



Taylor & Francis  
Taylor & Francis Group



---

Maximum Likelihood Variance Components Estimation for Binary Data

Author(s): Charles E. McCulloch

Source: *Journal of the American Statistical Association*, Mar., 1994, Vol. 89, No. 425 (Mar., 1994), pp. 330-335

Published by: Taylor & Francis, Ltd. on behalf of the American Statistical Association

Stable URL: <https://www.jstor.org/stable/2291229>

---

JSTOR is a not-for-profit service that helps scholars, researchers, and students discover, use, and build upon a wide range of content in a trusted digital archive. We use information technology and tools to increase productivity and facilitate new forms of scholarship. For more information about JSTOR, please contact [support@jstor.org](mailto:support@jstor.org).

Your use of the JSTOR archive indicates your acceptance of the Terms & Conditions of Use, available at <https://about.jstor.org/terms>



JSTOR

Taylor & Francis, Ltd. and American Statistical Association are collaborating with JSTOR to digitize, preserve and extend access to *Journal of the American Statistical Association*

# Maximum Likelihood Variance Components Estimation for Binary Data

Charles E. McCulloch\*

We consider a class of probit-normal models for binary data and describe ML and REML estimation of variance components for that class as well as best prediction for the realized values of the random effects. ML estimates are calculated using an EM algorithm; for complicated models EM includes a Gibbs step. The computations are illustrated through two examples.

KEY WORDS: EM algorithm; Fixed and random effects; Generalized linear models; Monte Carlo Markov chain; Probit.

## 1. INTRODUCTION

Models for continuous data that incorporate both fixed and random effects (mixed models) are commonly used in various disciplines, from ecology and medicine to the physical sciences. But the same is not true for binary data (Stram, Wei, and Ware 1988); usage has been limited to a large extent by the intractability of the computations involved in fitting many of the models. In this article we consider a class of probit-normal models. We describe maximum likelihood (ML) and restricted maximum likelihood (REML) estimation of the parameters in the model by use of the EM algorithm (Dempster, Laird, and Rubin 1977). Our version of the EM algorithm is very similar to that for the continuous normal linear model and offers a framework for computation of the ML and REML estimates. We demonstrate through two examples that the computations are feasible for any number and structure of random effects and an arbitrary number of fixed effects. This has not previously been possible; ML estimation has been described only in models with nested random effects.

Our focus will be on variance components estimation in mixed models and the analogs of best linear unbiased prediction (BLUP) of the observed values of the random effects. Thus our concentration differs somewhat from the usual one of repeated measures models, which is to treat the fixed effects as the primary quantities of interest, with the random effects introducing a "nuisance" correlation. We do not consider covariance components models.

A number of models for correlated binary data have been proposed. The beta-binomial distribution is a natural model to use (Crowder 1978; Williams 1975). This model hypothesizes a mixing distribution directly on the probability of success; however, it does not generalize easily to multiple random effects. Zeger and Liang (1986) and Liang and Zeger (1986) have proposed generalizations of quasi-likelihood methods, but their methods focus on the fixed effects and only estimate the variances and covariances as nuisance parameters. Prentice (1988) has considered extensions of the Zeger and Liang (1986) estimating equation approach, explicitly estimating the covariances as well. But like the beta-binomial models, Prentice's models are difficult to generalize to multiple random effects. For these reasons we consider correlated probit models that are generalizations of those of Ochi and Prentice (1984). These are similar to the logit-

normal models of Pierce and Sands (1975), Wong and Mason (1984), and Stiratelli, Laird, and Ware (1984) (although Stiratelli et al.'s models are intended only for the longitudinal data setting). Our model is essentially a simplified version of the threshold model considered by Harville and Mee (1984); but the computations for their general model were deemed "insurmountable" (p. 397), and they were forced to resort to ad hoc estimation methods. Zeger, Liang, and Albert (1988), Liang, Zeger, and Qaqish (1992), and Anderson, Gilmour, and Rae (1985) considered a generalized estimating equation approach, and Anderson and Aitken (1985) considered an iterative weighted logit analysis approach with models similar to ours. Other authors who have considered related models are Preisler (1989), Im and Gianola (1988), Gianola (1980), Quaas and Van Vleck (1980), and Manski and McFadden (1981).

## 2. THE MODEL

Our model is a threshold model where  $\mathbf{Y}$  represents an unobserved continuous variable and we observe only  $W_i = I_{\{Y_i > 0\}}$ ; that is, whether  $Y_i$  exceeds a threshold of 0. A flexible class of binary data models can be generated by assuming

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \boldsymbol{\varepsilon}, \quad (1)$$

$$W_i = I_{\{Y_i > 0\}}, \quad i = 1, 2, \dots, n,$$

where  $\mathbf{X}$  and  $\mathbf{Z}$  are known matrices,  $\mathbf{u} \sim \mathcal{N}(\mathbf{0}, \mathbf{D})$ , and  $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , independently of  $\mathbf{u}$ . It is unimportant whether we actually believe in the threshold model and the unobserved variable  $\mathbf{Y}$  or if we merely use it as a device to obtain estimates for the model. We will be primarily interested in estimating the elements of  $\mathbf{D}$ , the variances of the random effects.

