

Tests for Gaussian repeated measures with missing data in small samples

Diane J. Catellier^{*,†} and Keith E. Muller

Department of Biostatistics CB#7400, University of North Carolina, Chapel Hill, North Carolina 27599-7400, U.S.A.

SUMMARY

For small samples of Gaussian repeated measures with missing data, Barton and Cramer recommended using the EM algorithm for estimation and reducing the degrees of freedom for an analogue of Rao's F approximation to Wilks' test. Computer simulations led to the conclusion that the modified test was slightly conservative for total sample size of $N = 40$. Here we consider additional methods and smaller sample sizes, $N \in \{12, 24\}$. We describe analogues of the Pillai–Bartlett trace, Hotelling–Lawley trace and Geisser–Greenhouse corrected univariate tests which allow for missing data. Eleven sample size adjustments were examined which replace N by some function of the numbers of non-missing pairs of responses in computing error degrees of freedom. Overall, simulation results allowed concluding that an adjusted test can always control test size at or below the nominal rate, even with as few as 12 observations and up to 10 per cent missing data. The choice of method varies with the test statistic. Replacing N by the mean number of non-missing responses per variable works best for the Geisser–Greenhouse test. The Pillai–Bartlett test requires the stronger adjustment of replacing N by the harmonic mean number of non-missing pairs of responses. For Wilks' and Hotelling–Lawley, an even more aggressive adjustment based on the minimum number of non-missing pairs must be used. Copyright © 2000 John Wiley & Sons, Ltd.

1. INTRODUCTION

1.1. Motivation

Repeated measurements can be distinct variables, or a single variable measured at several points in time, with the spacing consistent across subjects. For ease of presentation, call the experimental unit a 'subject' and the metameter on which the measurements are indexed 'time'. Traditional linear models are particularly useful if the responses are at least approximately Gaussian and can

* Correspondence to: Diane J. Catellier, Department of Biostatistics CB#7400, University of North Carolina, Chapel Hill, North Carolina 27599-7400, U.S.A.

† E-mail: diane_catellier@UNC.EDU

Contract/grant sponsor: USFHHS; contract/grant number: 90CA1467

Contract/grant sponsor: IPRIC; contract/grant number: 5-326-12

Contract/grant sponsor: NIH; contract/grant numbers: PO1-CA47982-04, RO1-CA67183-01A1, RO1-CA72875, RO1-CA60193-04, MO1-RR000-46-33, NO1-ES-35356, MH33127

be explained by some linear function of predictors. When each subject is observed at the same p times and no missing observations, closed-form maximum likelihood (ML) estimates of the model parameters are often available. Often, especially in clinical trials, the response is not observed at all time points for every subject.

A large amount of research has been directed at estimation for linear repeated measures models with missing data. These appear to work well in both large and small samples. In contrast, much less effort has been directed towards methods for inference. Various asymptotic test statistics work well in large samples. However, in small samples the available methods for inference may work very poorly. In particular, the methods produce inflated type I error rates in small samples [1, 2].

We seek to develop hypothesis tests for Gaussian repeated measures with missing data, accurate in small samples. In doing so, we restrict attention to a particular range of studies. Using the terminology of Rubin [3], missing responses are said to be missing at random (MAR) if missingness of a particular response does not depend on its unobserved value, but can depend on the covariates or any of the observed responses. In contrast, describe data as missing completely at random (MCAR) if missingness of a particular response does not depend on its unobserved value or any of the observed responses or covariates. Likelihood-based estimation methods assume that the data are MAR. In the discussion in this paper we assume that the data are MAR, although the simulations all assume the more stringent property of MCAR.

1.2. General strategies for the analysis of repeated measures designs

Models in which the expected value of the response vector equals a linear function of the parameters have traditionally been described as (general) linear models. Most often, one of three strategies is used for linear models with repeated measures: the multivariate analysis of variance (MANOVA) approach; the univariate approach to repeated measures, or mixed model analysis. All three models account for the dependencies among the repeated measures, but differ in the special form assumed for the covariance matrix within subjects, Σ .

The mixed model has long been used for the analysis of continuous data, especially for missing and mistimed data. By virtue of modelling the subject as a random component, the mixed model can encompass a broad range of covariance structures. In this paper, we will compare existing inference methods for all three approaches to new ones.

In contrast to likelihood-based approaches, one can use the quasi-likelihood approach of Liang and Zeger [4] by solving the generalized estimating equations (GEE) to obtain estimates of the regression parameters. Park [5] compared the GEE approach to the ML approach for multivariate normal outcomes. He showed that with no missing data and an unstructured covariance matrix, the GEE and ML score equations are equivalent and lead to the same estimates of expected value and covariance parameters. With missing observations, however, the equivalence fails. For data missing completely at random, MCAR [3], the GEE solution produces consistent estimators. Three weaknesses, however, make GEE less desirable than ML estimation in the missing data setting. First, the GEE estimate of the working covariance matrix may not always be positive definite, while the ML estimate from the EM algorithm [6] is guaranteed to be positive definite [7]. Second, simulation studies allow concluding that ML estimators perform better in small samples than GEE estimators [5, 8–10]. In particular, ML estimates tend to have less bias, smaller mean squared errors and associated tests have more accurate test size than corresponding GEE estimates and tests. Third, under misspecification of the covariance, the GEE estimators will only be consistent provided that the missing observations are MCAR. ML procedures give unbiased estimates under

the weaker MAR assumption. The limitations of the GEE procedure in small samples, combined with the focus on Gaussian data, makes ML estimation more attractive and therefore the focus of this paper.

