

EXTENSION OF FACTORING LIKELIHOOD APPROACH TO NON-MONOTONE MISSING DATA[†]

JAE KWANG KIM¹

ABSTRACT

We address the problem of parameter estimation in multivariate distributions under ignorable non-monotone missing data. The factoring likelihood method for monotone missing data, termed by Rubin (1974), is extended to a more general case of non-monotone missing data. The proposed method is algebraically equivalent to the Newton-Raphson method for the observed likelihood, but avoids the burden of computing the first and the second partial derivatives of the observed likelihood. Instead, the maximum likelihood estimates and their information matrices for each partition of the data set are computed separately and combined naturally using the generalized least squares method.

AMS 2000 subject classifications. Primary 62F10; Secondary 62D05.

Keywords. EM algorithm, generalized least squares, Gauss-Newton method, missing at random.

1. INTRODUCTION

Missing data is quite common in practice. Statistical analysis of data with missing value is an important practical problem because we cannot simply ignore the missing data. When we simply ignore the missing part of the data, the resulting estimates will have nonresponse bias if the responding part of the data is systematically different from the nonresponding part of the data. In addition to the nonresponse bias, efficiency of the resulting estimates should be also considered because we might lose some information observed in the partially missing

Received April 2003; accepted August 2004.

[†]The work was partially supported by Yonsei University College of Business and Economics Research Fund of 2004.

¹Department of Applied Statistics, Yonsei University, Seoul 120-749, Korea (e-mail : kimj@yonsei.ac.kr)

data. Groves *et al.* (2001) and Little and Rubin (2002) provide comprehensive overview of the missing data problem.

To explain the basic idea of the existing methods, we use an example of bivariate normal data. Let (Y_{1i}, Y_{2i}) be a vector of bivariate normal random variable distributed as

$$\begin{pmatrix} Y_{1i} \\ Y_{2i} \end{pmatrix} \sim iid N \left[\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{pmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{12} & \sigma_{22} \end{pmatrix} \right], \quad (1.1)$$

where *iid* is the abbreviation of independently and identically distributed. Note that the five parameters, μ_1 , μ_2 , σ_{11} , σ_{12} , and σ_{22} , are needed to identify the bivariate normal distribution. We assume that the observations are missing at random (MAR) in the sense of Rubin (1976) so that the relevant likelihood is the observed likelihood, the marginal likelihood of the observed data. Under MAR, we can ignore the response mechanism when estimating the population parameters.

To estimate the parameters, a direct maximum likelihood method that maximizes the observed likelihood can be used. To do this, we need to compute the observed likelihood and its partial derivatives. The Newton-Raphson type solution to the likelihood equation also requires the computation of the second-order partial derivatives. The computation of the first and the second partial derivatives can be cumbersome.

The factoring likelihood method, termed by Rubin (1974), avoids the burden of computing the partial derivatives and still computes the maximum likelihood estimate (MLE) of the observed likelihood. The factoring likelihood method computes MLE's easily, but is applicable only to the monotone missing data. For the definition of monotone missing pattern and the non-monotone missing pattern, see Little and Rubin (2002, Section 1.2). Under monotone missing pattern, the observed likelihood can be factored into the marginal likelihood and the conditional likelihood so that the maximum likelihood estimates can be estimated separately at each likelihood. For example, assume that Y_1 is fully observed with n observations and Y_2 is subject to missing with r ($< n$) observations. Anderson (1957) first consider the estimation of parameters under this setup by using an alternative representation of the bivariate normal distribution as

$$\begin{aligned} Y_{1i} &\sim iid N(\mu_1, \sigma_{11}), \\ Y_{2i} | (Y_{1i} = y_{1i}) &\sim iid N(\beta_{20.1} + \beta_{21.1}y_{1i}, \sigma_{22.1}), \end{aligned} \quad (1.2)$$

where $\beta_{20.1} = \mu_2 - \beta_{21.1}\mu_1$, $\beta_{21.1} = \sigma_{11}^{-1}\sigma_{12}$ and $\sigma_{22.1} = \sigma_{22} - \beta_{21.1}^2\sigma_{11}$. The observed likelihood is then written as a product of marginal likelihood of a fully

observed variable Y_1 and the conditional likelihood of Y_2 given Y_1 . Thus, the parameters μ_1 and σ_{11} for the marginal distribution of Y_1 can be estimated with n observations and the other regression parameters, $\beta_{20\cdot1}$, $\beta_{21\cdot1}$, and $\sigma_{22\cdot1}$, can be estimated from the conditional distribution with r observations.