By taking  $\mathbf{u} = \mathbf{0}$ , the model simplifies to the usual probit analysis model. If we set  $\mathbf{X} = \text{diag}\{\mathbf{1}_{m_i}\}$ ,  $i = 1, 2, \dots, G$ ,  $\mathbf{Z} = \text{diag}\{\mathbf{1}_{n_{ij}}\}$ ,  $j = 1, 2, \dots, m_i$ ,  $\boldsymbol{\beta} = \boldsymbol{\mu}$ , this reduces to the Ochi and Prentice (1984) model, with the restriction that negative correlations cannot be modeled. Model (1) has the advantage over the Ochi and Prentice model that it does not require the mean to be constant within levels of the random effect.

This model is closely related to those of Pierce and Sands (1975) and Stiratelli et al. (1984). If  $\boldsymbol{\varepsilon}$  is assumed to have a

\* Charles E. McCulloch is Professor of Biological Statistics, Biometrics Unit and Statistics Center, Cornell University, Ithaca, NY 14853.

logistic distribution instead of a normal distribution, then generalizations of these authors' models are obtained.

Advantages of the probit-normal model (1) over the logit-normal models of Pierce and Sands and Stiratelli et al. are threefold:

1. With a single random effect and only one observation per level of the random effect, it reduces to the usual probit model, except with a different error term for  $\varepsilon$ . The logit-normal models do not reduce to the usual logit models. It is conceptually distasteful for a generalization of a simple model (the logit) not to reduce to the simple model when analyzing a data set appropriate for that model.

2. The marginal mean of  $W_i$  has a simple representation (Zeger et al. 1988):

$$E[W_i] = \Phi(\mathbf{x}_i' \boldsymbol{\beta} (\mathbf{z}_i' \mathbf{D} \mathbf{z}_i + 1)^{-1/2}). \quad (2)$$

3. The EM algorithm (Sec. 3) takes a form nearly identical to the continuous normal linear model.

Point 3 is perhaps the most important, because we exploit it using a Gibbs sampling approach (see Sec. 4) to find ML and REML estimates for arbitrarily complicated models of the form (3) below.

In what follows we will assume the standard analysis of variance (ANOVA) model for variance components estimation that is,

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \sum_{i=1}^r \mathbf{Z}_i \mathbf{u}_i + \varepsilon$$

and

$$\mathbf{u}_i \sim \text{independently } \mathcal{N}_{q_i}(\mathbf{0}, \theta_i \mathbf{I}). \quad (3)$$

These, along with  $W_i = I_{\{Y_i > 0\}}$ , define our basic model.

### 3. MAXIMUM LIKELIHOOD ESTIMATION

In this section we describe estimation of the fixed effects parameters and variance components via the EM algorithm and prediction of the realized values of the random effects. The EM algorithm is used for four reasons: (1) it offers a framework for estimation that is similar to the normal theory case; (2) it automatically constrains iterates to be in the parameter space; (3) it offers a natural extension for REML estimation; and (4) we have found in practice that for simple problems it tends to converge from a wider range of starting values than does a quasi-Newton algorithm (see Sec. 4). To use the EM algorithm we regard the complete data as  $\mathbf{Y}$  and  $\mathbf{u}_i$  ( $i = 1, 2, \dots, r$ ), as is typically done for the continuous normal linear model (Laird 1982). The advantage of the threshold model approach is that we can now appeal to standard results for normally distributed data.

The maximization step is quite simple, as shown by Laird (1982). The ML estimates for the  $\theta_i$  are  $\hat{\theta}_i = \mathbf{u}_i' \mathbf{u}_i / q_i$  and, given estimates of the  $\theta_i$ , the ML estimate of  $\boldsymbol{\beta}$  is  $(\mathbf{X}' \hat{\mathbf{V}}^{-1} \mathbf{X})^{-1} \mathbf{X}' \hat{\mathbf{V}}^{-1} \mathbf{Y}$ , where  $\hat{\mathbf{V}}^{-1}$  is  $\text{var}(\mathbf{Y})$  with  $\theta_i$  replaced by  $\hat{\theta}_i$  and

$$\text{var}(\mathbf{Y}) = \mathbf{I} + \sum_{i=1}^c \theta_i \mathbf{Z}_i \mathbf{Z}_i'.$$

The expectation step is also conceptually simple. We need to calculate  $E[\mathbf{Y} | \mathbf{W}]$  and  $E[\mathbf{u}_i' \mathbf{u}_i | \mathbf{W}]$ . This latter expectation can be calculated by iterated expectation, as was done by Pettit (1986) for censored data:

$$E[\mathbf{u}_i' \mathbf{u}_i | \mathbf{W}] = E[E[\mathbf{u}_i' \mathbf{u}_i | \mathbf{Y}] | \mathbf{W}]. \quad (4)$$

To calculate the inner expectation we can use the usual multivariate normal results:

$$E[\mathbf{u}_i' \mathbf{u}_i | \mathbf{Y}] = \theta_i^2 (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})' \mathbf{V}^{-1} \mathbf{Z}_i \mathbf{Z}_i' \mathbf{V}^{-1} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) + \text{tr}(\theta_i \mathbf{I} - \theta_i^2 \mathbf{Z}_i' \mathbf{V}^{-1} \mathbf{Z}_i). \quad (5)$$

