry working correlation, and its test size approaches to 0.05 as J increases. In summary, the robust Wald test and its bias-corrected version may have inflated sizes when J is small, and the problem is more severe for testing a joint hypothesis, which is consistent with previous empirical studies [18]; the robust score test is conservative when J is small, even more so for a joint hypothesis.

3.2. Correlated Poisson responses

Correlated Poisson responses, y_{jk} 's, are randomly generated from $\text{Poisson}(\mu_{jk})$, where μ_{jk} are generated from the following normal random-effects Poisson regression model:

$$\log(\mu_{jk} | b_j) = \beta_0 + x_{1jk}\beta_1 + x_{2jk}\beta_2 + x_{3jk}\beta_3 + b_j \quad (11)$$

where $\beta_0 = \beta_1 = \beta_2 = \beta_3 = 0$, $j = 1, \dots, J$ ($J = 10, 20, 30$), $k = 1, \dots, 50$. The covariates x_{1jk} , x_{2jk} and x_{3jk} are all iid from a Bernoulli distribution $\text{Bin}(1, \frac{1}{2})$, b_j 's are iid from $N(0, \frac{1}{19})$, and they are independent of each other. The intra-cluster correlation (ICC) is about 0.05; see appendix for a derivation.

A marginal Poisson regression model $\log(\mu_{jk}) = \beta_0 + x_{1jk}\beta_1 + x_{2jk}\beta_2 + x_{3jk}\beta_3$ was fitted using GEE for each simulated data set. The results are presented in Table II. (Again please ignore the columns S' and S'' for modified score tests, which will be discussed later.) The same conclusion on the operating characteristics of the various tests can be drawn as that for correlated binary data.

Table II. Empirical size for testing $H_0: \beta_1 = \beta_2 = \beta_3 = 0$ or $H_0: \beta_1 = 0$ at the nominal level of 0.05 in a mixed-effects Poisson regression model for correlated Poisson responses.

Tests		Independence					Compound symmetry				
H_0	J	W	W_{BC}	S	S'	S''	W	W_{BC}	S	S'	S''
$\beta_1 = 0$	10	0.103	0.069	0.037	0.099	0.058	0.095	0.068	0.043	0.099	0.055
$\beta_1 = 0$	20	0.067	0.054	0.042	0.068	0.049	0.067	0.050	0.039	0.068	0.047
$\beta_1 = 0$	30	0.061	0.057	0.052	0.064	0.055	0.067	0.062	0.051	0.066	0.059
$\beta_1 = \beta_2 = \beta_3 = 0$	10	0.233	0.170	0.007	0.222	0.027	0.228	0.166	0.007	0.223	0.021
$\beta_1 = \beta_2 = \beta_3 = 0$	20	0.152	0.109	0.030	0.143	0.041	0.146	0.112	0.030	0.144	0.044
$\beta_1 = \beta_2 = \beta_3 = 0$	30	0.100	0.083	0.039	0.093	0.046	0.102	0.079	0.042	0.099	0.046

4. THEORY FOR TWO SPECIAL CASES

It is natural to ask why the robust Wald test is too liberal and at the same time the robust score test is too conservative for small samples. Due to the lack of explicit and simple formulas of the two test statistics in general, we cannot compare them directly and provide an intuitive explanation. Here, we apply them to the special cases of one- and two-sample problems for testing one and two means, respectively, where we have simple forms of the two statistics. Their direct comparison provides an interesting explanation to the small-sample operating characteristics of the two tests, and thus motivates modifications to the robust score test to improve its small-sample performance for correlated data.

4.1. One-sample problem

Here we consider testing for a mean based on a given sample. The related robust Wald test and robust score test have been discussed in the literature under different contexts, such as Reference [28]. Here for completeness, we give a brief introduction. Given a sample y_1, y_2, \dots, y_J with the assumptions that

- (i) y_1, y_2, \dots, y_J are mutually independent, and
- (ii) $E(y_j) = \mu$ for $j = 1, 2, \dots, J$,

we want to test the null hypothesis $H_0: \mu = \mu_0$. Without loss of generality, we can take $\mu_0 = 0$ and use the alternative hypothesis $H_1: \mu \neq 0$. Note that we do not need to specify a full distribution for each y_j . In fact, we do not even need to assume that y_j 's are from the same population. For instance, it is possible that y_j 's do not have a common variance.