Note that the factoring likelihood approach consists of two steps. In the first step, the likelihood is factored, and in the second the MLE for each likelihood is computed separately. The advantage of the factoring likelihood method is that the MLE's are easily computed because the marginal and the conditional likelihoods are of known form and thus we can directly use the known solutions of the likelihood equations for each likelihood. For the monotone missing data, the MLE's for the conditional distribution are independent of the MLE's for the marginal distribution. This is because two sets of parameters, the parameters for the marginal likelihood and those for the conditional likelihood, are orthogonal. Following Cox and Reid (1987), two parameters, θ_1 and θ_2 , are called orthogonal with respect to the likelihood $l(\theta_1, \theta_2)$ if $E(\partial^2 \log l / \partial \theta_1 \partial \theta_2) = 0$. Thus, because of the orthogonality of the parameters, the MLE's for the conditional likelihood are not affected by the MLE's for the marginal likelihood. Rubin (1974) recommends the factoring likelihood approach as a general framework in the analysis of missing data with monotone missing pattern.

Under the non-monotone missing pattern, the factoring likelihood approach is not directly applicable because the parameters are no longer orthogonal. The EM algorithm, proposed by Dempster *et al.* (1977), can be used to compute the MLE's under the general missing pattern, but uses an iterative procedure and does not provide the information matrix directly.

In this paper, we consider an extension of the factoring likelihood method to the non-monotone missing data. Note that the factoring likelihood method ease the computation of the MLE's but the resulting estimators are no longer independent because of the non-orthogonality of the parameters. Thus, in addition of the two steps in the original factoring likelihood approach, we need another step to combine these separate MLE's computed within each likelihood to produce the final MLE's. The proposed method turns out to be essentially the same as the direct maximum likelihood method using the Newton-Raphson algorithm but has some computational advantages. The proposed method is described under the bivariate normal setup in Section 2. A justification of the proposed method under a more general setup is made in Section 3. The proposed method is applied to a categorical data example in Section 4. Concluding remarks are made in Section 5.

2. ONE-STEP ESTIMATOR

The proposed method can be described into three steps:

Step 1. Partition the original sample into several disjoint sets according to the missing pattern.

Step 2. Compute MLE for each identified parameter separately in each partition of the sample.

Step 3. Combine the estimators to get a set of final estimates using a generalized least squares (GLS) form.

In Step 1, with a non-monotone missing pattern with two variables, we have three types of respondents that contain information about the parameters. The first set H of units has both Y_1 and Y_2 observed, the second set K of units has Y_1 observed but Y_2 missing, and the third set L of units has Y_2 observed but Y_1 missing. That is, we partition the sample into several disjoint sets according to the pattern of missingness. We also define M to be the set of units that has both Y_1 and Y_2 missing. Let n_H, n_K, n_L , and n_M be the sample size of the set H, K, L , and M , respectively. Note that $n = n_H + n_K + n_L + n_M$.

In Step 2, we obtain the following estimators in each set: For set H , we get the ML estimates for the five parameters in (1.2): $\hat{\beta}_{20\cdot1,H}, \hat{\beta}_{21\cdot1,H}, \hat{\sigma}_{22\cdot1,H}, \hat{\mu}_{1,H}$, and $\hat{\sigma}_{11,H}$. For the set K , the ML estimates $\hat{\mu}_{1,K}$ and $\hat{\sigma}_{11,K}$ are obtained for μ_1 and σ_{11} , respectively. For the set L , the ML estimates $\hat{\mu}_{2,L}$ and $\hat{\sigma}_{22,L}$ are obtained for $\mu_2 = \beta_{20\cdot1} + \beta_{21\cdot1}\mu_1$ and $\sigma_{22} = \sigma_{22\cdot1} + \beta_{21\cdot1}^2\sigma_{11}$, respectively.

