

Implementations of the Monte Carlo EM algorithm

Richard A. Levine* George Casella†
University of California, Davis Cornell University

December 13, 1998

Abstract

The Monte Carlo EM (MCEM) algorithm is a modification of the EM algorithm where the expectation in the E-step is computed numerically through Monte Carlo simulations. The most flexible and generally applicable approach to obtaining a Monte Carlo sample in each iteration of an MCEM algorithm is through Markov chain Monte Carlo (MCMC) routines such as the Gibbs and Metropolis-Hastings samplers. While MCMC estimation presents a tractable solution to problems where the E-step is not available in closed form, two issues arise when implementing this MCEM routine: 1) how do we minimize the computational cost in obtaining an MCMC sample? and 2) how do we choose the Monte Carlo sample size? We address the first question through an application of importance sampling whereby samples drawn during previous EM iterations are recycled, rather than running an MCMC sampler each MCEM iteration. The second question is addressed through an application of regenerative simulation. We obtain approximate independent and identical samples by subsampling the generated MCMC sample during different renewal periods. Standard Central Limit Theorems may thus be used to gauge Monte Carlo error. In particular, we apply an automated rule for increasing the Monte Carlo sample size when the Monte Carlo error overwhelms the EM estimate at any given iteration. We illustrate our MCEM algorithm through analyses of three data sets fit by either generalized linear mixed models or survival models in the presence of censoring. As a part of these applications, we demonstrate the improvement in computational cost and efficiency of our routine over alternative MCEM strategies.

Keywords: Gibbs sampler, Metropolis-Hastings algorithm, importance sampling, Markov chain Monte Carlo, regenerative simulation, renewal theory, generalized linear mixed models, survival models.

*This research was done while Prof. Levine was visiting the Department of Biometrics and the Department of Statistical Science, Cornell University, supported by National Science Foundation Grant DMS-9625440. Address: One Shields Ave., Division of Statistics, University of California, Davis, 95616, USA. Email: rlevine@ucdavis.edu

†Supported by National Science Foundation Grant DMS-9625440. Email: gc15@cornell.edu
This paper is BU 1430-M of the Department of Biometrics technical report series, Cornell University.

1 Introduction

The EM algorithm provides a tool for obtaining maximum likelihood estimates under models that yield analytically formidable likelihood equations. The EM algorithm is an iterative routine requiring two primary calculations each iteration: Computation of a particular conditional Expectation of the log-likelihood (E-step) and Maximization of this expectation over the relevant parameters (M-step). The Monte Carlo EM (MCEM), introduced by Wei and Tanner (1990) is a modification of the EM algorithm where the expectation in the E-step is computed numerically through Monte Carlo simulations. While the Monte Carlo estimate presents a tractable solution to problems where the E-step is not available in closed form, we must account for the additional Monte Carlo (MC) error inherent in the approach and try to minimize the increased computational cost in obtaining the MC sample.

Booth and Hobert (1999) present a method for gauging Monte Carlo error in the MCEM algorithm through an automated routine to increase the number of MC samples as the algorithm converges. Their routine relies on MC estimation in the E-step via either independent samples from the distribution of interest or importance weighted random samples from a candidate distribution “close” to that distribution. Independent samples allow for computationally inexpensive and straightforward assessment of Monte Carlo error through the Central Limit Theorem. However, such sampling routines are often not available due either to the complexity of the target distribution or lack of a reasonable importance density (under which a large number of samples are needed to attain a good estimator).

Alternatively, McCulloch (1994, 1997) suggests obtaining a sample via Markov chain Monte Carlo techniques, in particular the Gibbs sampler and Metropolis-Hastings algorithm. Though the random variates are dependent in such a scenario, the E-step estimator is still unbiased and approaches the true value as the sample size increases. Furthermore, MCMC techniques are applicable to a wider range of distributions than approaches based on independent samples.

The greater flexibility introduced by MCMC sampling, however, is countered by greater computational cost and difficulties in assessing Monte Carlo error. Application of, say, a

Metropolis-Hastings algorithm each iteration of the EM algorithm to evaluate the E-step is significantly more expensive than comparable independent samplers, if such are available (Booth and Hobert, 1999). In situations where MCMC is the only alternative, the computational cost could be quite restrictive in that estimation of the E-step to a satisfactory number of significant digits requires an inordinate amount of time (McCulloch, 1997). Furthermore, validating Central Limit Theorems under Markov chain sampling could be a mathematically and computationally complex task.

