



Criteria to Select a Working Correlation Structure for the Generalized Estimating Equations Method in SAS

Masahiko Gosho
Aichi Medical University

Abstract

The generalized estimating equations (GEE) method is popular for analyzing clustered and longitudinal data. It is important to determine a proper working correlation matrix when applying the GEE method since an improper selection sometimes results in inefficient parameter estimates. In this paper, we provide the `CriteriaWorkCorr` macro in SAS to calculate the criteria proposed by Pan (2001), Hin, Carey, and Wang (2007), Hin and Wang (2009), and Gosho, Hamada, and Yoshimura (2011) for selecting the working correlation structure when the GEE method is applied. We illustrate the implementation and an example of the macro.

Keywords: generalized estimating equations, SAS, working correlation structure.

1. Introduction

The generalized estimating equations (GEE) method is one of the most popular ways to analyze clustered and longitudinal data. To apply the GEE method, a working correlation structure—independent, exchangeable, and first-order autoregressive (AR(1))—must be specified. If the working correlation structure is correctly specified, the GEE provides a best asymptotically normal (BAN) estimator of mean parameters. Fitzmaurice (1995) and Wang and Carey (2003) show, however, that the asymptotic relative efficiency of the parameter estimates of the GEE method is likely to be low when the working correlation structure is misspecified. Fitzmaurice (1995), Mancl and Leroux (1996), and Sutradhar and Das (2000) also point out that the misspecification of the correlation structure lowers the relative efficiency of the estimate even when the sample size is finite. To address this concern, some researchers have proposed new criteria to select a working correlation structure. Pan (2001) proposes a modification of Akaike’s information criterion (AIC) called the “quasi-likelihood under the

independence model criterion (QIC).” In addition, [Hin *et al.* \(2007\)](#) apply a method by [Rotnitzky and Jewell \(1990\)](#) as a criterion to select the working correlation structure. [Rotnitzky and Jewell \(1990\)](#) describe an approach to appraise the adequacy of the assumed correlation matrix using the fact that the asymptotic distribution of a modified working Wald statistic is the linear combination of independent χ_1^2 random variables. We call this criterion “Rotnitzky and Jewell’s criterion (RJC).” Furthermore, [Hin and Wang \(2009\)](#) propose a correlation information criterion (CIC) that modifies QIC and substantially improves its performance. Moreover, [Gosho *et al.* \(2011\)](#) devise an objective criterion for evaluating the appropriateness of the correlation structure. The proposed criterion measures the discrepancy between the covariance matrix estimator and the specified working correlation matrix. Hereafter, this criterion is referred as DEW. [Hin *et al.* \(2007\)](#), [Hin and Wang \(2009\)](#), and [Gosho *et al.* \(2011\)](#) compare the performances of these criteria for selecting the working correlation structure.

For most longitudinal data in biological applications, [Wang and Carey \(2003\)](#), [Ziegler and Vens \(2010\)](#), and [Vens and Ziegler \(2012\)](#) show that the AR(1) structure is preferable over banded correlation structures, e.g., the 1-dependent correlation structure. Furthermore, m -dependent correlation structures are not biologically plausible. [Ziegler and Vens \(2010\)](#) and [Vens and Ziegler \(2012\)](#) also point out that the investigators should choose a working correlation structure for both statistical and biological reasons. The statistical criteria for selecting the working correlation structure can be helpful tools to decide the most reasonable structure for the investigators.

In this paper, we present a SAS ([SAS Institute Inc. 2003a](#)) macro to calculate the values of the criteria (QIC, RJC, CIC, and DEW) for selecting the working correlation structure when the GEE method is applied to longitudinal data. The SAS macro (`CriteriaWorkCorr`) was developed and tested on Microsoft Windows XP and 7 operating systems and requires SAS 9.1 (at least the SAS/BASE, [SAS Institute Inc. 2003b](#), the SAS/IML, [SAS Institute Inc. 2003c](#), and the SAS/STAT, [SAS Institute Inc. 2003d](#), components) or above.