In Step 3, we use the GLS method to combine the nine estimators into an estimator for the five parameters. The nine estimates are

$$\hat{\eta} = \left(\hat{\beta}_{20\cdot1,H}, \hat{\beta}_{21\cdot1,H}, \hat{\sigma}_{22\cdot1,H}, \hat{\mu}_{1,H}, \hat{\sigma}_{11,H}, \hat{\mu}_{1,K}, \hat{\sigma}_{11,K}, \hat{\mu}_{2,L}, \hat{\sigma}_{22,L} \right)'. \quad (2.1)$$

The expected values of the nine estimates are

$$\eta(\theta) = (\beta_{20\cdot1}, \beta_{21\cdot1}, \sigma_{22\cdot1}, \mu_1, \sigma_{11}, \mu_1, \sigma_{11}, \beta_{20\cdot1} + \beta_{21\cdot1}\mu_1, \sigma_{22\cdot1} + \beta_{21\cdot1}^2\sigma_{11})' \quad (2.2)$$

and the asymptotic covariance matrix is

$$\mathbf{V} = \text{diag} \left\{ \Sigma_{bb}, \frac{2\sigma_{22\cdot1}^2}{n_H}, \frac{\sigma_{11}}{n_H}, \frac{2\sigma_{11}^2}{n_H}, \frac{\sigma_{11}}{n_K}, \frac{2\sigma_{11}^2}{n_K}, \frac{\sigma_{22}}{n_L}, \frac{2\sigma_{22}^2}{n_L} \right\}, \quad (2.3)$$

where $\theta = (\beta_{20\cdot1}, \beta_{21\cdot1}, \sigma_{22\cdot1}, \mu_1, \sigma_{11})'$ and

$$\Sigma_{bb} = \begin{pmatrix} n_H^{-1}\sigma_{22\cdot1}(1 + \sigma_{11}^{-1}\mu_1^2) & -n_H^{-1}\sigma_{11}^{-1}\sigma_{22\cdot1}\mu_1 \\ -n_H^{-1}\sigma_{11}^{-1}\sigma_{22\cdot1}\mu_1 & n_H^{-1}\sigma_{11}^{-1}\sigma_{22\cdot1} \end{pmatrix}.$$

Derivation for the asymptotic covariance matrix of the first five estimates in (2.1) can be found, for example, in Subsection 7.2.2 of Little and Rubin (2002). Independence between $\hat{\mu}_{1,K}$ and $\hat{\sigma}_{11,K}$ comes from the property of the normal distribution. Observations between different sets are independent because of the *iid* setup.

The GLS formulation in (2.2) and (2.3) is a nonlinear model of the five parameters. Using a Taylor expansion on the nonlinear model, a step of the Gauss-Newton method can be formulated as

$$\mathbf{e}_\eta = \mathbf{X}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_S) + \mathbf{u}, \tag{2.4}$$

where $\mathbf{e}_\eta = \hat{\boldsymbol{\eta}} - \boldsymbol{\eta}(\hat{\boldsymbol{\theta}}_S)$, $\hat{\boldsymbol{\theta}}_S$ is an initial estimator of $\boldsymbol{\theta}$, $\boldsymbol{\eta}(\hat{\boldsymbol{\theta}}_S)$ is the vector (2.2) evaluated at $\hat{\boldsymbol{\theta}}_S$,

$$\mathbf{X} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & \mu_1 & 2\beta_{21\cdot 1}\sigma_{11} \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 & 1 & 0 & \beta_{21\cdot 1} & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 1 & 0 & \beta_{21\cdot 1}^2 \end{pmatrix}', \tag{2.5}$$

and, approximately,

$$\mathbf{u} \sim (\mathbf{0}, \mathbf{V}),$$

where \mathbf{V} is the covariance matrix defined in (2.3). For a brief description of the Gauss-Newton method for the estimation of nonlinear models, see Fuller (1996, Section 5.5).

The procedure can be carried out iteratively until convergence, but we used a single step of the procedure. For a suitable choice of the initial estimates, the one-step estimator is a very good approximation to the maximum likelihood estimator. The estimator is

$$\hat{\boldsymbol{\theta}} = \hat{\boldsymbol{\theta}}_S + \left(\mathbf{X}'_S \hat{\mathbf{V}}_S^{-1} \mathbf{X}_S\right)^{-1} \mathbf{X}'_S \hat{\mathbf{V}}_S^{-1} \mathbf{e}_\eta, \tag{2.6}$$

where \mathbf{X}_S and $\hat{\mathbf{V}}_S$ are evaluated from \mathbf{X} in (2.5) and \mathbf{V} in (2.3), respectively, using the initial values of $\boldsymbol{\theta}$. The covariance matrix of the estimator in (2.6) can be estimated by

$$\mathbf{C} = \left(\mathbf{X}'_S \hat{\mathbf{V}}_S^{-1} \mathbf{X}_S\right)^{-1}. \tag{2.7}$$