2. KNOWN METHODS FOR ESTIMATION AND INFERENCE

2.1. Complete data

The repeated measures model can be represented as a special case of the general linear multivariate model (GLMM); see Davidson [11] or O'Brien and Kaiser [12] for general introductions. Muller *et al.* [13] briefly described the mathematical formulation of estimation and hypothesis testing and power analysis for the multivariate and univariate approaches to repeated measures ANOVA, assuming complete data. They described the assumptions behind the methods as well as the most widely used tests for both approaches. We essentially follow their notation. The key formulae for hypothesis testing are summarized here. For additional detail, see Muller *et al.* [13], as well as references they provide.

Suppose that the responses for subject i ($i \in \{1, \dots, N\}$) are measured at p times (the same for all subjects). Specify the GLMM as

$$\mathbf{Y} = \mathbf{X}\mathbf{B} + \mathbf{E} \quad (1)$$

with \mathbf{Y} the observed matrix of random observations ($N \times p$), \mathbf{X} the design matrix ($N \times q$, fixed and known, conditional upon having chosen the subjects) and \mathbf{B} the fixed and unknown parameters ($q \times p$). Indicate the i th row of \mathbf{X} as $\mathbf{X}_i = \text{row}_i(\mathbf{X})$. Assuming that $\text{row}_i(\mathbf{E})$ is $N_p(\mathbf{0}, \mathbf{\Sigma})$, we can test the general linear hypothesis (GLH)

$$H_0: \mathbf{\Theta} = \mathbf{C}\mathbf{B}\mathbf{U} = \mathbf{\Theta}_0 \text{ versus } H_1: \mathbf{\Theta} \neq \mathbf{\Theta}_0 \quad (2)$$

Each row of \mathbf{C} ($a \times q$) defines a between-subject contrast and each column of \mathbf{U} ($p \times b$) defines a within-subjects contrast. Define $v_E = N - \text{rank}(\mathbf{X})$ and $\mathbf{\Sigma}_* = \mathbf{U}'\mathbf{\Sigma}\mathbf{U}$. All tests of H_0 are based on

$$\hat{\mathbf{B}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y} \quad (3)$$

$$\hat{\mathbf{\Theta}} = \mathbf{C}\hat{\mathbf{B}}\mathbf{U} \quad (4)$$

$$\hat{\mathbf{H}} = (\hat{\mathbf{\Theta}} - \mathbf{\Theta}_0)'[\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}']^{-1}(\hat{\mathbf{\Theta}} - \mathbf{\Theta}_0) \quad (5)$$

and

$$\hat{\mathbf{E}} = \mathbf{U}'(\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}})'(\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}})\mathbf{U} \quad (6)$$

$\hat{\mathbf{H}}$ and $\hat{\mathbf{E}}$ have Wishart distributions $W(a, \mathbf{\Sigma}_*, \mathbf{\Omega})$ and $W(v_E, \mathbf{\Sigma}_*)$ respectively, with

$$\mathbf{\Omega} = (\mathbf{\Theta} - \mathbf{\Theta}_0)'[\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}']^{-1}(\mathbf{\Theta} - \mathbf{\Theta}_0)\mathbf{\Sigma}_*^{-1} \quad (7)$$

Let $s = \min(a, b)$ indicate the rank of $\hat{\mathbf{H}}$ and hence $\hat{\mathbf{H}}\hat{\mathbf{E}}^{-1}$.

The common multivariate tests may be constructed using the eigenvalues (l_1, \dots, l_s) of $\hat{\mathbf{H}}\hat{\mathbf{E}}^{-1}$. Specifically, for Wilks' lambda $W = \prod_{i=1}^s (1 + l_i)^{-1}$, for Pillai-Bartlett trace $V = \sum_{i=1}^s l_i(1 + l_i)^{-1}$, for Hotelling-Lawley trace $U = \sum_{i=1}^s l_i$, and for Roy's largest root $R = \max(l_i)$. When $s > 1$ (or $s > 2$ for Wilks'), closed form expressions for the distributions of these test statistics are not available and approximations are used [13].

In general, no uniformly most powerful test exists. Hence the optimal choice depends on the alternative hypothesis. Describe a situation with only one non-zero eigenvalue of $\mathbf{\Omega}$ as having concentrated non-centrality (also referred to as the ‘linear’ case). Describe situations with more than one non-zero eigenvalue of $\mathbf{\Omega}$ as having diffuse non-centrality. Roy’s test has the highest power of all statistics for concentrated non-centrality. The other three statistics typically have much higher power than Roy’s test for diffuse cases, and power lower, but not much lower, for concentrated cases. See Olson [14–16], Anderson (Reference [17], pp. 330–333), and Muller *et al.* [13] for detailed discussions. Roy’s test has much greater test size inflation under violation of the assumption of homogeneity of covariance, when compared to other multivariate tests. The poor robustness and low power for diffuse cases led us to not consider Roy’s test any further.