2. Summary of the generalized estimating equations method

Assume that an $n_i \times p$ matrix of covariate values $\mathbf{X}_i = (\mathbf{x}_{i1}, \dots, \mathbf{x}_{in_i})^\top$ is adjoined to the outcome vector $\mathbf{Y}_i = (Y_{i1}, \dots, Y_{in_i})^\top$ on clusters $i = 1, \dots, K$ and observations $t = 1, \dots, n_i$ per cluster. To simplify notation, we suppose that $n_i = n$. The expected value and variance of the outcome variable are assumed to be $\mu_{it} = \mathbf{E}(Y_{it}|\mathbf{x}_{it}) = h^{-1}(\mathbf{x}_{it}^\top \boldsymbol{\beta})$ and $\text{VAR}(Y_{it}|\mathbf{x}_{it}) = \phi v(\mu_{it})$, respectively, where h is a specified link function, $\boldsymbol{\beta}$ is a regression parameter (p -vector) to be estimated, ϕ is a scale parameter, and v denotes a variance function to indicate mean-variance relation. The working covariance matrix of \mathbf{Y}_i , \mathbf{V}_i is assumed to have the form $\phi \mathbf{A}_i^{\frac{1}{2}} \mathbf{R}_i(\boldsymbol{\alpha}) \mathbf{A}_i^{\frac{1}{2}}$, in which $\mathbf{A}_i = \text{diag}(v_{it})$ and $\mathbf{R}_i(\boldsymbol{\alpha})$ is the working correlation matrix parameterized by $\boldsymbol{\alpha}$, an association parameter (q -vector).

The GEE method identifies the estimator $\hat{\boldsymbol{\beta}}$ of the regression parameter $\boldsymbol{\beta}$ as the solution to Equation 1, substituting ϕ with a $K^{\frac{1}{2}}$ -consistent estimator $\hat{\phi}(Y, \boldsymbol{\beta})$ after replacing $\boldsymbol{\alpha}$ with a $K^{\frac{1}{2}}$ -consistent estimator $\hat{\boldsymbol{\alpha}}(\mathbf{Y}, \boldsymbol{\beta}, \phi)$.

$$U(\boldsymbol{\beta}) \equiv \sum_{i=1}^K \mathbf{D}_i^\top \mathbf{V}_i^{-1} \mathbf{S}_i = \mathbf{0}, \quad (1)$$

where \mathbf{D}_i is an $n \times p$ matrix defined by $\mathbf{D}_i = \partial \boldsymbol{\mu}_i / \partial \boldsymbol{\beta}$, $\mathbf{V}_i = \phi \mathbf{A}_i^{\frac{1}{2}} \mathbf{R}_i(\boldsymbol{\alpha}) \mathbf{A}_i^{\frac{1}{2}}$, $\mathbf{S}_i = \mathbf{Y}_i - \boldsymbol{\mu}_i$, and $\boldsymbol{\mu}_i = (\mu_{i1}, \dots, \mu_{in})^\top$.

A covariance estimator \mathbf{V}_r of $\hat{\boldsymbol{\beta}}$ by the GEE method, which is referred to as the robust variance, is given by Equation 2:

$$\mathbf{V}_r = \left(\sum_{i=1}^K \mathbf{D}_i^\top \mathbf{V}_i^{-1} \mathbf{D}_i \right)^{-1} \left(\sum_{i=1}^K \mathbf{D}_i^\top \mathbf{V}_i^{-1} \mathbf{S}_i \mathbf{S}_i^\top \mathbf{V}_i^{-1} \mathbf{D}_i \right) \left(\sum_{i=1}^K \mathbf{D}_i^\top \mathbf{V}_i^{-1} \mathbf{D}_i \right)^{-1}. \quad (2)$$

3. Criteria for selecting a working correlation structure

3.1. Quasi-likelihood under the independence model criterion

AIC is a well-known criterion for likelihood-based model selection. However, we cannot apply a criterion such as AIC to the GEE approach, since the GEE is not likelihood-based. Pan (2001) proposes a criterion based on quasi-likelihood, named QIC, to select the proper mean model or the working correlation structure.

The quasi-likelihood function on cluster i and observation t evaluated at the regression parameters $\boldsymbol{\beta}$ is given by $Q(\boldsymbol{\beta}, \phi; Y_{it}, \mathbf{x}_{it}) = Q_{it}/\phi$, where Q_{it} is listed for commonly used distributions in Table 1.