The initial values for the iterative procedure are $\hat{\beta}_{20\cdot 1,H}$, $\hat{\beta}_{21\cdot 1,H}$, $\hat{\sigma}_{11\cdot 2,H}$,

$$\hat{\mu}_1 = (n_H + n_K)^{-1} (n_H \bar{y}_{1,H} + n_K \bar{y}_{1,K}),$$

and

$$\tilde{\sigma}_{11} = (n_H + n_K - 2)^{-1} [(n_H - 1) s_{1,H}^2 + (n_K - 1) s_{1,K}^2],$$

where $\bar{y}_{1,H}$ and $\bar{y}_{1,K}$ are the sample means of Y_1 in the sets H and K , respectively, and $s_{1,H}^2$ and $s_{1,K}^2$ are the sample variances of Y_1 in the sets H and K .

3. JUSTIFICATION

Let the score function of a likelihood be defined as

$$S(\mathbf{y}; \boldsymbol{\theta}) = \frac{\partial \log l(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}},$$

where $l(\boldsymbol{\theta}) = \prod_i f(\mathbf{y}_i; \boldsymbol{\theta})$ is the likelihood function of parameter $\boldsymbol{\theta}$. The maximum likelihood estimator of $\boldsymbol{\theta}$ can be defined as a solution to the Newton-Raphson variant of scoring method

$$\boldsymbol{\theta}^{(k+1)} = \boldsymbol{\theta}^{(k)} + [\mathbf{I}(\boldsymbol{\theta}^{(k)})]^{-1} S(\mathbf{y}; \boldsymbol{\theta}^{(k)}), \quad (3.1)$$

where $\mathbf{I}(\boldsymbol{\theta}) = E[-\partial^2 \log l(\boldsymbol{\theta}) / \partial \boldsymbol{\theta}^2]$ is the expected information matrix for $\boldsymbol{\theta}$. It is known that if the starting value is a \sqrt{n} -consistent estimator of $\boldsymbol{\theta}$, then one-step iterate $\boldsymbol{\theta}^{(1)}$ in (3.1) is asymptotically equivalent to the maximum likelihood estimator of $\boldsymbol{\theta}$ (e.g. Lehmann, 1983, Theorem 3.1, p. 422).

Now, under the missing data structure in Section 2, we show that the one-step estimator in (2.6) is equivalent to the Newton-Raphson solution in (3.1). Note that the observed log-likelihood can be written as a sum of the log-likelihood in each set:

$$\log l(\boldsymbol{\theta}) = \log l_H(\boldsymbol{\theta}) + \log l_K(\boldsymbol{\theta}) + \log l_L(\boldsymbol{\theta}), \quad (3.2)$$

where $l_H = \prod_{i \in H} f(\mathbf{y}_i; \boldsymbol{\theta})$ is the likelihood function defined in set H , and l_K and l_L are defined similarly. Under MAR, l_H is the likelihood for the joint distribution of Y_1 and Y_2 , l_K is the likelihood for the marginal distribution of Y_1 , and l_L is the likelihood for the marginal distribution of Y_2 . By (3.2), the score function for the likelihood can be written as

$$S(\mathbf{y}; \boldsymbol{\theta}) = S_H(\mathbf{y}; \boldsymbol{\theta}) + S_K(\mathbf{y}; \boldsymbol{\theta}) + S_L(\mathbf{y}; \boldsymbol{\theta}) \quad (3.3)$$

and the expected information matrix also satisfies the additive decomposition:

$$\mathbf{I}(\boldsymbol{\theta}) = \mathbf{I}_H(\boldsymbol{\theta}) + \mathbf{I}_K(\boldsymbol{\theta}) + \mathbf{I}_L(\boldsymbol{\theta}), \quad (3.4)$$