The test statistics W , V and U can be accurately approximated by an F random variable. Rao’s [13, 18] approximation for W works well, even in very small samples. Although widely used in current statistical packages, Pillai’s [19] F approximation for V may be very conservative in small samples. Muller [20] developed an F approximation for V that provides substantially better accuracy. Hence Muller’s approximation will be used, with $v_1(V) = Kab$, $v_2(V) = Ks(v_E + s - b)$, with

$$K = \frac{1}{s(v_E + a)} \left[\frac{s(v_E + s - b)(v_E + a + 2)(v_E + a - 1)}{v_E(v_E + a - b)} \right]. \quad (8)$$

McKeon [21] provided a slightly better F approximation than the more widely used Pillai–Sampson approximation for the Hotelling–Lawley statistic. Write the McKeon approximation

$$F = \frac{(U/h)/(ab)}{1/v_2(U)} \quad (9)$$

with $v_1(U) = ab$, $v_2(U) = (4 + ab + 2)g'$

$$g' = \frac{v_E^2 - v_E(2b + 3) + b(b + 3)}{v_E(a + b + 1) - (a + 2b + b^2 - 1)} \quad (10)$$

and

$$h = \frac{v_2(U) - 2}{v_E - b - 1} \quad (11)$$

Computations for the ‘univariate’ approach to repeated measures follow easily from the multivariate set-up. Let $\hat{\lambda}_k$ indicate the k th eigenvalue of $\hat{\mathbf{\Sigma}}_* = \hat{\mathbf{E}}/v_E$, and

$$\hat{\varepsilon} = \frac{(\sum_{k=1}^b \hat{\lambda}_k)^2}{b \sum_{k=1}^b \hat{\lambda}_k^2} \quad (12)$$

the maximum likelihood estimate of ε , which indicates the deviation of $\mathbf{\Sigma}_*$ from sphericity. Also define the Geisser–Greenhouse corrected F statistic as

$$F_{GG} = \frac{\text{tr}(\hat{\mathbf{H}})/(ab)}{\text{tr}(\hat{\mathbf{E}})/(bv_E \hat{\varepsilon})} \quad (13)$$

with degrees of freedom $ab \hat{\varepsilon}$ and $bv_E \hat{\varepsilon}$. See Muller and Barton [22, 23], or Muller *et al.* [13] for additional details of the computation and approximation of the distribution of the various tests

for the ‘univariate’ approach to repeated measures. Muller and Barton [22] suggested that the Geisser–Greenhouse (GG) test provides the best compromise in controlling type I error rate with excellent power. Hence only the GG test statistic will be examined here.

2.2. Data missing at random

For the GLMM, both ML and REML estimation methods have been extensively investigated for MAR data. For a general review, see Little and Rubin (Reference [24], chapters 7–10). For some special case patterns of missing data, such as monotone missing data, the likelihood equations have a closed form solution [25]. Data are described as monotone missing if, for each sampling unit (subject), all observations after the first (in time) missing observation are always missing. As an example, assume no data are missing until a subject’s death, but all data after death are missing. For arbitrary missing data patterns, the solution must be obtained iteratively. The computational efficiency and simplicity of the EM algorithm [6, 26, 27] makes it an attractive choice for ML estimation in the setting of interest. Barton and Cramer [1] achieved accurate test size in simulations of methods of the sort we propose here. Barton and Cramer [1] studied the GLMM with as few as 40 sampling units. The best algorithm for more general (mixed) models is not as obvious, see for example, Mensah *et al.* [28] or Callahan and Harville [29].

Except in special cases, no known method exists for providing accurate and efficient inference in small multivariate normal samples with missing data. Hypothesis tests constructed from complete observations only, while accurate in small samples, are inefficient. A number of approximate methods have been proposed for the problem of testing equality of means for a bivariate normal sample with data missing on one variable [30–32].

Barton and Cramer [1] suggested an appealing technique for testing the general linear hypothesis in a GLMM with an arbitrary pattern of missing data. The approach involves using the EM algorithm for ML estimation, and modifying Rao’s F approximation to Wilks’ test, F_W , with adjusted error degrees of freedom. Let $N_{jj'}$ indicate the number of observations for which both Y_{ij} and $Y_{ij'}$, for $i \in \{1, \dots, N\}$, have non-missing values. Note that N_{jj} equals the number of cases observed for the j th response. All adjustments considered by Barton and Cramer [1], and in this paper, involve replacing N by N_* in $v_E = N - \text{rank}(\mathbf{X})$. In all cases N_* equals a function only of $\{N_{jj'}\}$. For samples of size 40 and up to 20 per cent missing data, test statistics with degrees of freedom based on the naive choice $N_* = N$ produced inflated type I error rates ranging from 0.10 to 0.23 assuming a nominal rate of 0.05. In contrast, choosing N_* as the number of complete cases gave very conservative rates (0.004–0.014). The best choice was the average number of non-missing pairs of responses. An analogue of Wilks’ test based on this adjustment produced acceptable test sizes across all simulated conditions.

The mixed model is often used for multivariate data with some missing observations. Let \mathbf{y}_i be an $(N_i \times 1)$ vector of measurements for the i th subject, and $N_+ = \sum_{i=1}^N N_i$. In the mixed model, $\mathbf{y}_+ = [\mathbf{y}'_1, \dots, \mathbf{y}'_N]'$ is modelled as

$$\mathbf{y}_+ = \mathbf{X}_+ \boldsymbol{\beta} + \mathbf{Z} \mathbf{b} + \mathbf{e}_+ \quad (14)$$

with \mathbf{X}_+ and \mathbf{Z} the known design matrices for the fixed and random effects respectively, and the \mathbf{b} the vector of unknown random effects. The key assumptions for inference are that \mathbf{b} and \mathbf{e}_+ are independent and multivariate Gaussian. Define $\text{vec}(\mathbf{M})$ as the vector created by stacking the columns of \mathbf{M} . Also let $\mathbf{A} \otimes \mathbf{B} = \{a_{ij} \mathbf{B}\}$ indicate the (left) Kronecker product. For the cases of interest, the mixed model may be written so that $\boldsymbol{\beta} = \text{vec}(\mathbf{B}')$, with \mathbf{B} from the GLMM, $\mathbf{Y} = \mathbf{X}\mathbf{B} + \mathbb{E}$.