Under the assumptions that the clusters and observations are independent, QIC can be expressed as

$$\text{QIC}(\mathbf{R}) = -2 \sum_{i=1}^K \sum_{t=1}^n Q(\boldsymbol{\beta}, \phi; Y_{it}, \mathbf{x}_{it}) + 2 \text{tr} \{ \boldsymbol{\Omega} \mathbf{V}_r(\mathbf{R}) \}, \quad (3)$$

where tr refers to the sum of the diagonal elements of the matrix and $\boldsymbol{\Omega} = \sum_{i=1}^K \mathbf{D}_i^\top \mathbf{A}_i^{-1} \mathbf{D}_i$.

3.2. Rotnitzky-Jewell's criterion

Rotnitzky and Jewell (1990) propose the test statistics to support the hypothesis that the vector of regression coefficients equals a given $\boldsymbol{\beta}$. In the theorem pertaining to the test

Distribution	Canonical link function	Variance function	Q_{it}
Normal	μ_{it}	1	$-\frac{1}{2}(y_{it} - \mu_{it})^2$
Binomial	$\ln \{ \mu_{it} / (1 - \mu_{it}) \}$	$\mu_{it}(1 - \mu_{it})$	$y_{it} \ln \{ \mu_{it} / (1 - \mu_{it}) \} + \ln(1 - \mu_{it})$
Poisson	$\ln \mu_{it}$	μ_{it}	$y_{it} \log \mu_{it} - \mu_{it}$
Gamma	$1/\mu_{it}$	μ_{it}^2	$-y_{it}/\mu_{it} - \log \mu_{it}$
Inverse Gaussian	$1/\mu_{it}^2$	μ_{it}^3	$-y_{it}/(2\mu_{it}^2) + 1/\mu_{it}$

Table 1: Canonical link function, variance function, and quasi-likelihood for commonly used exponential family distributions.

statistics, Ψ_0 , Ψ_1 , and Ψ were, respectively, defined as follows:

$$\begin{aligned}\Psi_0 &= \frac{1}{K} \sum_{i=1}^K \mathbf{D}_i^\top \mathbf{V}_i^{-1} \mathbf{S}_i \mathbf{S}_i^\top \mathbf{V}_i^{-1} \mathbf{D}_i, \\ \Psi_1 &= \frac{1}{K} \sum_{i=1}^K \mathbf{D}_i^\top \mathbf{V}_i^{-1} \mathbf{D}_i, \\ \Psi &= \Psi_0^{-1} \Psi_1.\end{aligned}$$

When the working correlation structure is correctly specified, Ψ should be close to an identity matrix. [Hin *et al.* \(2007\)](#) describe the Rotnitzky-Jewell's criterion (RJC) to select the working correlation structure as

$$\text{RJC}(\mathbf{R}) = \left[\{1 - \text{tr}(\Psi)/p\}^2 + \{1 - \text{tr}(\Psi^2)/p\}^2 \right]^{\frac{1}{2}}.$$

3.3. Correlation information criterion

[Hin and Wang \(2009\)](#) propose CIC in Equation 4 as a modification of QIC to improve its performance.

$$\text{CIC}(\mathbf{R}) = \text{tr} \{ \Omega \mathbf{V}_r(\mathbf{R}) \}. \quad (4)$$

CIC is constructed using only the second term that represents the penalty of QIC in Equation 3. The first term in QIC denotes the sum of quasi-likelihoods for all observations under the assumption that the subjects and time points are independent. It makes sense to ignore the first term when comparing different working correlation structures, since the term mostly does not depend on the specified \mathbf{R} .

3.4. Gosho's criterion

[Gosho *et al.* \(2011\)](#) propose to select the correlation structure that minimizes $\text{DEW}(\mathbf{R})$ as defined by the working correlation structure represented by Equation 5:

$$\text{DEW}(\mathbf{R}) = \text{tr} \left[\left\{ \left(\frac{1}{K} \sum_{i=1}^K \mathbf{S}_i \mathbf{S}_i^\top \right) \left(\frac{1}{K} \sum_{i=1}^K \mathbf{V}_i \right)^{-1} - \mathbf{I} \right\}^2 \right], \quad (5)$$

where \mathbf{I} is the identity matrix. In Equation 5, $\text{DEW}(\mathbf{R})$ is the criterion that directly measures the discrepancy between the covariance matrix estimator and the specified working covariance matrix.