Thus, using (4) we have

$$\begin{aligned} E[\mathbf{u}_i' \mathbf{u}_i | \mathbf{W}] &= \theta_i^2 E[(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})' \mathbf{V}^{-1} \mathbf{Z}_i \mathbf{Z}_i' \mathbf{V}^{-1} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) | \mathbf{W}] \\ &\quad + \text{tr}(\theta_i \mathbf{I} - \theta_i^2 \mathbf{Z}_i' \mathbf{V}^{-1} \mathbf{Z}_i) \\ &= \theta_i^2 \text{tr} E[\mathbf{V}^{-1} \mathbf{Z}_i \mathbf{Z}_i' \mathbf{V}^{-1} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})' | \mathbf{W}] \\ &\quad + \text{tr}(\theta_i \mathbf{I} - \theta_i^2 \mathbf{Z}_i' \mathbf{V}^{-1} \mathbf{Z}_i) \\ &= \theta_i^2 \text{tr} \mathbf{V}^{-1} \mathbf{Z}_i \mathbf{Z}_i' \mathbf{V}^{-1} (\mathbf{V}_{\mathbf{Y}|\mathbf{W}} + (\boldsymbol{\mu}_{\mathbf{Y}|\mathbf{W}} - \mathbf{X}\boldsymbol{\beta})(\boldsymbol{\mu}_{\mathbf{Y}|\mathbf{W}} - \mathbf{X}\boldsymbol{\beta})') \\ &\quad + \text{tr}(\theta_i \mathbf{I} - \theta_i^2 \mathbf{Z}_i' \mathbf{V}^{-1} \mathbf{Z}_i) \\ &= \theta_i^2 \text{tr} \mathbf{V}^{-1} \mathbf{Z}_i \mathbf{Z}_i' \mathbf{V}^{-1} \mathbf{V}_{\mathbf{Y}|\mathbf{W}} + \theta_i^2 (\boldsymbol{\mu}_{\mathbf{Y}|\mathbf{W}} - \mathbf{X}\boldsymbol{\beta})' \mathbf{V}^{-1} \mathbf{Z}_i \mathbf{Z}_i' \mathbf{V}^{-1} \\ &\quad \times (\boldsymbol{\mu}_{\mathbf{Y}|\mathbf{W}} - \mathbf{X}\boldsymbol{\beta}) + \text{tr}(\theta_i \mathbf{I} - \theta_i^2 \mathbf{Z}_i' \mathbf{V}^{-1} \mathbf{Z}_i), \end{aligned}$$

where  $\mathbf{V}_{\mathbf{Y}|\mathbf{W}} = \text{var}(\mathbf{Y} | \mathbf{W})$  and  $\boldsymbol{\mu}_{\mathbf{Y}|\mathbf{W}} = E[\mathbf{Y} | \mathbf{W}]$ .

This shows that the only extra computations needed for ML estimation for discrete data are those of  $\mathbf{V}_{\mathbf{Y}|\mathbf{W}}$  and  $\boldsymbol{\mu}_{\mathbf{Y}|\mathbf{W}}$ . By demonstrating that only the conditional mean and variance of  $\mathbf{Y}$  are needed, the EM algorithm offers a framework for relatively unrestricted computation of complicated mixed models for binary data. In Section 4, using both numerical integration and Gibbs sampling approaches, we show that the computations are feasible in practice.

We are now prepared to make a formal statement of the EM algorithm for ML estimation. In the statement of the algorithm, superscripts in parentheses on  $\mathbf{V}$ ,  $\mathbf{V}_{\mathbf{Y}|\mathbf{W}}$ , and  $\boldsymbol{\mu}_{\mathbf{Y}|\mathbf{W}}$  indicate that the current values of the parameters have been substituted.

#### 3.1 EM Algorithm for ML Estimation

0. Obtain starting values  $\boldsymbol{\beta}^{(0)}$  and  $\theta^{(0)}$ . Set  $m = 0$ .

1. (*E* Step) Calculate

$$\begin{aligned} \hat{t}_i^{(m)} &= E[\mathbf{u}_i' \mathbf{u}_i | \mathbf{W}, \boldsymbol{\beta} = \boldsymbol{\beta}^{(m)}, \theta = \theta^{(m)}] \\ &= \theta_i^{(m)2} \text{tr} \mathbf{V}^{(m)-1} \mathbf{Z}_i \mathbf{Z}_i' \mathbf{V}^{(m)-1} \mathbf{V}_{\mathbf{Y}|\mathbf{W}}^{(m)} \\ &\quad + \theta_i^{(m)2} (\boldsymbol{\mu}_{\mathbf{Y}|\mathbf{W}}^{(m)} - \mathbf{X}\boldsymbol{\beta}^{(m)})' \mathbf{V}^{(m)-1} \mathbf{Z}_i \mathbf{Z}_i' \mathbf{V}^{(m)-1} \\ &\quad \times (\boldsymbol{\mu}_{\mathbf{Y}|\mathbf{W}}^{(m)} - \mathbf{X}\boldsymbol{\beta}^{(m)}) \\ &\quad + \text{tr}(\theta_i^{(m)} \mathbf{I} - \theta_i^{(m)2} \mathbf{Z}_i' \mathbf{V}^{(m)-1} \mathbf{Z}_i). \end{aligned}$$

2. (*M* step) Set

$$\theta_i^{(m+1)} = \hat{t}_i^{(m)} / q_i$$

and

$$\beta^{(m+1)} = (\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{V}^{-1}\mu_{Y|\mathbf{W}}^{(m)}.$$

3. If convergence is reached, set  $\hat{\theta} = \theta^{(m+1)}$  and  $\hat{\beta} = \beta^{(m+1)}$ ; otherwise increase  $m$  by 1 and return to Step 1.