Inference on μ can be accomplished by solving the following estimating equation

$$U(\mu) = \sum_{j=1}^J w_j(y_j - \mu) = 0 \quad (12)$$

where w_j is the weight for the j th observation. For simplicity, we use $w_j = 1/\sigma^2$ with $\text{Var}(y_j) = \sigma^2$ under a working model, which may not be true. The (square-root of) robust

Wald statistic is

$$W = \frac{\bar{y}}{\sqrt{1/J^2 \sum_{j=1}^J (y_j - \bar{y})^2}} \quad (13)$$

and the (square-root of) robust score statistic is

$$S = \frac{U(0)}{\sqrt{\text{Var}(U(0))}} = \frac{\bar{y}}{\sqrt{1/J^2 \sum_{j=1}^J y_j^2}} \quad (14)$$

They both follow an asymptotic standard normal distribution $N(0, 1)$ under H_0 . Equivalently, their squares follow an asymptotic chi-squared distribution χ_1^2 under H_0 .

Since $\sum_{j=1}^J (y_j - \bar{y})^2 = \sum_{j=1}^J y_j^2 - J\bar{y}^2$, $\sqrt{1/J^2 \sum_{j=1}^J y_j^2} \geq \sqrt{1/J^2 \sum_{j=1}^J (y_j - \bar{y})^2}$. Thus $|W| \geq |S|$, where the equality holds if and only if $\bar{y} = 0$; i.e. $S = W = 0$. The extra term $1/J\bar{y}^2$ in the denominator of the robust score statistic provides an explicit explanation of why the robust score test tends to be more conservative than the robust Wald test. Furthermore, the extra term introduces a positive correlation between the numerator and the denominator of the robust score statistic for small samples, whereas the numerator and the denominator of the robust Wald test statistic are (nearly) uncorrelated (at least under the normality assumption for y_j 's). This positive correlation pushes the robust score test to be conservative.

4.2. Two-sample problem

Suppose we have two independent samples, $\{y_j, j = 1, \dots, J_1\}$ and $\{y_j, j = J_1 + 1, \dots, J_1 + J_2\}$. The sample means and variances are $\bar{y}_1 = 1/J_1 \sum_{j=1}^{J_1} y_j$, $\bar{y}_2 = 1/J_2 \sum_{j=J_1+1}^{J_1+J_2} y_j$, $s_1^2 = 1/(J_1 - 1) \sum_{j=1}^{J_1} (y_j - \bar{y}_1)^2$ and $s_2^2 = 1/(J_2 - 1) \sum_{j=J_1+1}^{J_1+J_2} (y_j - \bar{y}_2)^2$. We assume that

- (i) $y_1, \dots, y_{J_1}, y_{J_1+1}, \dots, y_{J_1+J_2}$ are mutually independent,
- (ii) $E(y_j) = \mu_1$ for $j = 1, \dots, J_1$, and
- (iii) $E(y_j) = \mu_2$ for $j = J_1 + 1, \dots, J_1 + J_2$.

Again note that we do not need to assume that all the y_j 's in each sample are iid. The goal is to test whether the two means are the same, i.e. $H_0: \mu_1 = \mu_2$ versus $H_1: \mu_1 \neq \mu_2$.

The estimating equations to be used are $U(\beta) = (U_{(1)}(\beta), U_{(2)}(\beta))^T = 0$ with

$$U_{(1)}(\beta) = \sum_{j=1}^{J_1} w_j(y_j - \mu_1) + \sum_{j=J_1+1}^{J_1+J_2} w_j(y_j - \mu_2) = 0$$

$$U_{(2)}(\beta) = \sum_{j=1}^{J_1} w_j(y_j - \mu_1) = 0$$

where $\beta = (a, b)^T$, $\mu_1 = a + b$, $\mu_2 = a$. Then our goal is to test $H_0: b = 0$ against $H_1: b \neq 0$. For simplicity, we take the weights to be $w_j = 1/\sigma_1^2$ for $j = 1, \dots, J_1$ and $w_j = 1/\sigma_2^2$ for $j = J_1 + 1, \dots, J_1 + J_2$, where σ_1^2 and σ_2^2 can be the assumed variances for the two samples under a working model, which may not be true.

The (square-root of) robust Wald statistic is derived as

$$W = \frac{\bar{y}_1 - \bar{y}_2}{\sqrt{[(J_1 - 1)/J_1^2]s_1^2 + [(J_2 - 1)/J_2^2]s_2^2}} \quad (15)$$

which follows an asymptotic standard normal distribution $N(0, 1)$ under H_0 . The robust Wald statistic has a similar form to the two-sample t statistic (with unequal variances) because $[(J_1 - 1)/J_1^2]s_1^2 = (1/J_1)\tilde{s}_1^2$ with $\tilde{s}_1^2 = \sum_{j=1}^{J_1} (y_j - \bar{y}_1)^2/J_1$, the biased sample variance estimator, and similarly for the second sample. This fact has been pointed out in the literature (e.g. Reference [10]).