where $\mathbf{I}_H(\boldsymbol{\theta}) = E[-\partial^2 \log l_H(\boldsymbol{\theta})/\partial\boldsymbol{\theta}^2]$, and $\mathbf{I}_K(\boldsymbol{\theta})$ and $\mathbf{I}_L(\boldsymbol{\theta})$ are defined similarly. Let $\boldsymbol{\eta}_H = \boldsymbol{\eta}_H(\boldsymbol{\theta})$ be a parametrization that $\mathbf{I}_H(\boldsymbol{\eta}_H)$ matrix is easy to compute. One such parametrization is $\boldsymbol{\eta}_H = (\boldsymbol{\eta}_{H1}, \boldsymbol{\eta}_{H2})$, where $\boldsymbol{\eta}_{H1}$ is the parameters for the conditional distribution and $\boldsymbol{\eta}_{H2}$ is the parameters for the marginal distribution. Since the parameters for the conditional distribution are orthogonal to those for the marginal distribution, the parametrization $\boldsymbol{\eta}_H = (\boldsymbol{\eta}_{H1}, \boldsymbol{\eta}_{H2})$ makes the $\mathbf{I}_H(\boldsymbol{\eta}_H)$ matrix block-diagonal. The parametrization for the set H need not be the same as that for the set K nor for the set L , providing more flexibility in choosing the parametrization. Separate orthogonal parametrization in each set will lead to computational advantages over the direct maximum likelihood method.

The equation in (3.4) can be written as

$$\begin{aligned} \mathbf{I}(\boldsymbol{\theta}) &= \left(\frac{\partial\boldsymbol{\eta}_H}{\partial\boldsymbol{\theta}}\right) \mathbf{I}_H(\boldsymbol{\eta}_H) \left(\frac{\partial\boldsymbol{\eta}_H}{\partial\boldsymbol{\theta}}\right)' + \left(\frac{\partial\boldsymbol{\eta}_K}{\partial\boldsymbol{\theta}}\right) \mathbf{I}_K(\boldsymbol{\eta}_K) \left(\frac{\partial\boldsymbol{\eta}_K}{\partial\boldsymbol{\theta}}\right)' \\ &\quad + \left(\frac{\partial\boldsymbol{\eta}_L}{\partial\boldsymbol{\theta}}\right) \mathbf{I}_L(\boldsymbol{\eta}_L) \left(\frac{\partial\boldsymbol{\eta}_L}{\partial\boldsymbol{\theta}}\right)' \\ &= \mathbf{X}'\hat{\mathbf{V}}^{-1}\mathbf{X}, \end{aligned} \tag{3.5}$$

where $\mathbf{X}' = (\partial\boldsymbol{\eta}_H/\partial\boldsymbol{\theta}, \partial\boldsymbol{\eta}_K/\partial\boldsymbol{\theta}, \partial\boldsymbol{\eta}_L/\partial\boldsymbol{\theta})$ and $\hat{\mathbf{V}}^{-1} = \text{diag}\{\mathbf{I}_H(\boldsymbol{\eta}_H), \mathbf{I}_K(\boldsymbol{\eta}_K), \mathbf{I}_L(\boldsymbol{\eta}_L)\}$. Now, consider the score function in (3.3). Using the chain rule, the score function can be written as

$$S(\mathbf{y}; \boldsymbol{\theta}) = \left(\frac{\partial\boldsymbol{\eta}_H}{\partial\boldsymbol{\theta}}\right) S_H(\mathbf{y}; \boldsymbol{\eta}_H) + \left(\frac{\partial\boldsymbol{\eta}_K}{\partial\boldsymbol{\theta}}\right) S_K(\mathbf{y}; \boldsymbol{\eta}_K) + \left(\frac{\partial\boldsymbol{\eta}_L}{\partial\boldsymbol{\theta}}\right) S_L(\mathbf{y}; \boldsymbol{\eta}_L). \tag{3.6}$$

Let $\hat{\boldsymbol{\eta}}_H$ be the MLE of the likelihood l_H . Taking a Taylor expansion of $S_H(\mathbf{y}; \boldsymbol{\eta}_H)$ about $\hat{\boldsymbol{\eta}}_H$ leads to

$$S_H(\mathbf{y}; \boldsymbol{\eta}_H) \doteq S_H(\mathbf{y}; \hat{\boldsymbol{\eta}}_H) - \mathcal{I}_H(\hat{\boldsymbol{\eta}}_H)(\boldsymbol{\eta}_H - \hat{\boldsymbol{\eta}}_H),$$

where $\mathcal{I}_H(\boldsymbol{\eta}_H) = -\partial^2 \log l_H(\boldsymbol{\eta}_H)/\partial\boldsymbol{\eta}_H^2$. Using $S_H(\mathbf{y}; \hat{\boldsymbol{\eta}}_H) = 0$ and the weak convergence of the observed information matrix to the expected information matrix, we have

$$S_H(\mathbf{y}; \boldsymbol{\eta}_H) \doteq -\mathcal{I}_H(\hat{\boldsymbol{\eta}}_H)(\boldsymbol{\eta}_H - \hat{\boldsymbol{\eta}}_H).$$