This implementation of the EM algorithm is identical to the continuous case, except that  $(\mathbf{Y} - \mathbf{X}\beta)(\mathbf{Y} - \mathbf{X}\beta)'$  and  $\mathbf{Y}$  are replaced by their expected values given  $\mathbf{W}$ . Most of the computational effort is expended in the calculation of  $\mathbf{V}_{Y|\mathbf{W}}$  and  $\mu_{Y|\mathbf{W}}$ . This is discussed in more detail in Section 4.

We are now in a position to give a version of REML estimation. The basic idea behind REML is to maximize a portion of the likelihood that depends only on the variance components and not on the fixed effects. (See Searle, Casella, and McCulloch 1992, secs. 6.6 and 9.2b for further details.) We use the approach of Laird (1982) and obtain REML estimators by treating the fixed effects as random effects whose variance tends to infinity. This approach is motivated by adopting a Bayesian viewpoint and letting the prior information about the fixed effects tend to 0 (variance tends to infinity); see Harville (1974). Using the same device, as for ML estimation, equation (4), we calculate

$$\begin{aligned} E[\mathbf{u}'_i \mathbf{u}_i | \mathbf{W}] &= E[E[\mathbf{u}'_i \mathbf{u}_i | \mathbf{Y}] | \mathbf{W}] \\ &= \theta_i^2 E[\mathbf{Y}'\mathbf{P}\mathbf{Z}_i\mathbf{Z}'_i\mathbf{P}\mathbf{Y} | \mathbf{W}] + \text{tr}(\theta_i \mathbf{I} - \theta_i^2 \mathbf{Z}'_i\mathbf{P}\mathbf{Z}_i) \\ &= \theta_i^2 \text{tr} \mathbf{P}\mathbf{Z}_i\mathbf{Z}'_i\mathbf{P}(\mathbf{V}_{Y|\mathbf{W}} + \mu_{Y|\mathbf{W}}\mu'_{Y|\mathbf{W}}) \\ &\quad + \text{tr}(\theta_i \mathbf{I} - \theta_i^2 \mathbf{Z}'_i\mathbf{P}\mathbf{Z}_i) \\ &= \theta_i^2 \text{tr} \mathbf{P}\mathbf{Z}_i\mathbf{Z}'_i\mathbf{P}\mathbf{V}_{Y|\mathbf{W}} + \theta_i^2 \mu'_{Y|\mathbf{W}}\mathbf{P}\mathbf{Z}_i\mathbf{Z}'_i\mathbf{P}\mu_{Y|\mathbf{W}} \\ &\quad + \text{tr}(\theta_i \mathbf{I} - \theta_i^2 \mathbf{Z}'_i\mathbf{P}\mathbf{Z}_i), \end{aligned}$$

where  $\mathbf{P} = \mathbf{V}^{-1} - \mathbf{V}^{-1}\mathbf{X}(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{V}^{-1}$ . An analog of REML can be defined for discrete data, using an EM algorithm as follows.

### 3.2 EM Algorithm for REML Estimation

0. Obtain starting values  $\theta^{(0)}$ . Set  $m = 0$ .

1. (*E* Step) Calculate

$$\begin{aligned} \hat{t}_i^{(m)} &= E[\mathbf{u}'_i \mathbf{u}_i | \mathbf{W}, \theta = \theta^{(m)}] \\ &= \theta_i^{(m)2} \text{tr} \mathbf{P}^{(m)}\mathbf{Z}_i\mathbf{Z}'_i\mathbf{P}^{(m)}\mathbf{V}_{Y|\mathbf{W}}^{(m)} \\ &\quad + \theta_i^{(m)2} \mu_{Y|\mathbf{W}}^{(m)'}\mathbf{P}^{(m)}\mathbf{Z}_i\mathbf{Z}'_i\mathbf{P}^{(m)}\mu_{Y|\mathbf{W}}^{(m)} \\ &\quad + \text{tr}(\theta_i^{(m)} \mathbf{I} - \theta_i^{(m)2} \mathbf{Z}'_i\mathbf{P}^{(m)}\mathbf{Z}_i). \end{aligned}$$

2. (*M* step) Set

$$\theta_i^{(m+1)} = \hat{t}_i^{(m)} / q_i.$$

3. If convergence is reached, set  $\hat{\theta} = \theta^{(m+1)}$ ; otherwise increase  $m$  by 1 and return to Step 1.

A major difference between ML and REML estimation is that for REML the limiting values of  $\mathbf{V}_{Y|\mathbf{W}}$  and  $\mu_{Y|\mathbf{W}}$  as the variance of the fixed effects tends to infinity must be used.

The prediction of the observed values of the random effects,  $\mathbf{u}_i$ , is often of interest in applied work (Mabry, Beny-

ske, Johnson, and Little 1987). For continuous data the BLUP methodology is often used, giving rise to  $\hat{\mathbf{u}}_i = \hat{\theta}_i^2 \mathbf{Z}'_i \hat{\mathbf{V}}^{-1}(\mathbf{Y} - \mathbf{X}\hat{\beta}) = \hat{\theta}_i^2 \mathbf{Z}'_i \hat{\mathbf{P}}\mathbf{Y}$ , which is an estimate of  $E[\mathbf{u}_i | \mathbf{Y}]$ . The corresponding calculation for discrete data is  $\hat{\mathbf{u}}_i = \hat{\theta}_i^2 \mathbf{Z}'_i \hat{\mathbf{V}}^{-1}(\hat{\mu}_{Y|\mathbf{W}} - \mathbf{X}\hat{\beta}) = \hat{\theta}_i^2 \mathbf{Z}'_i \hat{\mathbf{P}}\hat{\mu}_{Y|\mathbf{W}}$ , which is an estimate of  $E[\mathbf{u}_i | \mathbf{W}]$ . The form of the estimator is the same whether we use ML or REML estimation, although the estimates generally will be different due to different values for the variance components and  $\hat{\mu}_{Y|\mathbf{W}}$ .