4. Program description and usage

The `CriteriaWorkCorr` macro is included with this article in the file `CriteriaWorkCorr.sas`. The arguments taken by the macro are summarized in Table 2. As shown in Table 2, the first five arguments, i.e., the dataset name (`INDS`), the name of the variable identifying each cluster (`ID`), the name of the visit variable in each cluster (`VISIT`), the name of the outcome variable

Argument	Description	Note
INDS	Name of the SAS dataset.	Specify the input dataset that contains the cluster/subject, visit, outcome variables, and any additional covariates. The data structure is illustrated in Table 3.
ID	Name of the variable identifying each cluster (subject).	The types character and numerical are available (SAS Institute Inc. 2003b).
VISIT	Name of the visit variable (within a cluster).	The types character and numerical are available (SAS Institute Inc. 2003b).
OUTCOME	Name of the outcome variable.	The outcome can be continuous, binary, or count.
DIST	Name of the distribution of the outcome variable.	Specify one of <code>binomial</code> , <code>gamma</code> , <code>igaussian</code> , <code>normal</code> , or <code>poisson</code> .
COVCONT	List of continuous explanatory variables.	Do not separate variable names by comma.
COVNOMI	List of nominal explanatory variables.	Do not separate variable names by comma. Dummy variables with two levels (0 or 1) are automatically generated.
SCALEPAR	Name of the estimate method of ϕ .	Specify one of <code>fixed</code> , <code>pearson</code> , or <code>deviance</code> . <code>fixed</code> , fixed of 1; <code>pearson</code> , based on Pearson residuals; <code>deviance</code> , based on deviance residuals.

Table 2: Arguments for implementing `CriteriaWorkCorr` macro.

(`OUTCOME`), and the name of the distribution of the outcome variable (`DIST`), are essential for implementing the macro. In addition, `COVCONT` and `COVNOMI` are the names of the continuous and nominal explanatory variable lists, respectively. A user should generally at least specify either `COVCONT` or `COVNOMI` to implement the macro. `SCALEPAR` requests the estimate method of the scale parameter ϕ .

The `CriteriaWorkCorr` macro consists of three nested macros, `DataHandling`, `CalCri`, and `ResultDs`. The `DataHandling` macro generates the analysis dataset. When one or more nominal explanatory variables are specified as `COVNOMI`, corresponding dummy variables are automatically generated by the macro. The `CalCri` macro provides the regression parameter estimates, the robust standard errors, the 95% confidence intervals, and the p values of the z test in the case of applying the GEE method. It also provides the values of the criteria (QIC, RJC, CIC, and DEW) for selecting the working correlation structure that are given in Section 3. The `ResultDs` macro creates the combined output of the results derived from the `CalCri` macro for each working correlation structure.

A user can call the `CriteriaWorkCorr` macro by inputting the arguments listed in Table 2. Before a user implements the macro, he/she needs to prepare a SAS dataset for analysis. Part of an example dataset (Wheeze data) provided by Hardin and Hilbe (2003) is shown in Table 3. These data study the effect of air pollution on the health of 16 children.

case	t	wheeze	kingston	age	smoke
1	1	1	0	9	0
1	2	1	0	10	0
1	3	1	0	11	0
1	4	0	0	12	0
2	1	1	1	9	1
2	2	1	1	10	2
2	3	0	1	11	2
2	4	0	1	12	2

Table 3: Example dataset for two subjects from the Wheeze data in [Hardin and Hilbe \(2003\)](#).

The **Wheeze** data include the case number, **case**; a within-subject observation identifier (time point), **t**; a binary indicator for whether the subject wheezes, **wheeze**; a binary indicator for whether the observation is in Kingston, **kingston**; the age of the subject in years, **age**; and a measure of the smoking habits of the subject's mother, **smoke**. In this case, **wheeze** is an outcome variable and **kingston**, **age**, and **smoke** are explanatory variables. The nominal explanatory variables **kingston** and **smoke**, which take the value zero or one and the value zero, one, or two, respectively, are turned into dummy variables. The dummy variable **kingston1**, derived from **kingston**, takes a value of zero if **kingston** is zero and a value of one otherwise. In addition, the two dummy variables **smoke1** and **smoke2**, derived from **smoke**, take a value of one if **smoke** is equal to one and zero otherwise and a value of one if **smoke** is equal to two and zero otherwise, respectively. The minimum value of a nominal explanatory variable becomes a reference level for the corresponding dummy variables in the macro.