In this paper we overcome the computational burdens of MCMC sampling in the MCEM algorithm by applying results from importance sampling and renewal theory. First, we show how to recycle variates generated from previous iterations of the EM algorithm through importance weighting. Hence, the MCEM algorithm generates only one sample for evaluation of subsequent E-steps, rather than obtaining an MCMC sample for each iteration of the EM algorithm. Second, we construct a Central Limit Theorem based on renewals of the Markov chain. Markov chains often exhibit *regeneration times* during which the chain essentially restarts. Excursions between these renewal times are hence independent and identically distributed (see Mykland, Tierney, and Yu, 1995). Here, we implement a method of Robert et al. (1998) to identify samples in each of these excursions. Consequently, we are able to collect an iid sample, being a subsample of our original MCMC sample. We may then apply a Central Limit Theorem based on these iid observations as in Booth and Hobert (1999) to assess MC error.

The paper unfolds as follows. In Section 2 we formalize the MCEM routine and present the importance sampling modification. We also discuss the application of regenerative simulation for assessing MC error. We conclude the section with an outline of our MCEM algorithm. In Section 3 we apply our MCEM to fitting three models: the logit-normal model of McCulloch (1997) fit to simulated data, the lifetime distribution model of Fan et al. (1998) fit to the Australian twin data of Duffy et al. (1990), and the probit-normal model of McCulloch (1994) fit to the salamander data of McCullagh and Nelder (1989). As part of these applications, we compare our model with the performance of other MCEM routines utilizing MCMC for evaluating the E-step.

2 Monte Carlo EM

Let $\mathbf{y} = (y_1, \dots, y_n)'$ denote observed data with distribution $f(\mathbf{y} | \Psi)$ characterized by the s -vector of parameters Ψ . Here $'$ denotes vector transpose. The EM algorithm is an iterative routine for computing the maximum likelihood estimate of Ψ , denoted $\hat{\Psi}$. The driving force behind the EM algorithm is that the MLE is simpler to compute on the data \mathbf{y} augmented by a set of latent variables $\mathbf{u} = (u_1, \dots, u_q)'$. For example, we may think of the latent variables as unknown random effects in a generalized linear model or missing observations from individuals with censored survival times.

The EM algorithm thus works on the augmented log-likelihood $\ln f(\mathbf{y}, \mathbf{u} | \Psi)$ to obtain the MLE of Ψ over the distribution $f(\mathbf{y} | \Psi)$ where it is assumed that $f(\mathbf{y} | \Psi) = \int f(\mathbf{y}, \mathbf{u} | \Psi) d\mathbf{u}$. More specifically, the EM algorithm iterates between a calculation of the expected complete-data likelihood

$$Q(\Psi | \hat{\Psi}^{(r)}) = E_{\hat{\Psi}^{(r)}} \{ \ln f(\mathbf{y}, \mathbf{u} | \Psi) | \mathbf{y} \} \quad (1)$$

and the maximization of $Q(\Psi | \hat{\Psi}^{(r)})$ over Ψ , where the maximum value of Ψ is denoted by $\hat{\Psi}^{(r+1)}$ and $\hat{\Psi}^{(r)}$ denotes the maximum of Ψ at the r th iteration. Wu (1984) shows that under regularity conditions, the sequence of values $\{\hat{\Psi}^{(r)}\}$ converges to the MLE $\hat{\Psi}$.