As pointed out by Wu (1983), EM is not guaranteed to converge to a global maximum. Our experience has shown that multimodal likelihoods are possible for models such as these; so the best we can hope for in this setting is that EM will converge to a local maximum. Unfortunately the regularity conditions of Wu (1983) do not apply; a realistic compactification of the parameter space by including infinite variance components leads to identifiability problems. Truncation of the parameter space to exclude extremely large values would allow the regularity conditions to be met. Then, because  $Q((\theta^*, \beta^*) | (\theta, \beta)) = E[\log f(\mathbf{Y} | (\theta^*, \beta^*)) | \mathbf{W}, (\theta, \beta)]$  is continuous in both  $(\theta^*, \beta^*)$  and  $(\theta, \beta)$ , theorem 2 of Wu (1983) applies and EM is guaranteed to converge to a stationary point. For any particular data set a local maximum would need to be verified by numerically calculating the second derivative matrix via numerical integration or techniques like that of Meng and Rubin (1991).

## 4. EXAMPLES

We applied the methods derived in Section 3 to the data analyzed by Ochi and Prentice (1984), from Weil (1970), and to the salamander data from McCullagh and Nelder (1989, sec. 14.5).

### 4.1 The Weil Data

The Weil data set has a treatment and control group and a single nested random effect. The response is survival of rats, and the random effect is litter. The model would be

$$Y_{ij} = \mu_i + u_{ij} + \varepsilon_{ijk}$$

$$W_{ij} = I_{\{Y_{ij} > 0\}},$$

where  $i$  indexes treatment/control,  $j$  indexes litter, and  $k$  indexes rat within litter, so  $\mu_i$  is the treatment mean on the latent scale and the  $u_{ij}$  are the random litter effects. To find ML estimates for the Weil data set, it is most efficient to numerically evaluate integrals of the form

$$\int_{-\infty}^{\infty} \alpha \Psi(\alpha; n, s, \mu, \sigma) \phi(\alpha) d\alpha,$$

where

$$\Psi(\alpha; n, s, \mu, \sigma)$$

$$= \frac{\Phi(\mu + \sigma\alpha)^s [1 - \Phi(\mu + \sigma\alpha)]^{n-s}}{\int_{-\infty}^{\infty} \Phi(\mu + \sigma x)^s [1 - \Phi(\mu + \sigma x)]^{n-s} \phi(x) dx}$$

and  $\Phi(\cdot)$  and  $\phi(\cdot)$  are the standard normal cdf and pdf. As pointed out by Pettit (1986), integrals such as these have been well studied and are relatively easy to evaluate numerically.



ically. We used Hermite quadrature with 20 evaluation points (Abramowitz and Stegun 1964, table 25.10,  $n = 10$ ). We noticed none of the accuracy problems reported by Ochi and Prentice (1984) and in fact were able to reproduce the true values in their table 1 exactly.

We used the matrix language GAUSS (Aptech Systems 1990) on an IBM PC and fit each group separately. (Previous analyses have shown unequal values of the variance in the two groups.) For the treated group ( $i = 2$ ) the algorithm converged in 28 iterations and under 1 minute. For the control group ( $i = 1$ ) the EM algorithm required 207 iterations and about  $1\frac{1}{2}$  minutes. The ML estimates and their standard errors (calculated from the observed information matrix) were  $\hat{\mu}_1 = 1.306$  (standard error .169),  $\text{stddev}(u_{1j}) \equiv \hat{\sigma}_1 = .240$  (standard error .301),  $\hat{\mu}_2 = .946$  (standard error .319), and  $\text{stddev}(u_{2j}) \equiv \hat{\sigma}_2 = 1.023$  (standard error .291). These estimates agree substantially with those of Ochi and Prentice (1984); slight differences are to be expected, because Ochi and Prentice used the approximation due to Mendell and Elston (1974). For example, in group 2 Ochi and Prentice obtained  $\hat{\gamma}_2 \equiv \hat{\mu}_2 / \sqrt{1 + \hat{\sigma}_2^2} = .651$ , whereas our estimates give  $\hat{\gamma}_2 = .661$ . The large number of iterations required by the EM algorithm for the control group is typical of problems for which the estimates lie near the boundary of the parameter space. When the likelihood can be evaluated numerically, as in this example, it is straightforward to conduct likelihood ratio tests and to evaluate derivatives of the likelihood function for calculating standard errors.

We also fitted this data set using a quasi-Newton algorithm (Aptech Systems 1990, p. 207). Convergence was achieved to essentially the same parameter values, and each group was fitted in less than 1 minute. A small amount of experimentation with the starting values showed that the EM algorithm converged from a wider range of starting values than did the quasi-Newton algorithm.

## 4.2 The Gibbs Sampler and the Salamander Data

In a design with a more complicated random effects structure (e.g., crossed effects), the computations become too burdensome for direct numerical calculation; for example, the algorithm of Leppard and Tallis (1989) works only for small dimensions. To illustrate the flexibility of the framework of Section 3, we use a Gibbs sampling approach (Gelfand and Smith 1990) to calculate  $E[\mathbf{Y}|\mathbf{W}]$  and  $\text{var}(\mathbf{Y}|\mathbf{W})$ . Tanner (1991) suggested a similar Monte Carlo EM algorithm. By using the Gibbs sampler, arbitrarily complicated designs can be accommodated easily. We apply this approach to the salamander data of McCullagh and Nelder (1989, sec. 14.5), which has two crossed random effects and four fixed effects.