Similar results hold for the sets K and L . Thus, (3.6) becomes

$$\begin{aligned} S(\mathbf{y}; \boldsymbol{\theta}) &\doteq \left(\frac{\partial\boldsymbol{\eta}_H}{\partial\boldsymbol{\theta}}\right) \mathcal{I}_K(\hat{\boldsymbol{\eta}}_H)(\hat{\boldsymbol{\eta}}_H - \boldsymbol{\eta}_H) + \left(\frac{\partial\boldsymbol{\eta}_K}{\partial\boldsymbol{\theta}}\right) \mathcal{I}_K(\hat{\boldsymbol{\eta}}_K)(\hat{\boldsymbol{\eta}}_K - \boldsymbol{\eta}_K) \\ &\quad + \left(\frac{\partial\boldsymbol{\eta}_L}{\partial\boldsymbol{\theta}}\right) \mathcal{I}_L(\hat{\boldsymbol{\eta}}_L)(\hat{\boldsymbol{\eta}}_L - \boldsymbol{\eta}_L) \\ &= \mathbf{X}'\hat{\mathbf{V}}^{-1}(\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}), \end{aligned} \tag{3.7}$$

TABLE 3.1 2×2 table with supplemental margins for both variables

Set	Y_1	Y_2	Count
H	1	1	100
	1	2	50
	2	1	75
	2	2	75
K	1		30
	2		60
L		1	28
		2	60

where $\boldsymbol{\eta} = (\boldsymbol{\eta}'_H, \boldsymbol{\eta}'_K, \boldsymbol{\eta}'_L)'$ and $\hat{\boldsymbol{\eta}} = (\hat{\boldsymbol{\eta}}'_H, \hat{\boldsymbol{\eta}}'_K, \hat{\boldsymbol{\eta}}'_L)'$. Therefore, inserting (3.5) and (3.7) into (3.1), we have

$$\boldsymbol{\theta}^{(k+1)} = \boldsymbol{\theta}^{(k)} + [\mathbf{X}'\hat{\mathbf{V}}^{-1}\mathbf{X}]^{-1} \mathbf{X}'\hat{\mathbf{V}}^{-1} [\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}(\boldsymbol{\theta}^{(k)})], \quad (3.8)$$

which is equivalent to the expression in (2.6).

4. A NUMERICAL EXAMPLE

For a numerical example, we consider the data set originally presented by Little (1982) and also discussed in Little and Rubin (2002). Table 3.1 gives the data for a 2×2 table with supplemental margins for both the classifying variables. According to Little (1982), the final probabilities of classification obtained from EM algorithm are

$$\hat{\pi}_{11} = 0.28, \quad \hat{\pi}_{12} = 0.17, \quad \hat{\pi}_{21} = 0.24, \quad \hat{\pi}_{22} = 0.31, \quad (4.1)$$

where $\pi_{ij} = \Pr(Y_1 = i, Y_2 = j)$, $i, j = 1, 2$.

For the orthogonal parametrization, we use

$$\boldsymbol{\eta}_H = (\pi_{1|1}, \pi_{1|2}, \pi_{+1})'$$

where

$$\pi_{1|1} = \Pr(Y_1 = 1 \mid Y_2 = 1),$$

$$\pi_{1|2} = \Pr(Y_1 = 1 \mid Y_2 = 2),$$

$$\pi_{+1} = \Pr(Y_2 = 1).$$

We also set $\theta = \eta_H$. Note that the validity of the proposed method does not depend on the choice of the parametrization. A suitable parametrization will make the computation of the information matrix simple.