Under $H_0: b = 0$, the (square-root of) robust score statistic is derived as

$$S = \frac{1/\sigma_1^2 \sum_{j=1}^{J_1} (y_j - \tilde{a})}{\sqrt{1/\sigma_1^4 (1 - [J_1\sigma_2^2/(J_1\sigma_2^2 + J_2\sigma_1^2)])^2 \sum_{j=1}^{J_1} (y_j - \tilde{a})^2 + 1/\sigma_2^4 ([J_1\sigma_2^2/(J_1\sigma_2^2 + J_2\sigma_1^2)])^2 \sum_{j=J_1+1}^{J_1+J_2} (y_j - \tilde{a})^2}} \quad (16)$$

where $\tilde{a} = (J_1\sigma_2^2\bar{y}_1 + J_2\sigma_1^2\bar{y}_2)/(J_1\sigma_2^2 + J_2\sigma_1^2)$. It has an asymptotic standard normal distribution $N(0, 1)$ under H_0 .

Under the working assumption of $\sigma_1^2 = \sigma_2^2$, (16) becomes

$$S_1 = \frac{\bar{y}_1 - \bar{y}_2}{\sqrt{[(J_1 - 1)/J_1^2]s_1^2 + [(J_2 - 1)/J_2^2]s_2^2 + (J_1^2 - J_1J_2 + J_2^2)/(J_1J_2(J_1 + J_2))(\bar{y}_1 - \bar{y}_2)^2}} \quad (17)$$

If we do not assume $\sigma_1^2 = \sigma_2^2$ and estimate σ_1^2 and σ_2^2 with s_1^2 and s_2^2 , respectively, we have

$$S_2 = \frac{\bar{y}_1 - \bar{y}_2}{\sqrt{[(J_1 - 1)/J_1^2]s_1^2 + [(J_2 - 1)/J_2^2]s_2^2 + (J_1^3s_2^4 + J_2^3s_1^4)/(J_1J_2(J_1s_2^2 + J_2s_1^2)^2)(\bar{y}_1 - \bar{y}_2)^2}} \quad (18)$$

When we compare the robust score statistic S_1 or S_2 to the robust Wald statistic W , we notice that they look similar except that there is an extra non-negative term in the denominator of S_1 or S_2 , $(J_1^2 - J_1J_2 + J_2^2)/(J_1J_2(J_1 + J_2))(\bar{y}_1 - \bar{y}_2)^2$ or $(J_1^3s_2^4 + J_2^3s_1^4)/(J_1J_2(J_1s_2^2 + J_2s_1^2)^2)(\bar{y}_1 - \bar{y}_2)^2$. Each term is zero if and only if $\bar{y}_1 = \bar{y}_2$, in which case the two robust statistics are equivalent. Otherwise, these two terms are always positive, and are of the order $O(J^{-1})$ if $J_1 = J_2 = J$. Therefore, as the sample size increases, the two robust tests are more and more similar, in accordance with the theoretical asymptotic equivalence. However, when the sample size is small, the extra terms in the robust score statistics can have significant influence on the testing results: they introduce a positive correlation between the numerator and the denominator of the robust score statistic for small samples and cause the test to be conservative.

5. MODIFIED ROBUST SCORE TESTS

From the simulation study based on correlated Bernoulli or Poisson responses, we can see that the robust Wald test has inflated test sizes, whereas the robust score test has smaller test sizes than the specified nominal level. Neither of them has satisfactory performance when J is small. Mancini and DeRouen [18] found that the bias-correction method was helpful to improve the small-sample property of the sandwich estimator to bring down the inflated test size of the robust Wald test. But this method cannot be directly applied to the robust score test to improve its small-sample performance since the bias-corrected robust score test will

be even more conservative. Similarly, extending the idea of Pan and Wall [19] to adopt a t - or F -distribution as the reference for the robust score statistic does not work because the resulting test will be even more conservative.

A problem with the robust score test is that, as shown for one- and two-sample test situations, the numerator and the denominator of the robust score statistic are correlated, which may have caused it to be conservative for small-samples. In the robust score statistic (9), the correlation between $\sum_{j=1}^J \tilde{U}_j$ and $\sum_{j=1}^J \tilde{U}_j \tilde{U}_j^T$ may be high for small samples and thus cause the robust score test to be conservative with small samples. Therefore, we try to modify the robust score statistic by reducing the correlation to achieve better small-sample performance. A simple method to reduce the correlation is to use the sample variance estimator $\sum_{j=1}^J (\tilde{U}_j - \bar{U})(\tilde{U}_j - \bar{U})^T$ with $\bar{U} = \sum_{j=1}^J \tilde{U}_j / J$, rather than $\sum_{j=1}^J \tilde{U}_j \tilde{U}_j^T$, to estimate $\text{Cov}(U)$. Note that we have $E(U) = 0$ under H_0 . Thus, under H_0 , the two covariance estimators are essentially the same for large J . However, for small J , the use of the correspondingly modified robust score statistic

$$S' = \left[\tilde{C} \sum_{j=1}^J \tilde{U}_j \right]^T \left[\tilde{C} \left(\sum_{j=1}^J (\tilde{U}_j - \bar{U})(\tilde{U}_j - \bar{U})^T \right) \tilde{C}^T \right]^{-1} \left[\tilde{C} \sum_{j=1}^J \tilde{U}_j \right] \quad (19)$$

leads to quite different operating characteristics. The size of the robust score test S' is inflated when J is small, for both correlated Bernoulli data (Table I) and Poisson data (Table II). As J increases, its test size approaches the nominal level of 0.05 and closer to the size of the usual robust score test S . In summary, S' over-corrects the conservativeness of the robust score test S , and similar to the robust Wald test, it behaves too liberally for small J .