In situations where the E-step is analytically troublesome, we may estimate the quantity (1) from Monte Carlo simulations. Note the expectation in (1) is over the latent variable \mathbf{u} . In particular,

$$E_{\hat{\Psi}^{(r)}} \{ \ln f(\mathbf{y}, \mathbf{u} | \Psi) | \mathbf{y} \} = \int \ln f(\mathbf{y}, \mathbf{u} | \Psi) g(\mathbf{u} | \mathbf{y}, \hat{\Psi}^{(r)}) d\mathbf{u}$$

where $g(\mathbf{u} | \mathbf{y}, \Psi)$ is the conditional distribution of the latent variables given the observed data and Ψ . If we obtain a sample $\mathbf{u}_1^{(r)}, \dots, \mathbf{u}_m^{(r)}$ from the distribution $g(\mathbf{u} | \mathbf{y}, \hat{\Psi}^{(r)})$, this expectation may be estimated by the Monte Carlo sum

$$Q_m(\Psi | \hat{\Psi}^{(r)}) = \frac{1}{m} \sum_{t=1}^m \ln f(\mathbf{y}, \mathbf{u}_t^{(r)} | \Psi) \quad (2)$$

where the subscript m denotes the dependence of this estimator on the MC sample size. By the law of large numbers, the estimator in (2) converges to the theoretical expectation in

(1). The EM algorithm can thus be modified into an *MCEM* whereby the E-step is replaced by the estimated quantity from (2). The M-step maximizes then the sum (2) over Ψ . See Chan and Ledolter (1995) for more details.

In this paper we focus solely on the problem of an intractable E-step. While the M-step may also require sophisticated numerical routines, our applications in Section 3 yield straightforward maximizations in the M-step. Implementation of the MCEM algorithm then presents two important issues for us: how do we obtain a random sample from the distribution $g(\mathbf{u} \mid \mathbf{y}, \Psi)$ and how do we choose m ? We will discuss these issues in the following two subsections. The third subsection summarizes our MCEM algorithm.

2.1 Importance sampling

We consider the situation where the sample of latent variables $\mathbf{u}_1, \dots, \mathbf{u}_m$ in the E-step is obtained from a Markov chain Monte Carlo routine such as the Gibbs sampler or Metropolis-Hastings algorithm with stationary distribution $g(\mathbf{u} \mid \mathbf{y}, \Psi)$. As mentioned in Section 1, drawing an MCMC sample each iteration of the EM algorithm could be prohibitively costly particularly for large m .

The computational expense of the MCMC based MCEM algorithm can be substantially improved through an application of importance sampling. We initialize the algorithm with a sample $\mathbf{u}_1, \dots, \mathbf{u}_m$ from the distribution $g(\mathbf{u} \mid \mathbf{y}, \Psi^{(0)})$ where $\Psi^{(0)}$ is the initial value of the parameter Ψ at the start of the EM algorithm. At each iteration r , rather than obtaining a new sample from $g(\mathbf{u} \mid \mathbf{y}, \hat{\Psi}^{(r)})$ with the most recent iterate $\hat{\Psi}^{(r)}$, we can importance weight the original sample through the updated distribution $g(\mathbf{u} \mid \mathbf{y}, \hat{\Psi}^{(r)})$. That is,

$$Q_m(\Psi \mid \hat{\Psi}^{(r)}) = \sum_{t=1}^m w_t \ln f(\mathbf{y}, \mathbf{u}_t \mid \Psi) \Big/ \sum_{t=1}^m w_t, \quad (3)$$

where the original sample is corrected for the new information we have at iteration r through the weights

$$w_t = \frac{g(\mathbf{u}_t \mid \mathbf{y}, \hat{\Psi}^{(r)})}{g(\mathbf{u}_t \mid \mathbf{y}, \Psi^{(0)})}.$$

(See Robert and Casella, 1999, Chapter 3 for a discussion of importance sampling.)

The cost in obtaining the weights w_t is less than obtaining a new sample. The reason for this expense saving is that the weights are not dependent on the unknown likelihood $L(\Psi | \mathbf{y})$. Note that

$$w_t = \frac{L(\hat{\Psi}^{(r)} | \mathbf{u}_t, \mathbf{y}) / L(\hat{\Psi}^{(r)} | \mathbf{y})}{L(\Psi^{(0)} | \mathbf{u}_t, \mathbf{y}) / L(\Psi^{(0)} | \mathbf{y})}.$$