The macro creates an output table that includes the regression parameter estimates, the robust standard errors, the 95% confidence intervals, and the p values of the z test in the case where the GEE method is applied and three working correlation structures (independent, exchangeable, and AR(1) structures) are specified using the **GENMOD** procedure in SAS. In addition, the output table contains the values of the criteria (QIC, RJC, CIC, and DEW) for selecting the working correlation structure for each such structure using the **IML** procedure in SAS.

There are some limitations of the **CriteriaWorkCorr** macro. The macro can be applied to incomplete longitudinal data only when the incompleteness follows a monotone missing pattern; that is, a subject missing in one follow-up must also fail to participate in subsequent follow-ups. Another limitation is that the macro assumes the link function is the canonical link as listed in Table 1.

5. Example

In this section, the implementation of the **CriteriaWorkCorr** macro is demonstrated using the **Wheeze** data. As mentioned earlier, the outcome variable is a binary indicator for whether or not the subject wheezed, and is measured consistently four times yearly at ages 9, 10, 11, and 12. We fitted the following logistic model to the data:

$$\text{logit}\{E(Y_{it})\} = \beta_0 + \beta_1 \text{age}_{it} + \beta_2 \text{kingston1}_i + \beta_3 \text{smoke1}_{it} + \beta_4 \text{smoke2}_{it},$$

WorkingCorr	Parm	Level	Estimate	RobustSE	LowerCL	UpperCL	ProbZ	QIC	RJC	CIC	DEW
Independent	Intercept		1.9838	3.0408	-3.9761	7.9437	0.5142	85.5221	0.9356	5.9122	0.8762
	age		-0.2459	0.2813	-0.7972	0.3055	0.3821				
	kingston1	1	0.2105	0.6810	-1.1242	1.5452	0.7572				
	smoke1	1	-0.9709	0.6200	-2.1860	0.2442	0.1173				
Exchangeable	smoke2	1	0.2003	0.6221	-1.0189	1.4196	0.7474				
	Intercept		1.9434	2.9708	-3.8793	7.7661	0.5130	85.0896	0.4169	5.6732	0.7476
	age		-0.2444	0.2736	-0.7806	0.2918	0.3716				
	kingston1	1	0.1605	0.6741	-1.1607	1.4817	0.8118				
AR(1)	smoke1	1	-0.8517	0.4928	-1.8175	0.1142	0.0839				
	smoke2	1	0.2163	0.6386	-1.0353	1.4680	0.7348				
	Intercept		1.7133	2.8224	-3.8184	7.2451	0.5438	84.8718	0.1530	5.3881	0.3250
	age		-0.2420	0.2622	-0.7559	0.2719	0.3561				
	kingston1	1	0.3400	0.6466	-0.9273	1.6073	0.5990				
	smoke1	1	-0.6091	0.4577	-1.5062	0.2879	0.1832				
	smoke2	1	0.4130	0.6731	-0.9062	1.7322	0.5395				

Table 4: Parameter estimates, standard errors, 95% confidence intervals, p values, and criteria (QIC, RJC, CIC, and DEW) values for each working correlation structure in `Wheeze` data.

where Y_{it} is the binary indicator for whether or not subject i wheezed at time t ; $\mathbf{age}_{it} \in \{9, 10, 11, 12\}$ denotes the child's age; $\mathbf{kingston1}_i \in \{0, 1\}$ indicates whether the child is a resident of Portage or Kingston, respectively; and $\mathbf{smoke1}_{it}$ and $\mathbf{smoke2}_{it}$ are dummy variables for the smoking habits of the child's mother, that take a value of one if \mathbf{smoke} is equal to one and zero otherwise and a value of one if \mathbf{smoke} is equal to two and zero otherwise, respectively. Three structures—independent, exchangeable, and AR(1)—are adopted as candidates for the working correlation structure.