The modified robust score test S' shows inflated rather than deflated type I error rate. Next, we propose a modification that has a better performance. We have

$$\begin{aligned} \sum_{j=1}^J (\tilde{U}_j - \bar{U})(\tilde{U}_j - \bar{U})^T &= \sum_{j=1}^J \tilde{U}_j \tilde{U}_j^T - \frac{1}{J} \left(\sum_{j=1}^J \tilde{U}_j \tilde{U}_j^T + \sum_{j \neq l} \tilde{U}_j \tilde{U}_l^T \right) \\ &= \frac{J-1}{J} \sum_{j=1}^J \tilde{U}_j \tilde{U}_j^T - \frac{1}{J} \sum_{j \neq l} \tilde{U}_j \tilde{U}_l^T \end{aligned}$$

Since $E(U_j U_l) = 0$ for any $j \neq l$ under H_0 , we have that the last term in the above equation is close to zero, and thus can be ignored. This leads to using $(J-1)/J \sum_{j=1}^J \tilde{U}_j \tilde{U}_j^T$ to estimate $\text{Cov}(U)$. When J is small, the factor $(J-1)/J$ can help improve the performance of the robust score test. When J is large, it is almost the same as $\sum_{j=1}^J \tilde{U}_j \tilde{U}_j^T$. The corresponding robust score statistic is

$$S'' = \left[\tilde{C} \sum_{j=1}^J \tilde{U}_j \right]^T \left[\tilde{C} \left(\frac{J-1}{J} \sum_{j=1}^J \tilde{U}_j \tilde{U}_j^T \right) \tilde{C}^T \right]^{-1} \left[\tilde{C} \sum_{j=1}^J \tilde{U}_j \right] \quad (20)$$

which has the same value as $J/(J-1)S$. When we apply this proposed modification to the robust score statistic (14) for the special case of one-sample problem, we get

$$S'' = \frac{\bar{y}}{\sqrt{(J-1)^2/J^4 \sum_{j=1}^J y_j^2}}$$

Note that there is a key difference between S' and S or S'' : as we have shown at the end of Section 4.1, for the one-dimensional case under the Normality, using $\sum_j (U_j - \bar{U})^2$ to estimate the variance of U_j 's eliminates the positive correlation between the numerator and the denominator of S' , whereas using $\sum_j U_j^2$ in S or S'' retains the positive correlation, leading to that S or S'' tends to be smaller than S' . In addition, the adjustment factor $(J-1)/J$ in S'' was not motivated from the point of correcting the bias of the covariance estimator; as mentioned earlier, such a bias-correction for the score test will lead to even a more conservative test.

S'' has better performance than S or S' when J is small, and similar performance to S or S' when J is large. For the individual hypothesis test with $H_0: \beta_1 = 0$ for the correlated Bernoulli data (Table I), when J is 10, the test size of S is 0.036 for independence working correlation and 0.034 for compound symmetry correlation; the test size of S'' is 0.051 for independence working correlation and 0.045 for compound symmetry correlation. For the joint hypothesis test $H_0: \beta_1 = \beta_2 = \beta_3 = 0$, the improvement is more obvious (Table I). When J is 10, the test size of S is 0.015 for independence working correlation and 0.012 for compound symmetry working correlation; the test size of S'' is 0.036 for independence working correlation and 0.038 for compound symmetry working correlation. We also see improvement for the correlated Poisson data. For the individual hypothesis test $H_0: \beta_1 = 0$, when J is 10, the test size of S is 0.037 for independence working correlation and 0.043 for compound symmetry correlation; the test size of S'' is 0.058 for independence working correlation and 0.055 for compound symmetry correlation (Table II). For the joint hypothesis test $H_0: \beta_1 = \beta_2 = \beta_3 = 0$, when J is 10, the test size of S is 0.007 for both working correlation structures; the test sizes of S'' are 0.027 and 0.021, respectively (Table II).