From the data in Table 3.1, the five observations for three parameters are

$$\begin{aligned} \hat{\eta} &= (\hat{\pi}_{1|1,H}, \hat{\pi}_{1|2,H}, \hat{\pi}_{+1,H}, \hat{\pi}_{1+,K}, \hat{\pi}_{+1,L})' \\ &= \left(\frac{100}{175}, \frac{50}{125}, \frac{175}{300}, \frac{30}{90}, \frac{28}{88} \right)' \end{aligned}$$

with the expectations

$$\eta(\theta) = (\pi_{1|1}, \pi_{1|2}, \pi_{+1}, \pi_{1|1}\pi_{+1} + \pi_{1|2} - \pi_{1|2}\pi_{+1}, \pi_{+1})'$$

and the variance-covariance matrix

$$\mathbf{V} = \text{diag} \left\{ \frac{\pi_{1|1}(1 - \pi_{1|1})}{n_H}, \frac{\pi_{1|2}(1 - \pi_{1|2})}{n_H}, \frac{\pi_{+1}(1 - \pi_{+1})}{n_H}, \frac{\pi_{1+}(1 - \pi_{1+})}{n_K}, \frac{\pi_{+1}(1 - \pi_{+1})}{n_L} \right\}.$$

The Gauss-Newton method as in (2.4) can be used to solve the nonlinear model of three parameters, where the initial estimator of θ is $\hat{\theta}_S = (100/175, 50/125, 203/388)'$ and the \mathbf{X} matrix is

$$\mathbf{X} = \begin{pmatrix} 1 & 0 & 0 & \pi_{+1} & 0 \\ 0 & 1 & 0 & 1 - \pi_{+1} & 0 \\ 0 & 0 & 1 & \pi_{1|1} - \pi_{1|2} & 1 \end{pmatrix}'. \tag{4.2}$$

The resulting one-step estimates are

$$\hat{\pi}_{11} = 0.29, \quad \hat{\pi}_{12} = 0.18, \quad \hat{\pi}_{21} = 0.23, \quad \hat{\pi}_{22} = 0.30,$$

which is close to the final results in (4.1) obtained from the EM algorithm.

5. CONCLUDING REMARKS

The proposed method is shown to be algebraically equivalent to the scoring method for maximum likelihood estimation but avoids the burden of obtaining the observed likelihood. Instead, the MLE's separately computed from each partition of the marginal likelihoods and the full likelihoods are combined in a natural

way. The way we combine the information takes the form of GLS estimation and thus can be easily implemented using the existing software. The estimated covariance matrix is obtained automatically as a by-product of the computation. The proposed method is not restricted to the bivariate normal distribution. It can be applied to any parametric multivariate distribution as long as the computation for the marginal likelihood and the full likelihood are relatively easier than that of the observed likelihood.

The proposed method assumes ignorable response mechanism. A more realistic situation will be the case where the probability of Y_2 missing depends on the value of Y_1 . In this case, the assumption of missing at random no longer holds and we have to take the response mechanism into account. Other existing likelihood-based methods, such as EM algorithm, cannot handle the situation, either.

ACKNOWLEDGEMENTS

We thank anonymous referees and the editor for helpful comments, which greatly improved the presentation of the paper.

REFERENCES

- ANDERSON, T. W. (1957). "Maximum likelihood estimates for a multivariate normal distribution when some observations are missing", *Journal of the American Statistical Association*, **52**, 200–203.
- COX, D. R. AND REID, N. (1987). "Parameter orthogonality and approximate conditional inference (with discussion)", *Journal of the Royal Statistical Society*, **B49**, 1–39.
- DEMPSTER, A. P., LAIRD, N. M. AND RUBIN, D. B. (1977). "Maximum likelihood from incomplete data via the EM algorithm (with discussion)", *Journal of the Royal Statistical Society*, **B39**, 1–38.
- FULLER, W. A. (1996). *Introduction to Statistical Times Series*, 2nd ed., John Wiley & Sons, New York.
- GROVES, R. M., DILLMAN, D. A., ELTINGE, J. L. AND LITTLE, R. J. A. (2001). *Survey Nonresponse*, John Wiley & Sons, New York.
- LEHMANN, E. L. (1983). *Theory of Point Estimation*, John Wiley & Sons, New York.
- LITTLE, R. J. A. (1982). "Models for nonresponse in sample surveys", *Journal of the American Statistical Association*, **77**, 237–250.
- LITTLE, R. J. A. AND RUBIN, D. B. (2002). *Statistical Analysis with Missing Data*, 2nd ed., John Wiley & Sons, New York.
- RUBIN, D. B. (1974). "Characterizing the estimation of parameters in incomplete-data problems", *Journal of the American Statistical Association*, **69**, 467–474.
- RUBIN, D. B. (1976). "Inference and missing data", *Biometrika*, **63**, 581–592.