For complete data $\mathbf{X}_+ = \mathbf{X} \otimes \mathbf{I}_p$. For missing data merely delete each row corresponding to a missing response. Let Σ_i , of dimension p_i , be the submatrix of Σ with rows and columns corresponding to data observed for subject i . Let $\Lambda = \text{Dg}(\Sigma_1, \dots, \Sigma_N)$ indicate the block-diagonal matrix created by placing Σ_i in the upper left diagonal, etc. Then $\mathbf{e}_+ = \mathcal{N}_{N_+}(\mathbf{0}, \Lambda)$. Likelihood-based estimation of β and Λ typically requires iterative methods. The software used here (SAS[®], PROC MIXED) employs a method of Lindstrom and Bates [33].

Exact methods are not available to test hypotheses about β . The approximate large-sample test statistics can be unreliable in small samples. Schluchter and Elashoff [2] examined the test size of various ML Wald-type statistics. Their small sample simulation results led them to suggest approximating a modified Wald statistic by an F distribution with denominator degrees of freedom based on the number of complete cases. A simple χ^2 approximation to the LR statistic [34] proves unreliable in small samples [35–37]. The version of PROC MIXED studied here used

$$F = \frac{\hat{\theta}'(\mathbf{X}_+^t \hat{\Lambda}^{-1} \mathbf{X}_+)^{-1} \hat{\theta}}{\text{rank}(\mathbf{C})} \quad (15)$$

with numerator degrees of freedom equal to $\text{rank}(\mathbf{C})$ [38]. Although several approximations are available for the denominator degrees of freedom [38], only the Satterthwaite approximation was considered in this paper.

3. NEW TESTS FOR DATA MISSING AT RANDOM

The success of the Barton and Cramer [1] strategy encouraged us to consider a number of generalizations. First, the approach will be applied to other test statistics. Second, some additional functions of the sample sizes merit consideration. Third, even smaller sample sizes will be studied. In all cases the EM algorithm will be used for estimation.

In addition to W, the U, V and GG tests may be applied to missing data. First, use the EM algorithm to compute maximum likelihood estimates of \mathbf{B} and Σ . Second, use the estimates to compute an analogue of $\hat{\mathbf{H}}\hat{\mathbf{E}}^{-1}$, as well as analogues of U, V and GG statistics. Third, compute F approximations, changing only the error degrees of freedom by replacing N with some form of N_* .

Overall, eleven forms for N_* will be examined for each of four test statistics. They are listed in Table I in rank order from smallest to largest, with the exception that N_{*3} can be either less than or greater than N_{*6} (and hence N_{*4} and N_{*5}). In all cases $N/N_* \rightarrow 1$. Consequently, in large samples (as $N \rightarrow \infty$, with fixed p, q and proportion missing) the choice of N_* has less and less effect. The form of the results of Rothenberg (which assume a sequence of local alternatives), as cited in Anderson (Reference [17], Section 8.6.5) support this position.

4. NUMERICAL EVALUATIONS

4.1. Methods

All simulations involved a small range of research designs. In all cases, the designs included: (i) one within-subject factor with $p \in \{3, 6\}$ levels (such as $p = 3$ repeated measures); (ii) one between-subject factor with $q = 4$ levels; (iii) $N \in \{12, 24\}$; (iv) 0, 5 or 10 per cent of the data missing. No subject's data were allowed to be completely missing. The procedure for producing

Table I. Sample size adjustments for error degrees of freedom.

Name	Function of $\{N_{jj'}\}$
N_{*1}	= number of complete cases
N_{*2}	= $\min(\{N_{jj'}\})$
N_{*3}	= $\min(\{N_{jj}\})$
N_{*4}	= harmonic mean($\{N_{jj'}\}$)
N_{*5}	= geometric mean($\{N_{jj'}\}$)
N_{*6}	= arithmetic mean($\{N_{jj'}\}$)
N_{*7}	= harmonic mean($\{N_{jj}\}$)
N_{*8}	= geometric mean($\{N_{jj}\}$)
N_{*9}	= arithmetic mean($\{N_{jj}\}$)
N_{*10}	= $\max(\{N_{jj}\})$
N_{*11}	= N

Table II. Test size for mixed model F (5000 replications, ± 0.006).

N	$\rho_{jj'}$	σ_j^2	% Missing	$p = 3$	$p = 6$
12	Low	=	0	0.126	0.58
12	Low	\neq	0	0.134	0.60
12	High	\neq	0	0.125	0.59
24	Low	\neq	0	0.069	0.16
24	High	\neq	0	0.074	0.16
12	Low	=	5	0.182	
12	Low	\neq	5	0.171	
12	High	\neq	5	0.172	
24	Low	\neq	5	0.080	0.21
24	High	\neq	5	0.081	0.23
12	Low	=	10	0.250	
12	Low	\neq	10	0.244	
12	High	\neq	10	0.262	
24	Low	\neq	10	0.086	0.303
24	High	\neq	10	0.095	0.322

missing data generated data that are MCAR. Other factors considered are the relative error variance of response variables (equal, unequal), and the error correlation structure (medium, high). See Tables I and II in Barton and Cramer [1] for details. In addition, a third level was added to the correlation structure factor, which allowed assessing the effect of very low correlation between the responses. The structure specified equal correlation ($\rho = 0.1$) for each pair of responses. The overall test for the presence of a trend (linear, quadratic or cubic) with respect to the between-subject factor for each of the response measures was of primary interest. Under the null, of course, $\mathbf{B} = \mathbf{0}$. For 5000 replications and assuming a true type I error rate of 0.05, the 95 per cent confidence bounds around the type I error rate estimates are approximately ± 0.006 .