In addition, the robust score test S'' performs better than the bias-corrected robust Wald test W_{BC} . The bias-correction alone is not enough to bring down the inflated test size of the robust Wald test to the nominal level; it still has inflated test size even when J is 30. For instance, for the joint hypothesis test of the correlated Bernoulli data ($J=30$), the size of the bias-corrected robust Wald test is 0.084 for independence working correlation and 0.091 for compound symmetry working correlation; the sizes of S'' are 0.055 and 0.054, respectively (Table I). For the joint hypothesis test of the correlated Poisson data, the size of the bias-corrected robust Wald test is 0.083 for independence working correlation and 0.079 for compound symmetry working correlation; the sizes of S'' are both 0.046 (Table II).

The test results based on the two working correlation matrices are similar. For simplicity, we only use independence working correlation in the following sections.

6. INTERVAL ESTIMATION

Confidence intervals for the regression coefficients can be easily constructed from the robust Wald statistic W or W_{BC} . By inverting the robust score statistic S or S'' , we can also construct the corresponding confidence interval. Since the confidence interval based on the robust Wald statistic is often too narrow, we enlarge it generously by taking five times the robust standard error of the estimated regression coefficient to obtain a reasonable starting point to find the score-test-based confidence interval. We then use the bisection algorithm [29] to search within this enlarged interval for the two endpoints of the confidence interval based on the robust score statistic S or S'' . We compared 95 per cent equal-tail confidence intervals based on the robust Wald statistic W , the bias-corrected robust Wald statistic W_{BC} , the robust score

Table III. Average length and coverage probability of 95 per cent equal-tail confidence intervals based on the statistics W , W_{BC} , S and S'' for correlated Bernoulli responses.

		W				S			
Set-up		Average	Average	Average	Coverage	Average	Average	Average	Coverage
β_1	J	lower bound	upper bound	length	probability	lower bound	upper bound	length	probability
0	10	-0.421	0.422	0.843	0.910	-0.545	0.545	1.090	0.961
0	20	-0.307	0.307	0.614	0.933	-0.343	0.343	0.686	0.955
0	30	-0.198	0.196	0.394	0.950	-0.212	0.210	0.422	0.959

		W_{BC}				S''			
β_1	J	Average	Average	Average	Coverage	Average	Average	Average	Coverage
		lower bound	upper bound	length	probability	lower bound	upper bound	length	probability
0	10	-0.476	0.476	0.952	0.944	-0.500	0.500	1.000	0.945
0	20	-0.326	0.326	0.652	0.946	-0.333	0.332	0.665	0.950
0	30	-0.205	0.203	0.408	0.955	-0.208	0.206	0.414	0.957

statistic S and modified robust score statistic S'' for correlated Bernoulli data generated from the following normal random-effects logistic model:

$$\text{logit}(\mu_{jk} | b_j) = \beta_0 + x_{jk}\beta_1 + b_j \quad (21)$$

where $\beta_0 = 0$, $j = 1, \dots, J$ ($J = 10, 20, 30$), and $k = 1, \dots, 50$. The covariate x_{jk} is from $\text{Bin}(1, \frac{1}{2})$, b_j 's are iid from $N(0, \frac{1}{5})$, and they are independent of each other. For each data configuration, 1000 simulations are run to obtain the average length of the intervals and their coverage probabilities (Table III). The coverage probability for the robust Wald confidence interval is smaller than the specified confidence level of 0.95, whereas it is larger than 0.95 for the confidence interval based on the robust score statistic S . By using the bias-correction method, the interval based on W is appropriately widened to accomplish better coverage probability. The method we proposed can also adjust the interval based on the robust score statistic S to obtain better coverage probability.

7. OTHER COMPARISONS

Here we briefly compare the performance of our proposed robust score test with that of an F -test, which has been shown to improve over the robust Wald test [19]. The test statistic of the F -test is the same as the robust Wald statistic, however, instead of using a Chi-squared distribution, it uses an F -distribution as its reference distribution to calculate the P -value. To facilitate the comparison, we did a simulation study with the set-up as that in Reference [19]. Specifically, the data were generated from a normal random-effects logistic model as in (10) except that (1) the cluster size is 20 (i.e. $k = 1, \dots, 20$) and (2) the distribution of the random effect b_j is $N(0, 1)$, introducing a larger ICC. The empirical test sizes using the independence working model are shown in Table IV, where the results for the F -test are taken from Reference [19]. It can be seen that sometimes the F -test seems slightly more conservative than the modified score test S'' , and in general, the modified score test is competitive.

Table IV. Empirical size for testing $H_0: \beta_1 = \beta_2 = \beta_3 = 0$ or $H_0: \beta_1 = 0$ at the nominal level of 0.05 in a mixed-effects logistic regression model for correlated Bernoulli responses.