Hence the likelihood $L(\Psi | \mathbf{y})$ cancels in the formulation (3) so that

$$Q_m(\Psi | \hat{\Psi}^{(r)}) = \frac{\sum_{t=1}^m w'_t \ln f(\mathbf{y}, \mathbf{u}_t | \Psi)}{\sum_{t=1}^m w'_t}$$

where

$$w'_t = \frac{L(\hat{\Psi}^{(r)} | \mathbf{u}_t, \mathbf{y})}{L(\Psi^{(0)} | \mathbf{u}_t, \mathbf{y})}.$$

The importance sampling estimator may, alternatively, be calculated as

$$Q_m(\Psi | \hat{\Psi}^{(r)}) = \frac{1}{m} \sum_{t=1}^m w_t \ln f(\mathbf{y}, \mathbf{u}_t | \Psi) \quad (4)$$

where $\sum_{t=1}^m w_t$ is replaced by m in the denominator. This choice will not affect the EM algorithm since the unknown normalizing constant

$$\frac{L(\Psi^{(0)} | \mathbf{y})}{L(\hat{\Psi}^{(r)} | \mathbf{y})}$$

depends only on the known values $\Psi^{(0)}$ and $\hat{\Psi}^{(r)}$ and not the unknown value of Ψ . This constant does not then come into play in the maximization step. However, we choose the estimator (3) in order to avoid calculation of this normalizing constant in the routine described in the next subsection for choosing m . The estimator (3) may be further rationalized by the fact that $E[w] = 1$ where the expectation is taken with respect to $g(\mathbf{u} | \mathbf{y}, \Psi^{(0)})$. Furthermore, the estimator (3) often has smaller mean squared error than (4); see Liu (1996) and Casella and Robert (1996).

Importance sampling estimators are not without drawbacks. Most notably, if the importance density $g(\mathbf{u}_t | \mathbf{y}, \hat{\Psi}^{(r)})$ is not close enough to $g(\mathbf{u}_t | \mathbf{y}, \Psi^{(0)})$, the weights w_t will vary widely giving many samples little weight and allowing for a few variates to be overinfluential. Consequently, the estimator (3) will be imprecise (Robert and Casella, 1999, Section 3.3.2; also see Geyer, 1991). In our setting, if the initial values $\Psi^{(0)}$ are poor, the importance

sampling estimator will take a long time to converge. We alleviate this problem by initiating a burn-in whereby for the first few iterations, a new sample is obtained from $g(\mathbf{u}_t | \mathbf{y}, \hat{\Psi}^{(r)})$ rather than importance weighting. McCulloch (1997) shows that such an algorithm reaches the neighborhood of the MLE quickly. Following the burn-in, the target density should be close enough to the distribution $g(\mathbf{u}_t | \mathbf{y}, \hat{\Psi})$ based on the MLE $\hat{\Psi}$ to ensure well-behaved weights; i.e. weights with small variance.

The burn-in idea is further rationalized by the choice of m . Tanner (1993; Section 4.5) suggests increasing m as the EM algorithm progresses. Hence, the first few iterations are implemented with a relatively small MC sample size, perhaps as small as $m = 10$ as we will see in Section 3. Generating a sample of size ten each iteration is a computationally inexpensive task. In later iterations, when m is of the order of a thousand to tens of thousands, the cost of generating a sample is too much to perform each iteration. Importance sampling thus becomes particularly crucial for feasible running of the MCEM under MCMC sampling at these later stages.

2.2 Central limit theorem

The choice of Monte Carlo sample size m is equally important to the sampling method chosen for the MCEM algorithm. We do not want to start with a large value of m when the iterates $\hat{\Psi}^{(r)}$ are far from the MLE $\hat{\Psi}$. The tradeoff between improving accuracy and the computational cost in obtaining more samples favors starting the algorithm with small MC samples. However, as the EM algorithm progresses, we may wish to increase m as the approximation $\hat{\Psi}^{(r)}$ approaches the true MLE (Tanner, 1993; Section 4.5).