Table III. Test size for GLMM F tests (0 per cent missing, 5000 replications, ± 0.006).

N	$\rho_{jj'}$	σ_j^2	F_W		F_U		F_V		F_{GG}	
			$p=3$	$p=6$	$p=3$	$p=6$	$p=3$	$p=6$	$p=3$	$p=6$
12	Low	=	0.050	0.046	0.048	0.051	0.044	0.046	0.027	0.013
12	Low	\neq	0.052	0.054	0.050	0.055	0.049	0.059	0.040	0.025
12	High	\neq	0.053	0.053	0.048	0.053	0.042	0.045	0.054	0.058
24	Low	\neq	0.049	0.051	0.050	0.050	0.048	0.049	0.041	0.039
24	High	\neq	0.053	0.049	0.051	0.048	0.051	0.049	0.049	0.053

Table IV. Adjusted degree of freedom test size for F_W (5 per cent, 10 per cent missing, 5000 replications, ± 0.006).

N	$\rho_{jj'}$	σ_j^2	% Missing	N_{*2}		N_{*4}		N_{*11}	
				$p=3$	$p=6$	$p=3$	$p=6$	$p=3$	$p=6$
12	Low	=	5	0.029		0.073		0.145	
12	Low	\neq	5	0.033		0.072		0.134	
12	High	\neq	5	0.032		0.074		0.141	
24	Low	\neq	5	0.030	0.017	0.049	0.067	0.087	0.142
24	High	\neq	5	0.031	0.022	0.053	0.067	0.089	0.146
12	Low	=	10	0.047		0.148		0.345	
12	Low	\neq	10	0.042		0.144		0.335	
12	High	\neq	10	0.051		0.152		0.354	
24	Low	\neq	10	0.020	0.078	0.051	0.171	0.133	0.379
24	High	\neq	10	0.025	0.094	0.062	0.199	0.158	0.411

4.2. Results

On average, higher levels of correlation within subjects were associated with modestly higher type I error rates. Since this pattern was consistent for each of the test statistics, only the results for the low and high correlation conditions will be presented.

The empirical type I error rates for the mixed model F statistic are given in Table II. The results indicate that the test has poor small sample properties, producing inflated type I error rates, even when none of the data were missing. For $N=24$, test sizes increased from slightly greater (0.07–0.10) to considerably greater (0.16–0.32) than the nominal $\alpha=0.05$ level as the number of repeated measures increased from 3 to 6.

For the conditions with no missing data, all four univariate and multivariate test statistics succeeded in controlling the type I error rate at or below the nominal rate (see Table III). This illustrates that the sample sizes, while small, are large enough that the approximate F tests are essentially unbiased for complete data. Hence any discrepancy from the desired test size may be attributed to the influence of missing responses, and not to any inaccuracy in test approximations for complete data.

Tables IV, V, VI and VII summarize the empirical test sizes for W, U, V, and GG for 5 per cent and 10 per cent missing data. All tables give results for tests based on $N_{*2} = \min\{N_{jj'}\}$ and $N_{*11} = N$ in order to define bounds on test size.

The EM algorithm failed roughly 90 per cent of the time for the condition with $p=6$, $N=12$, and even 5 per cent missing data. Estimates are well defined for complete data. The table cells for

Table V. Adjusted degree of freedom test size for F_U (5 per cent, 10 per cent missing, 5000 replications, ± 0.006).

N	$\rho_{jj'}$	σ_j^2	% missing	N_{*2}		N_{*4}		N_{*11}	
				$p=3$	$p=6$	$p=3$	$p=6$	$p=3$	$p=6$
12	Low	=	5	0.032		0.072		0.142	
12	Low	\neq	5	0.033		0.067		0.130	
12	High	\neq	5	0.034		0.074		0.134	
24	Low	\neq	5	0.028	0.021	0.048	0.069	0.084	0.144
24	High	\neq	5	0.031	0.025	0.052	0.068	0.088	0.148
12	Low	=	10	0.052		0.163		0.341	
12	Low	\neq	10	0.047		0.153		0.333	
12	High	\neq	10	0.055		0.169		0.352	
24	Low	\neq	10	0.021	0.093	0.052	0.184	0.135	0.379
24	High	\neq	10	0.029	0.086	0.061	0.217	0.158	0.417

Table VI. Adjusted degree of freedom test size for F_V (5 per cent, 10 per cent missing, 5000 replications, ± 0.006).

N	$\rho_{jj'}$	σ_j^2	% missing	N_{*2}		N_{*4}		N_{*11}	
				$p=3$	$p=6$	$p=3$	$p=6$	$p=3$	$p=6$
12	Low	=	5	0.023		0.052		0.102	
12	Low	\neq	5	0.021		0.051		0.099	
12	High	\neq	5	0.022		0.050		0.096	
24	Low	\neq	5	0.029	0.015	0.046	0.054	0.081	0.125
24	High	\neq	5	0.030	0.017	0.049	0.052	0.088	0.128
12	Low	=	10	0.011		0.057		0.206	
12	Low	\neq	10	0.010		0.051		0.215	
12	High	\neq	10	0.013		0.058		0.214	
24	Low	\neq	10	0.017	0.019	0.044	0.106	0.127	0.334
24	High	\neq	10	0.020	0.022	0.055	0.120	0.148	0.354

these conditions were left blank. Not surprisingly, the results indicate that the worst accuracy tends to occur with more repeated measures, fewer subjects, more missing data and more correlation within subjects.