H_0	J	W	W_{BC}	S	S'	S''	F
$\beta_1 = 0$	10	0.093	0.056	0.044	0.091	0.053	0.036
$\beta_1 = 0$	20	0.068	0.050	0.035	0.069	0.048	0.032
$\beta_1 = 0$	30	0.062	0.049	0.043	0.064	0.049	0.044
$\beta_1 = 0$	40	0.064	0.057	0.053	0.065	0.055	0.056
$\beta_1 = \beta_2 = \beta_3 = 0$	10	0.201	0.131	0.006	0.216	0.025	0.034
$\beta_1 = \beta_2 = \beta_3 = 0$	20	0.116	0.087	0.037	0.128	0.047	0.026
$\beta_1 = \beta_2 = \beta_3 = 0$	30	0.098	0.075	0.040	0.099	0.047	0.026
$\beta_1 = \beta_2 = \beta_3 = 0$	40	0.079	0.068	0.048	0.078	0.052	0.036

8. EXAMPLE

We illustrate the application of the robust score test S'' using the data taken from the TACOS study, which was a two-year, group-randomized, school-based nutrition intervention trial to evaluate an environmental intervention to increase the sales of lower fat foods in secondary school cafeteria a la carte areas [1, 2]. At the baseline (Fall 2000), 20 secondary schools were randomly assigned to either a control group with no intervention or to an environmental intervention group with 10 schools in each group. The intervention consisted of increasing the availability of lower fat foods in the cafeteria a la carte areas and peer influence *via* school-wide, student-based promotions for these lower fat foods. Student survey data were collected at the baseline (Fall 2000) and Spring 2001 through mail surveys of a random sample of students from each school with respect to students' food choices, attitudes, perceived environment and behavioural intentions, demographic variables and, etc. It turned out that increasing the availability and promotion of lower fat foods had a positive effect on secondary school students' food purchase and their perceptions about the food environment at school.

Here, we analyse the survey data using GEE to study whether the environmental intervention has any effect on students' weight. The number of observations from different schools are different, which ranges from 64 to 129. The estimated within-school correlation is about 0.07. We consider the following covariates, Grade, Sex (= 1 for males; = 0 for females), Trt (= 1 for intervention; = 0 for control), Race (7 indicator variables, Ami for American Indians, Blk for Blacks, Mex for Mexicans, Ric for Puerto Ricans, Asn for Asians, Oth for others, and Mlt for being in more than two of the categories; the reference category is White), Time (= 1 for Spring 2001; = 0 for Fall 2000). The interaction between Time and Trt is not statistically significant (not shown).

The robust Wald test, the robust score test and the modified robust score test are conducted to study whether the covariates have any effect on the response individually or jointly. The estimated regression coefficients and P -values are summarized in Table V. As we have expected, the P -values from S'' are smaller than those from S . This is consistent with the results from simulations. The individual effects of Grade, Sex, Asn and Time are statistically significant at the nominal level of 0.05 according to the results based on the modified robust score test. The effect of Race, manifested by testing the seven coefficients, is statistically significant based on the robust Wald test (P -value < 0.001), but not based on the robust score

Table V. *P*-values of the various tests for the regression coefficients for the TACOS data. *n* is the total number of students in each group at the baseline or follow-up.

Tested covariate	Coefficient estimate	<i>W</i>	<i>S</i>	<i>S''</i>
Grade	6.016	<0.001	<0.001	<0.001
Sex	29.152	<0.001	<0.001	<0.001
Trt	1.346	0.448	0.454	0.442
Time	4.462	<0.001	0.002	0.002
Race	—	<0.001	0.314	0.279
White (<i>n</i> = 1838)	—	—	—	—
Ami (<i>n</i> = 37)	2.672	0.628	0.635	0.627
Blk (<i>n</i> = 64)	5.737	0.002	0.064	0.058
Mex (<i>n</i> = 34)	−3.976	0.407	0.417	0.405
Ric (<i>n</i> = 13)	−6.669	0.312	0.313	0.300
Asn (<i>n</i> = 104)	−18.140	<0.001	0.007	0.005
Oth (<i>n</i> = 68)	−1.208	0.760	0.761	0.755
Mlt (<i>n</i> = 11)	−5.713	0.508	0.534	0.523

test (*P*-value = 0.314), highlighting the operating difference between the two tests. Because of the liberal behaviour of the robust Wald test and the conservative behaviour of the robust score test, the truth probably lies between the two, as indicated by the modified robust score test (*P*-value = 0.279).

Note that, although all the three tests give a statistically significant difference between whites and Asian, it is possible that an overall test on Race is not significant because of the non-significance between whites and other groups; in particular, we note that most of the other groups have much smaller sample sizes.