We now outline the use of the Gibbs sampler. It rests on a result of Robert (1992) for sampling from a truncated multivariate normal and is similar to the treatment of Gelfand, Smith, and Lee (1992). The basic idea is that fast acceptance-rejection methods exist (see, for example, Marsaglia 1964) for sampling from a truncated univariate normal. By cycling through the conditional distributions of  $Y_i | Y_j, j \neq i$ , we only ever need to simulate truncated univariate normals. Here is

an outline of how the Gibbs sampler is used to generate a sample of  $\mathbf{Y}$ 's from the conditional distribution of  $\mathbf{Y}|\mathbf{W}$ :

1. For each  $i$  calculate

$$\sigma_{i|(i)}^2 = \text{var}(Y_i | Y_j, j \neq i)$$

and

$$\beta_{i|(i)} = \text{cov}(Y_i, Y_{(i)}),$$

where  $Y_{(i)} = (Y_1, Y_2, \dots, Y_{i-1}, Y_{i+1}, \dots, Y_n)'$ .

2. For each  $i$  calculate

$$\mu_{i|(i)} = E[Y_i | Y_j, j \neq i] = \mathbf{x}_i \boldsymbol{\beta} + \beta'_{i|(i)} (\mathbf{Y}_{(i)} - \mathbf{X}_{(i)} \boldsymbol{\beta}),$$

where  $\mathbf{X}_{(i)} = \mathbf{X}$  with row  $i$  deleted and  $\mathbf{x}_i$  is the  $i$ th row of  $\mathbf{X}$ .

3. Simulate  $Y_i$  from a truncated normal distribution with mean  $\mu_{i|(i)}$  and standard deviation  $\sigma_{i|(i)}$ . If  $W_i = 1$ , simulate  $Y_i$  truncated above 0. If  $W_i = 0$ , simulate  $Y_i$  to be truncated below 0.

Repeat Steps 2 and 3 a large number of times,  $k$ , to obtain  $\mathbf{Y}^{(1)}, \mathbf{Y}^{(2)}, \dots, \mathbf{Y}^{(k)}$ . Discard a suitable number of the  $\mathbf{Y}^{(j)}$  from the beginning of the sequence (the burn-in period), and after that use every  $m$ th one to estimate  $E[\mathbf{Y}|\mathbf{W}]$  and  $\text{var}(\mathbf{Y}|\mathbf{W})$ . Because of the iterative nature of the EM algorithm and the desire to take as few Gibbs samples as possible (especially at the beginning of EM), we settled on a burn-in period of  $i$ , skipped integer  $(i/10) + 1$  samples, and used  $i + 1$  replications, where  $i$  is the iteration in the EM algorithm. These numbers are small in relation to those recommended in the literature, but we noticed no problems. We tried larger values with no improvement. For simpler cases where quasi-Newton estimation was possible, we compared EM-Gibbs and quasi-Newton for a number of simulated data sets and had success with the smaller number of Gibbs samples in each case.

In the Gibbs sampler most of the computational effort is expanded in repeating Steps 2 and 3 a sufficiently large number of times. Thus complicated random effects structures have little impact on the computational time because they affect only Step 1.

The salamander data consist of three experiments, each with  $n = 120$  matings.  $W_i = 1$  if the  $i$ th mating is successful and 0 otherwise. There were 20 males and 20 females, 10 of each of two species, and four types of crosses in the matings: species R female-species R male, species R female-species W male, species W female-species R male, and species W female-species W male. For each experiment the model is

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}_f \mathbf{U}_f + \mathbf{Z}_m \mathbf{U}_m + \boldsymbol{\varepsilon},$$

$$W_i = I_{\{Y_i > 0\}},$$

$$\mathbf{U}_f \sim \mathcal{N}_{20}(0, \theta_f \mathbf{I}), \text{ the female effects,}$$

$$\mathbf{U}_m \sim \mathcal{N}_{20}(0, \theta_m \mathbf{I}), \text{ the male effects,}$$

$$\boldsymbol{\varepsilon} \sim \mathcal{N}_{120}(0, \mathbf{I}),$$

$$\mathbf{X} = \text{indicator matrix for the type of cross,}$$

$$\mathbf{Z}_f = \text{indicator matrix for the females,}$$

$$\mathbf{Z}_m = \text{indicator matrix for the males,}$$

and

$$\boldsymbol{\beta} = (\beta_{R/R}, \beta_{R/W}, \beta_{W/R}, \beta_{W/W})' \text{ effect for type of cross.}$$

$\mathbf{U}_f$  and  $\mathbf{U}_m$  represent the consistent effect that individual

males and females have across matings on the latent variable  $Y$ , which governs mating success. They are the random effects assumed to be iid with variances given by  $\theta_f$  and  $\theta_m$ . Figure 1 shows the convergence of the parameter estimates for experiment 1. The final estimates were  $\hat{\beta}_{R/R} = .819$ ,  $\hat{\beta}_{R/W} = .538$ ,  $\hat{\beta}_{W/R} = -.978$ ,  $\hat{\beta}_{W/W} = .707$ ,  $\hat{\theta}_f = .600$ , and  $\hat{\theta}_m = .067$ . These are relatively similar to the estimates that Karim and Zeger (1992) obtained in a Bayesian analysis using the Gibbs sampler and a logit-normal model. Table 1 shows the Bayesian and ML estimates of the variance components for the three experiments. When the likelihood is not directly evaluated, as in this example, it is much more complicated to calculate standard errors. Techniques based directly on EM (see, for example, Meng and Rubin 1991) are necessary.