From Tables IV and V, it is evident that the adjusted F_W and F_U tests based on N_{*11} give inflated test sizes, and those based on N_{*4} , while accurate for $N=24$ were inflated for $N=12$. On the other hand, tests based on N_{*2} controlled test size at or below the nominal rate under all simulated conditions, with the exception of the condition with $p=6$, $N=12$ and 10 per cent missing data, in which case test size was as high as 0.09.

Table VI contains test size for modified F_V tests. The test based on N_{*11} provided inflated type I error rates, and the N_{*2} -adjusted test was conservative. Test sizes for N_{*4} were extremely accurate, with the exception of the condition with $p=6$, $N=12$ and 10 per cent missing data where the type I error rates were approximately 0.1.

Table VII. Adjusted degree of freedom test size for F_{GG} (5 per cent, 10 per cent missing data, 5000 replications, ± 0.006).

N	$\rho_{ij'}$	σ_j^2	% missing	N_{*2}		N_{*9}		N_{*11}	
				$p=3$	$p=6$	$p=3$	$p=6$	$p=3$	$p=6$
12	Low	=	5	0.010		0.035		0.053	
12	Low	\neq	5	0.019		0.043		0.059	
12	High	\neq	5	0.026		0.050		0.066	
24	Low	\neq	5	0.030	0.010	0.047	0.037	0.061	0.057
24	High	\neq	5	0.049	0.018	0.049	0.041	0.061	0.053
12	Low	=	10	0.002		0.030		0.073	
12	Low	\neq	10	0.008		0.038		0.078	
12	High	\neq	10	0.009		0.040		0.075	
24	Low	\neq	10	0.018	0.004	0.051	0.044	0.091	0.095
24	High	\neq	10	0.017	0.009	0.048	0.043	0.076	0.070

Table VIII. Adjusted degree of freedom test size for multivariate F test with $s = \min(a, b) = 1$ (5 per cent, 10 per cent missing, 5000 replications, ± 0.006).

N	$\rho_{ij'}$	σ_j^2	% missing	$a = 1, b = 3$			$a = 3, b = 1$		
				N_{*2} $p=3$	N_{*4} $p=3$	N_{*11} $p=3$	N_{*2} $p=3$	N_{*4} $p=3$	N_{*11} $p=3$
12	Low	=	5	0.038	0.072	0.114	0.031	0.051	0.078
12	Low	\neq	5	0.034	0.066	0.105	0.029	0.048	0.071
12	High	\neq	5	0.039	0.074	0.110	0.027	0.042	0.065
24	Low	\neq	5	0.036	0.052	0.077	0.036	0.046	0.067
24	High	\neq	5	0.037	0.050	0.078	0.030	0.043	0.064
12	Low	=	10	0.058	0.132	0.253	0.018	0.043	0.096
12	Low	\neq	10	0.054	0.131	0.251	0.016	0.040	0.010
12	High	\neq	10	0.056	0.137	0.258	0.014	0.036	0.088
24	Low	\neq	10	0.028	0.052	0.102	0.026	0.043	0.089
24	High	\neq	10	0.028	0.048	0.104	0.018	0.036	0.073

The results given in Table VII suggest that N_{*9} is a reasonable adjustment function for the F_{GG} test. When $N = 12$, this test is slightly conservative, however this corresponds to the modest conservatism found when no data are missing (Table I in Muller and Barton [22]).

When $s = \min(a, b) = 1$, all of the multivariate test statistics are equivalent. This can occur if the rank of the \mathbf{C} contrast matrix is one ($a = 1$) or if the rank of U is one ($b = 1$). The empirical test sizes shown in Table VIII suggest that when $a = 1$, the best adjustment for the degrees of freedom is based on N_{*2} , while N_{*4} appears to work well when $b = 1$.

4.3. Example

Table IX contains data from a study that examined the effects of choline deprivation on plasma choline concentration over 35 days in healthy male subjects [39]. Subjects were given a standard

Table IX. Choline measurements over 5-week period in male subjects.

Treatment	Day					
	0	7	14	21	28	35
Control	9.93	12.29	9.30	9.51	10.84	9.24
	9.77	8.14	11.43	9.44	11.10	10.56
	12.56	10.90	11.19	12.31	9.95	
	10.15	10.32	8.86	9.23	8.56	12.78
	11.00	9.20	8.78	9.37	7.54	12.39
	10.46	8.72	8.13	8.14	11.76	9.74
Deficient	12.15	9.52	9.05	9.07	6.76	9.39
	12.88	9.66	7.71	7.29	6.37	10.61
	7.94	9.86	7.87	8.89	8.69	12.28
	9.42	12.82	7.17	8.18	8.30	12.61
	9.57	10.95	9.01	8.98	6.56	9.66
	11.54	10.43	8.66	8.60	7.87	9.69
	11.65	10.64	9.81	8.04	7.52	8.76
	8.73	8.08	7.70	6.44	6.42	8.93

diet which included 500 mg/day of choline for one week, and then were randomly assigned into two diet groups, one that contained choline and one that did not. During the fifth week of study all subjects again consumed a diet containing choline. Blood samples for choline analyses were obtained before the start of the study (day 0) and on days 7, 14, 21, 28 and 35. One subject had data missing for day 35.