9. DISCUSSION

The robust Wald test is constructed from the comparison between regression coefficient estimates and their estimated covariance. It performs well when sample size is sufficiently large. Due to its easy use and wide availability in SAS, S-Plus and other popular statistical packages, it has been widely, and indeed almost exclusively, used to conduct hypothesis testing and construct confidence intervals. When the sample size is small, the robust Wald test does not perform well. Our simulation studies have confirmed that it has dramatically inflated type I error and the corresponding confidence interval has lower coverage probability than a specified confidence level. In contrast, the robust score test is too conservative for small sample sizes. In this paper, the operating characteristics of the robust Wald test and the robust score test are confirmed in GEE analyses for correlated data. It is also demonstrated that the bias-correction in the robust Wald test can alleviate, but not completely avoid, the problem. Considerations of two special cases, one- and two-sample comparisons, shed light on the operating difference between the two tests. We have proposed and studied several modifications to the robust score test, among which a simple adjustment to the usual robust score statistic by a factor of $J/(J-1)$ (where *J* is the number of clusters) gives test size closest to the specified nominal level in our simulations. However, there are rooms for improvement. For example, our modified robust score test may still be conservative when the number of the groups or clusters is

small. In addition, the independence between the adjustment factor $J/(J-1)$ and cluster size implies that the adjustment is probably over-simplified. Therefore, a more rigorous justification on the use of the adjustment factor $J/(J-1)$ or other modifications will be helpful.

In contrast to extensive studies on the score test with small samples for independent data (e.g. Reference [30] and references therein), the robust score test for correlated data has received scant consideration; it remains an open question whether some ideas can be borrowed from that for independent data [31] and extended successfully to the context for correlated data.

APPENDIX A

A.1. Derivation of ICC for a mixed-effects logistic regression model

According to model (10), we have

$$y_{ij} | b_i \sim \text{Bin}(1, \pi_i), \quad \pi_i = 1/(1 + e^{-b_i}), \quad b_i \sim N(0, \sigma^2)$$

Using a first-order Taylor approximation, we have

$$E(y_{ij}) = E[E(y_{ij} | b_i)] = E(\pi_i) = E\left(\frac{1}{1 + e^{-b_i}}\right) \approx \frac{1}{2}$$

and

$$\begin{aligned} \text{Var}(y_{ij}) &= E[\text{Var}(y_{ij} | b_i)] + \text{Var}[E(y_{ij} | b_i)] = E[\pi_i(1 - \pi_i)] + \text{Var}(\pi_i) \\ &= E(\pi_i)[1 - E(\pi_i)] \approx 1/4 \end{aligned}$$

For any $j \neq k$,

$$E(y_{ij}y_{ik}) = E[E(y_{ij}y_{ik} | b_i)] = E(\pi_i^2)$$

and

$$\text{Cov}(y_{ij}, y_{ik}) = E(y_{ij}y_{ik}) - E(y_{ij})E(y_{ik}) = \text{Var}(\pi_i)$$

Using the Delta-method, we have

$$\text{Var}(\pi_i) = \text{Var}\left(\frac{1}{1 + e^{-b_i}}\right) \approx \left(\frac{e^{-b_i}}{(1 + e^{-b_i})^2}\right)^2 \bigg|_{b_i=0} \text{Var}(b_i) = \sigma^2/16$$

Hence the intra-class correlation

$$\text{ICC} = \frac{\text{Cov}(y_{ij}, y_{ik})}{\text{Var}(y_{ij})} \approx \sigma^2/4$$

A.2. Derivation of ICC for a mixed-effects Poisson regression model

According to model (11), we have

$$y_{ij} | b_i \sim \text{Pois}(\pi_i), \quad \pi_i = e^{b_i}, \quad b_i \sim N(0, \sigma^2)$$

Using the same method as shown in A.1, we have

$$\text{Var}(y_{ij}) = E[\text{Var}(y_{ij} | b_i)] + \text{Var}[E(y_{ij} | b_i)] = E(\mu_i) + \text{Var}(\mu_i) \approx 1 + \sigma^2$$

$$\text{Cov}(y_{ij}, y_{ik}) = E(y_{ij}y_{ik}) - E(y_{ij})E(y_{ik}) = E(\mu_i^2) - [E(\mu_i)]^2 = \text{Var}(\mu_i) \approx \sigma^2$$

Hence, $\text{ICC} \approx \sigma^2/(1 + \sigma^2)$.

ACKNOWLEDGEMENTS

We thank the reviewers and the editor for many helpful and constructive comments. X.G. and W.P. were supported in part by NIH grant R01-HL65462 and a Minnesota Medical Foundation grant, X.G. and J.E.C. by NIH grant 5210-HL59275 and contract N01-HR-96140, and P.J.H. and S.A.F. by NIH grant R18-HL61305.