The estimates of the marginal probabilities [see (2)] are almost exactly equal to the observed proportions:

Cross	Estimated marginal proportion $\Phi(\hat{\beta}/(\hat{\theta}_f + \hat{\theta}_m + 1)^{1/2})$	Observed proportion
R/R	$\Phi(.819/ (.6 + .067 + 1)^{1/2}) = .737$	$22/30 = .733$
R/W	.661	$20/30 = .667$
W/R	.224	$7/30 = .233$
W/W	.708	$21/30 = .7$

Although this approach is computationally intensive, it is not prohibitive. On a fast (33 mHz 486) IBM PC-compatible computer using the GAUSS language (Aptech Systems 1990), 50 iterations were completed in 90 minutes and 80 iterations were completed in 250 minutes. (Later iterations do more Gibbs sampling.) These times could undoubtedly be improved by more efficient programming and computational techniques.

This application of the Gibbs sampler is unusual in that it is used to solve directly for ML estimates rather than using a Bayesian framework. It would seem to be of broad utility for models that contain a latent multivariate normal component.

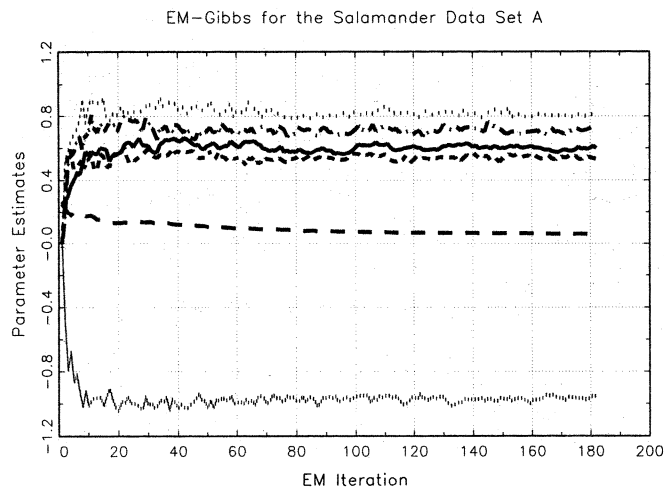


Figure 1. EM Iterations Using a Gibbs Sampler for the Salamander Data Set 1. —  $\theta_f$ ; - -  $\theta_m$ ; ...  $\beta_{R/R}$ ; - ·  $\beta_{W/R}$ ; · ·  $\beta_{R/W}$ ; and · · ·  $\beta_{W/W}$ .

Table 1. Comparison of ML and Bayes Estimates of the Female ( $\theta_f$ ) and Male ( $\theta_m$ ) Variance Components from the Three Salamander Data Sets (McCullagh and Nelder 1989)

Estimate	Variance:	Data set					
		1		2		3	
		$\theta_f$	$\theta_m$	$\theta_f$	$\theta_m$	$\theta_f$	$\theta_m$
ML		.60	.06	.49	.45	.10	.44
Bayes		.81	.05	1.03	.49	.11	1.00

NOTE: The Bayesian estimates are taken from Karim and Zeger (1992, table 4) and are divided by  $(\pi/\sqrt{3})/(15/16)^2$  for comparability (Johnson and Kotz 1970, p. 6).

5. CONCLUSIONS

We have developed a framework for ML and REML estimation of variance components from binary data using the EM algorithm. This framework is very similar to the EM algorithm for the continuous normal linear model. For simple settings the ML computations can be performed by numerical integration. For more complicated problems this framework can be used with Gibbs sampling approach to calculate ML and REML estimates. This has not been previously possible in designs with complicated (e.g., crossed) random effects.

[Received April 1991. Revised April 1993.]

REFERENCES

Abramowitz, M., and Stegun, I. (1964), *Handbook of Mathematical Functions*, Washington, DC: National Bureau of Standards.

Anderson, D. A., and Aitken, M. (1985), "Variance Components Models with Binary Response: Interviewer Variability," *Journal of the Royal Statistical Society, Ser. B*, 47, 203-210.

Aptech Systems (1990), *GAUSS 2.1 Users' Manual*, Kent, WA: Author.

Crowder, M. J. (1978), "Beta-Binomial ANOVA for Proportions," *Applied Statistics*, 27, 34-37.

Dempster, A. P., Laird, N. M., and Rubin, D. B. (1977), "Maximum Likelihood From Incomplete Observations," *Journal of the Royal Statistical Society, Ser. B*, 39, 1-38.

Gelfand, A. E., and Smith, A. F. M. (1990), "Sampling-Based Approaches to Calculating Marginal Densities," *Journal of the American Statistical Association*, 85, 105-114.

Gelfand, A. E., Smith, A. F. M., and Lee, T.-M. (1992), "Bayesian Analysis of Constrained Parameter and Truncated Data Problems Using Gibbs Sampling," *Journal of the American Statistical Association*, 87, 523-532.

Gianola, D. (1980), "Genetic Evaluation of Animals for Traits With Categorical Responses," *Journal of Animal Science*, 51, 1272-1276.

Gilmour, A. R., Anderson, R. D., and Rae, A. L. (1985), "The Analysis of Binomial Data by a Generalized Linear Mixed Model," *Biometrika*, 72, 593-599.