A **multivariate analysis of covariance model** allowed testing the effects of diet on the plasma choline concentration over time, while controlling for treatment group differences in baseline choline levels. The null hypothesis of interest is a test of the trends by treatment interaction. This particular design has two treatment groups. With no missing data, all multivariate F tests would coincide. Hence we chose to use a single choice of N_* . The value of the unadjusted (multivariate) F statistic ($N_* = N$) is 2.77, with $v_1 = 4$, $v_2 = 8$ and p value of 0.102. For the same hypothesis, unadjusted $F_{GG} = 2.67$, $v_1 = 2.95$, $v_2 = 32.4$ and $p = 0.065$. Using N_{*2} gives an adjusted (multivariate) F of 2.65 with $v_1 = 4$ and $v_2 = 7$, leading to $p = 0.144$ and using N_{*9} gives $F_{GG} = 2.61$, $v_1 = 2.95$, $v_2 = 31.86$ and $p = 0.069$. Both missing data analysis methods led to a smaller p -value than the analysis based on 13 complete cases ($F = 2.00$, $v_1 = 4$, $v_2 = 7$ and $p = 0.20$; $F_{GG} = 2.40$, $v_1 = 2.8$, $v_2 = 28.3$ and $p = 0.092$).

5. CONCLUSIONS

A number of weaknesses in the performance of some of the methods should be noted. Not surprisingly, for all tests considered, accuracy decreases with more repeated measures, fewer subjects, more missing data and more correlation within subjects. In particular, the mixed model **F statistic used by PROC MIXED in SAS[®] with Satterthwaite approximated denominator degrees of freedom gives very liberal test size for $N \leq 24$, even with complete data.** We know of a number of statisticians seeking to find testing methods which perform better in small samples for mixed models. The inherent appeal of the methods lies in their great flexibility. The corresponding vast

range of such models makes finding universal solutions extremely difficult. We hope our results encourage more such research. In our simulations, for 5 per cent missing data, 6 responses and 12 subjects, the simple version of the EM algorithm used here failed roughly 90 per cent of the time. We expect that more sophisticated algorithms, especially better starting values, would avoid this problem most of the time. Of course, a small set of data with many holes may simply not contain enough information to support estimating a complex model involving means and covariances. This should not be seen as a failure of the method, but as a limitation of the data collection.

Despite the limitations just discussed, the new methods described here perform extremely well in small samples, and should be adopted. We base this conclusion on the observation that a degree of freedom adjustment always controlled test size at or below the nominal level in our simulations, even for conditions as extreme as $N = 12$ and 10 per cent missing data. Note that the choice of adjustment varies with the test. Using N_{*2} , the $\min N_{jj'}$, works best for the Wilks' and Hotelling–Lawley tests. In contrast, N_{*4} , the harmonic mean of $N_{jj'}$, works best for the Pillai–Bartlett test. Finally N_{*9} , the mean N_{jj} , works best for the Geisser–Greenhouse test.

A free copy of software which implements the new methods may be obtained from the website <http://www.bios.unc.edu/~muller> (choose LINMOD version 4.0). The code is written in SAS IML[®], and assumes the ability to describe and understand multivariate linear models in matrix notation. Earlier versions have been used for many years for data analysis of complex multivariate models and related teaching.

6. FUTURE RESEARCH

The simulated data were generated in such a way as to create data that are MCAR. Our motivation was twofold. First, good methods were not available even in the simplest case. Second, the study designs we simulate, although simple in some ways, occur very often in practice. It seems reasonable to hope that the techniques described here will also work with MAR data. Hence such data deserve attention in further simulations.

Techniques for power analysis, given that test size can be controlled, would be very useful. The approach taken here is more intuitive than analytical. Nevertheless, we believe we have succeeded in describing an approach to ensure that test size does not exceed the nominal rate in small, missing data samples for the GLMM. A more formal approach must necessarily involve a rather sophisticated attack, due to the complexity of the distributions for the multivariate test statistics, even with complete data.

ACKNOWLEDGEMENTS

An earlier version of this paper was one portion of the dissertation research for the first author's Dr PH degree program. The authors thank Drs J.D. Hosking, G. G. Koch, P. W. Stewart and D. J. Weber for helpful comments. Catellier's work supported in part by USFHS grant 90CA1467 and IPRIC grant 5-326-12. Muller's work supported in part by NIH grants PO1-CA47982-04, RO1-CA67183-01A1, RO1-CA72875, RO1-CA60193-04, MO1-RR000-46-33, NO1-ES-35356 and MH33127.

REFERENCES

1. Barton CN, Cramer EC. Hypothesis testing in multivariate linear models with randomly missing data. *Communications in Statistics – Simulations* 1989; **18**:875–895.