To implement the `CriteriaWorkCorr` macro, a user inputs values for the various arguments as shown in the code below (also see Table 2) and invokes the macro.

```
\%CriteriaWorkCorr (INDS = wheeze, ID = case, VISIT = t, OUTCOME = wheeze,
  DIST = binomial, COVCONT = age, COVNOMI = kingston smoke,
  SCALEPAR = fixed);
```

Table 4 shows the output results of the macro. In Table 4, “WorkingCorr” is the specified working correlation structure, and “Parm” and “Level” are the names of the parameters and the parameterized level, respectively. “Estimate” and “RobustSE” refer to the parameter estimates and the robust standard errors, respectively. “LowerCL” and “UpperCL” are the lower and upper limits of the 95% confidence intervals, respectively. “ProbZ” gives the p values of the z test. “QIC,” “RJC,” “CIC,” and “DEW” refer to the criteria values for selecting the working correlation structure given in Section 3.

6. Conclusion

In this paper, we provide the `CriteriaWorkCorr` macro to calculate the values of the criteria (QIC, RJC, CIC, and DEW) for selecting a working correlation structure at the time of applying the GEE method to analyze clustered and longitudinal data.

Acknowledgments

The author would like to thank Emeritus Professor Isao Yoshimura of Tokyo University of Science for his valuable comments. We are grateful to the Editors and two anonymous referees for their helpful comments and suggestions.

References

- Fitzmaurice GM (1995). “A Caveat Concerning Independence Estimating Equations with Multivariate Binary Data.” *Biometrics*, **51**(1), 309–317.
- Gosho M, Hamada C, Yoshimura I (2011). “Criterion for the Selection of a Working Correlation Structure in the Generalized Estimating Equation Approach for Longitudinal Balanced Data.” *Communications in Statistics – Theory and Methods*, **40**(21), 3839–3856.
- Hardin J, Hilbe J (2003). *Generalized Estimating Equations*. Chapman and Hall, London.

- Hin LY, Carey VJ, Wang YG (2007). “Criteria for Working-Correlation-Structure Selection in GEE: Assessment via Simulation.” *The American Statistician*, **61**(4), 360–364.
- Hin LY, Wang YG (2009). “Working-Correlation-Structure Identification in Generalized Estimating Equations.” *Statistics in Medicine*, **28**(4), 642–658.
- Mancl LA, Leroux BG (1996). “Efficiency of Regression Estimates for Clustered Data.” *Biometrics*, **52**(2), 500–511.
- Pan W (2001). “Akaike’s Information Criterion in Generalized Estimating Equations.” *Biometrics*, **57**(1), 120–125.
- Rotnitzky A, Jewell NP (1990). “Hypothesis Testing of Regression Parameters in Semiparametric Generalized Linear Models for Cluster Correlated Data.” *Biometrika*, **77**(3), 485–497.
- SAS Institute Inc (2003a). *The SAS System, Version 9.1*. SAS Institute Inc., Cary, NC. URL <http://www.sas.com/>.
- SAS Institute Inc (2003b). *SAS/BASE Software, Version 9.1*. SAS Institute Inc., Cary, NC. URL <http://www.sas.com/>.
- SAS Institute Inc (2003c). *SAS/IML 9.1 User’s Guide*. SAS Institute Inc., Cary, NC. URL <http://www.sas.com/>.
- SAS Institute Inc (2003d). *SAS/STAT Software, Version 9.1*. SAS Institute Inc., Cary, NC. URL <http://www.sas.com/>.
- Sutradhar BC, Das K (2000). “On the Accuracy of Efficiency of Estimating Equation Approach.” *Biometrics*, **56**(2), 622–625.
- Vens M, Ziegler A (2012). “Generalized Estimating Equations and Regression Diagnostics for Longitudinal Controlled Clinical Trials: A Case Study.” *Computational Statistics & Data Analysis*, **56**(5), 1232–1242.
- Wang YG, Carey V (2003). “Working Correlation Structure Misspecification, Estimation and Covariate Design: Implications for Generalised Estimating Equations Performance.” *Biometrika*, **90**(1), 29–41.
- Ziegler A, Vens M (2010). “Generalized Estimating Equations: Notes on the Choice of the Working Correlation Matrix.” *Methods of Information in Medicine*, **49**(5), 421–425.

Affiliation:

Masahiko Goshō

Unit of Biostatistics, Advanced Medical Research Center

Aichi Medical University
1-1, Yazakokarimata, Nagakute
Aichi, 480-1195, Japan
E-mail: mgosho@aichi-med-u.ac.jp