Harville, D. A. (1974), "Bayesian Inference for Variance Components Using Only Error Contrasts," *Biometrika*, 61, 383-385.

Harville, D. A., and Mee, R. W. (1984), "A Mixed Model Procedure for Analyzing Ordered Categorical Data," *Biometrics*, 40, 393-408.

Im, S., and Gianola, D. (1988), "Mixed Models for Binomial Data With an Application to Lamb Mortality," *Applied Statistics*, 37, 196-204.

Karim, M. R., and Zeger, S. L. (1992), "Generalized Linear Models With Random Effects; Salamander Mating Revisited," *Biometrics*, 48, 681-694.

Laird, N. L. (1982), "Computation of Variance Components using the EM Algorithm," *Journal of Statistical Computation and Simulation*, 14, 295-303.

Leppard, P., and Tallis, G. M. (1989), "Evaluation of the Mean and Covariance of the Truncated Multinormal Distribution," *Applied Statistics*, 38, 543-553.

Liang, K.-Y., and Zeger, S. L. (1986), "Longitudinal Data Analysis Using Generalized Linear Models," *Biometrika*, 73 13-22.

- Liang, K.-Y., Zeger, S. L., and Qaqish, B. (1992), "Multivariate Regression Analyses for Categorical Data," *Journal of the Royal Statistical Society, Ser. B*, 54, 3–40.
- Mabry, J. W., Benyskek, L. L., Johnson, M. H., and Little, D. E. (1987), "A Comparison of Methods for Ranking Boars From Different Central Test Stations," *Journal of Animal Science*, 65, 56–62.
- Manski, C. F., and McFadden, D. (1981), "Structural Analysis of Discrete Data With Econometric Applications," Cambridge, MA: MIT Press.
- Marsaglia, G. (1964), "Generating a Variable From the Tail of the Normal Distribution," *Technometrics*, 6, 101–102.
- Mendell, N. R., and Elston, R. C. (1974), "Multifactorial Qualitative Traits: Genetic Analysis and Prediction of Recurrence Risks," *Biometrics*, 30, 41–57.
- Meng, X.-L., and Rubin, D. B. (1991), "Using EM to Obtain Asymptotic Variance–Covariance Matrices: The SEM Algorithm," *Journal of the American Statistical Association*, 86, 899–909.
- Ochi, Y., and Prentice, R. L. (1984), "Likelihood Inference in a Correlated Probit Regression Model," *Biometrika*, 71, 531–543.
- Pettit, A. N. (1986), "Censored Observations, Repeated Measures, and Mixed Effects Models: An Approach Using the EM Algorithm and Normal Errors," *Biometrika*, 73, 634–643.
- Pierce, D. A., and Sands, B. R. (1975), "Extra-Bernoulli Variation in Binary Data," Technical Report 46, Oregon State University, Dept. of Statistics.
- Preisler, H. K. (1989), "Analysis of a Toxicological Experiment Using a Generalized Linear Model With Nested Random Effects," *International Statistical Review*, 57, 145–159.
- Prentice, R. L. (1988), "Correlated Binary Regression With Covariates Specific to Each Binary Observations," *Biometrics*, 4, 1033–1048.
- Quaas, D. L., and Van Vleck, L. D. (1980), "Categorical Trait Sire Evaluation by Best Linear Unbiased Prediction of Future Progeny Category Frequency," *Biometrics*, 36, 117–122.
- Robert, C. R. (1992), "Simulation of Truncated Normal Variables," Technical Report No. 161, LSTA, University of Paris 6.
- Searle, S. R., Casella, G., and McCulloch, C. E. (1992), *Variance Components*, New York: John Wiley.
- Stratelli, R., Laird, N. M., and Ware, J. H. (1984), "Random-Effects Models for Serial Observations With Binary Response," *Biometrics*, 40, 961–971.
- Stram, D. O., Wei, L. J., and Ware, J. H. (1988), "Analysis of Repeated Ordered Categorical Outcomes With Possibly Missing Observations and Time-Dependent Covariates," *Journal of the American Statistical Association*, 83, 631–637.
- Tanner, M. A. (1991), *Tools for Statistical Inference: Observed Data and Data Augmentation Methods*, Berlin: Springer-Verlag.
- Tierney, L. (1991), "Markov Chains for Exploring Posterior Distributions," Technical Report No. 560, University of Minnesota, School of Statistics.
- Weil, C. S. (1970), "Selection of the Valid Number of Sampling Units and Consideration of Their Combination in Toxicological Studies Involving Reproduction, Teratogenesis, or Carcinogenesis," *Food and Cosmetic Toxicology*, 8, 177–182.
- Williams, D. A. (1975), "The Analysis of Binary Responses From Toxicological Experiments Involving Reproduction and Teratogenicity," *Biometrics*, 31, 949–952.
- Wong, G. Y., and Mason, W. M. (1985), "The Hierarchical Logistic Regression Model for Multilevel Analysis," *Journal of the American Statistical Association*, 80, 513–524.
- Wu, C.-F. (1983), "On the Convergence Properties of the EM Algorithm," *The Annals of Statistics*, 11, 95–103.
- Zeger, S. L., and Liang, K.-Y. (1986), "Longitudinal Data Analysis for Discrete and Continuous Outcomes," *Biometrics*, 42, 121–130.
- Zeger, S. L., Liang, K.-Y., and Albert, P. S. (1988), "Models for Longitudinal Data: A Generalized Estimating Equation Approach," *Biometrics*, 44, 1049–1060.