2. Schluchter MD, Elashoff JD. Small-sample adjustments to tests with unbalanced repeated measures assuming several covariance structures. *Journal of Statistical Computing – Simulations* 1990; **37**:69–87.
3. Rubin DB. Inference and missing data. *Biometrika* 1976; **63**:581–592.
4. Liang KY, Zeger SL. Longitudinal data analysis using generalized linear models. *Biometrika* 1986; **73**:13–22.
5. Park T. A comparison of the generalized estimating equation approach with the maximum likelihood approach for repeated measurements. *Statistics in Medicine* 1993; **12**:1723–1732.
6. Dempster AP, Laird NM, Rubin DB. Maximum likelihood from incomplete data via the EM algorithm (with Discussion). *Journal of the Royal Statistical Society-B* 1977; **39**:1–38.
7. Laird NM, Lange N, Stram D. Maximum likelihood computations with repeated measures: Application of the EM algorithm. *Journal of the American Statistical Association* 1987; **82**:97–105.
8. Stiger TR, Kosinski AS, Barnhart HX, Kleinbaum DG. ANOVA for repeated ordinal data with small sample size? A comparison of ANOVA, MANOVA, WLS and GEE methods by simulation. In *Joint Statistical Meetings of the American Statistical Association Abstract Book*. 1997; 246.
9. Qu YS, Piedmonte MR, Williams GW. Small sample validity of latent variable models for correlated binary data. *Communications in Statistics – Simulations* 1994; **23**:243–269.
10. 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.
11. Davidson ML. Univariate versus multivariate tests in repeated measures experiments. *Psychological Bulletin* 1972; **77**:446–452.
12. O'Brien RG, Kaiser MK. MANOVA method for analyzing repeated measures designs: An extensive primer. *Psychological Bulletin* 1985; **97**:316–333.
13. Muller KE, LaVange LM, Ramey SL, Ramey CT. Power calculations for general linear multivariate models including repeated measures applications. *Journal of the American Statistical Association* 1992; **87**:1209–1226.
14. Olson CL. Comparative robustness of six tests in multivariate analysis of variance. *Journal of the American Statistical Association* 1974; **69**:894–908.
15. Olson CL. Choosing a test statistic in multivariate analysis. *Psychological Bulletin* 1976; **86**:579–586.
16. Olson CL. Practical considerations in choosing a MANOVA test statistic: a rejoinder to Stevens. *Psychological Bulletin* 1979; **86**:1350–1352.
17. Anderson TW. *An Introduction To Multivariate Statistical Analysis*. Wiley: New York, 1984.
18. Rao CR. *Linear Statistical Inference and Its Applications*. Wiley: New York, 1973.
19. Pillai KCS. On some distribution problems in multivariate analysis. In *Institute of Statistics Mimeo Series No. 88*. University of North Carolina: Chapel Hill, 1954.
20. Muller KE. A new F approximation for the Pillai-Bartlett trace under H_0 . *Journal of Computational and Graphical Statistics* 1998; **7**:131–137.
21. McKeon JJ. F approximations to the distribution of Hotelling's T_0^2 . *Biometrika* 1974; **61**:381–383.
22. Muller KE, Barton CN. Approximate power for repeated measures for repeated measures anova lacking sphericity. *Journal of the American Statistical Association* 1989; **84**:549–555.
23. Muller KE, Barton CN. Correction to 'Approximate power for repeated measures for repeated measures anova lacking sphericity'. *Journal of the American Statistical Association* 1991; **86**:255–256.
24. Little RJA, Rubin DB. *Statistical Analysis with Missing Data*. Wiley: New York, 1987.
25. Rubin DB. Characterization of estimation of parameters in incomplete data problems. *Journal of the American Statistical Association* 1974; **69**:467–474.
26. Orchard T, Woodbury MA. A missing information principle: theory and applications. In *Proceedings of the 6th Berkeley Symposium on Mathematical Statistics and Probability 1*. University of California Press: Berkeley, 1972; 697–715.
27. Beale EML, Little RJA. Missing values in multivariate analysis. *Journal of the Royal Statistical Society, Series B* 1975; **37**:129–145.
28. Mensah RD, Elswick RK, Chinchilli VM. Consistent estimators of the variance-covariance matrix of the GMANOVA model with missing data. *Communications in Statistics, Series A* 1993; **22**:1495–1514.
29. Callahan TP, Harville DA. Some new algorithms for computing maximum likelihood estimates of variance components. *Journal of Statistical Computation and Simulation* 1990; **38**:239–259.
30. Morrison DF. A test for equality of means of correlated variates with missing data on one response. *Biometrika* 1973; **60**:101–105.
31. Little RJA. Inference about means from incomplete multivariate data. *Biometrika* 1976; **63**:593–604.
32. Little RJA. Approximate calibrated small sample inference about means from bivariate normal data with missing values. *Computational Statistics and Data Analysis* 1988; **7**:161–178.
33. Lindstrom MJ, Bates DM. Newton-Raphson and EM algorithms for linear mixed-effects models for repeated measures data. *Journal of the American Statistical Association* 1988; **83**:1014–1022.
34. Hocking RR. *The Analysis of Linear Models*. Brooks/Cole: Monterey, 1985.
35. Woolson RF, Leeper JD. Hypothesis testing in multivariate linear models with randomly missing data. *Communications in Statistics – Theory and Methods* 1980; **A9**:1491–1513.

36. Woolson RF, Leeper JD, Clarke WR. Analysis of incomplete data from longitudinal and mixed longitudinal studies. *Journal of the Royal Statistical Society, Series A* 1978; **141**:242–252.
37. Leeper JD, Woolson RF. Testing hypotheses for the growth curve model when the data are incomplete. *Journal of Statistical Computation and Simulation* 1982; **15**:97–107.
38. SAS Institute Inc. *SAS/STAT Software: Changes and enhancements, Release 6.12*. SAS Institute Inc.: Cary, 1997; 607, 644.
39. Zeisel SH, DaCosta K, Franklin PD, Alexander EA, Lamont JT, Sheard NF, Beiser A. Choline, an essential nutrient for humans. *FASEB* 1991; **5**:2093–2098.