REFERENCES

1. French SA, Story M, Fulkerson JA. An environmental intervention to promote lower fat food choices in secondary schools: outcomes of the TACOS study. *American Journal of Public Health* 2004; **94**:1507–1512.
2. French SA, Story M, Fulkerson JA, Gerlach AF. Food environment in secondary schools: a la carte, vending machines, and food policies and practices. *American Journal of Public Health* 2003; **93**:1–7.
3. Murray DM. *Design and Analysis of Group-Randomization Trials*. Oxford University Press: New York, 1988.
4. Donner A, Klar N. *Design and Analysis of Cluster Randomization Trials in Health Research*. Arnold: London, 2000.
5. Preisser JS, Young ML, Zaccaro DJ, Wolfson M. An integrated population-averaged approach to the design, analysis and sample size determination of cluster-unit trials. *Statistics in Medicine* 2003; **22**:1235–1254.
6. Murray DM, Varnell SP, Blitstein JL. Design and analysis of group-randomized trials: a review of recent methodological developments. *American Journal of Public Health* 2004; **94**:423–432.
7. Liang KY, Zeger SL. Longitudinal data analysis using generalized linear models. *Biometrika* 1986; **73**:13–22.
8. Prentice RL, Zhao LP. Estimating equations for parameters in means and covariances of multivariate discrete continuous responses. *Biometrics* 1991; **47**:825–839.
9. Rotnitzky A, Jewell NP. Hypothesis testing of regression parameters in semiparametric generalized linear models for cluster data. *Biometrika* 1990; **77**:485–489.
10. Drum M, McCullagh P. Comment. *Statistical Science* 1993; **8**:300–301.
11. Emrich LJ, Piedmonte MR. On some small sample properties of generalized estimating equation estimates for multivariate dichotomous outcomes. *Journal of Statistical Computation and Simulation* 1992; **41**:19–29.
12. Gunsolley JC, Getchell C, Chinchilli VM. Small sample characteristics of generalized estimating equations. *Communications in Statistics—Simulation and Computation* 1995; **24**:869–878.
13. Pan W. On the robust variance estimator in generalized estimating equations. *Biometrika* 2001; **88**:901–906.
14. Feng ZD, Braun TM. Small sample inference for clustered data. *Proceedings of the Seattle Symposium in Biostatistics*, 2002.
15. Daniels MJ, Kass RE. Shrinkage estimators for covariance matrices. *Biometrics* 2001; **57**:1173–1184.
16. Fay MP, Graubard BI. Small-sample adjustments for Wald-type tests using sandwich estimators. *Biometrics* 2001; **57**:1198–1206.
17. Fay MP, Graubard BI, Freedman LS, Midthune DN. Conditional logistic regression with sandwich estimators: application to meta analysis. *Biometrics* 1998; **54**:195–208.
18. Mancl LA, DeRouen TA. A covariance estimator for GEE with improved small-sample properties. *Biometrics* 2001; **57**:126–134.
19. Pan W, Wall M. Small-sample adjustments in using the sandwich variance estimator in generalized estimating equations. *Statistics in Medicine* 2002; **21**:1429–1441.
20. Kauermann G, Carroll RJ. A note on the efficiency of sandwich covariance matrix estimation. *Journal of the American Statistical Association* 2001; **96**:1387–1396.
21. Dennis BD. On generalized score tests. *The American Statistician* 1992; **46**(4):327–333.
22. Breslow N. Tests of hypotheses in overdispersed Poisson regression and other quasi-likelihood models. *Journal of the American Statistical Association* 1990; **85**:565–571.
23. Lefkopoulou M, Ryan L. Global tests for multiple binary outcomes. *Biometrics* 1993; **49**:975–988.

24. Legler J, Lefkopoulou M, Ryan L. Efficiency and power of tests for multiple binary outcomes. *Journal of the American Statistical Association* 1995; **90**:680–693.
25. Barnhart HX, Williamson JM. Goodness-of-fit tests for GEE modeling with binary responses. *Biometrics* 1998; **54**:720–729.
26. Zeger SL, Liang KY, Albert PS. Models for longitudinal data: a generalized estimating equation approach. *Biometrics* 1988; **44**:1049–1060.
27. Neuhaus JM, Kalbfleisch JD, Hauck WW. A comparison of cluster-specific and population-averaged approaches for analysing correlated binary data. *International Statistical Review* 1991; **59**:25–35.
28. Viraswami K, Reid N. Higher-order asymptotics under model misspecification. *The Canadian Journal of Statistics* 1996; **24**:263–278.
29. Press WH, Teukolsky SA, Vetterling WT, Flannery BP. *Numerical Recipes in C, The Art of Scientific Computing* (2nd edn). Cambridge: New York, 1992.
30. Cordeiro GM, Botter DA, Barroso CP, Ferrari SLP. Three corrected score tests for generalized linear models with dispersion covariates. *Statistica Neerlandica* 2003; **57**:391–409.
31. Cordeiro GM. Improved score tests for generalized linear models. *Journal of the Royal Statistical Society, Series B* 1993; **55**:661–674.