Up until recently, the choice of m was purely ad hoc; that is, one would arbitrarily increase m at predetermined iterations of the EM algorithm. Booth and Hobert (1999) present an automated routine where the MC sample size is chosen through considerations of the Monte Carlo error inherent in the Monte Carlo sum (2). In particular, since the random variates $\mathbf{u}_1, \dots, \mathbf{u}_m$ are generated independently, the Central Limit Theorem (CLT) ensures that, conditional on the iterate $\hat{\Psi}^{(r)}$, the current iterate $\hat{\Psi}^{(r+1)}$ is approximately

normally distributed. Thus we may estimate Monte Carlo error via the normal distribution. In particular, if the past value $\hat{\Psi}^{(r)}$ lies in an approximate confidence interval about $\hat{\Psi}^{(r+1)}$, then the Monte Carlo error is said to swamp the EM step and the number of simulations, m , is increased (see Booth and Hobert, 1999 for more details).

This idea may be generalized to dependent samples $\mathbf{u}_1, \dots, \mathbf{u}_m$ generated by MCMC routines. The difficulty, of course, is obtaining an appropriate CLT and then estimating the corresponding Monte Carlo variance for the requisite confidence interval. The MCMC literature presents a number of CLTs depending on the induced Markov chain (see for example, Kipnis and Varadhan, 1986; Tierney, 1994; Robert, 1995a; and Robert and Casella, 1999 to name a few). Based on this work, application of a CLT appears feasible in the MCMC setting upon checking the appropriate regularity conditions. However, computation of the variance of the asymptotic distribution is particularly difficult under these CLTs. The dependency between MC samples forces a variance estimation equivalent to estimating the spectral density function at frequency zero (see Geyer, 1992 and Tierney, 1994).

As an alternative approach, renewal theory and regenerative simulation presents methods for extracting independent subsets of the Markov chain under which the classical CLT may be applied (see Robert, 1995a and Mykland et al., 1995). We will apply ideas of Robert et al. (1998) along these lines.

Recall we have a sample $\mathbf{u}_1, \dots, \mathbf{u}_m$ taken from the distribution $g(\mathbf{u} | \mathbf{y}, \Psi)$ under stationarity of the Markov chain. We will assess Monte Carlo error in our MCEM algorithm through a confidence interval about some function $h(\mathbf{u})$ as defined later. The confidence interval will require three pieces: a CLT to ensure normal critical values, and estimates of both $E_g[h(\mathbf{u})]$ and $VAR_g[h(\mathbf{u})]$. We can estimate the expectation with

$$\hat{\mu}_m = \frac{\sum_{t=1}^m w_t h(\mathbf{u}_t)}{\sum_{t=1}^m w_t} \quad (5)$$

and the variance by

$$\hat{v}_m = \left\{ \sum_{t=1}^m w_t [h(\mathbf{u}_t)]^2 / \sum_{t=1}^m w_t \right\} - \hat{\mu}_m^2 \quad (6)$$

where the weights w_t are as defined in Section 2.1. Implementation of a Central Limit Theorem, however, is complicated by the correlation between sample points from the Markov

chain.

We utilize a subsampling scheme from Robert et al. (1998) to overcome Markov chain correlation issues. Specifically, choose the sequence X_1, \dots, X_N such that $X_k - 1 \sim \text{Poisson}(\nu_k)$ where $\nu_k = \nu k^d$ for some $\nu \geq 1$ and $d > 0$ as in Robert et al. (1998). The sums $t_k = x_1 + \dots + x_k$ are used as the subsampling points; that is we consider $\mathbf{u}_{t_1}, \dots, \mathbf{u}_{t_N}$ where $N = \sup\{k : t_k \leq m\}$ is the number of subsamples taken from the m random variates. In the applications of Section 3, $\nu = 1$ and $d = 0.5$.

Let

$$S_m = \frac{1}{\sqrt{N\hat{v}_m}} \sum_{k=1}^N [h(\mathbf{u}_{t_k}) - \hat{\mu}_m] \quad (7)$$

for $j = 1, \dots, s$. Recall that a Markov chain $\mathbf{u}_1, \dots, \mathbf{u}_m$ is *strongly mixing* or α -*mixing* if

$$\alpha(t) = \sup_{A,B} |P(\mathbf{u}_t \in A, \mathbf{u}_0 \in B) - P(\mathbf{u}_t \in A)P(\mathbf{u}_0 \in B)|$$

approaches zero as t goes to infinity and $\mathbf{u}_0 \sim g(\mathbf{u} | \mathbf{y}, \Psi)$. If $\alpha(t) \leq C\gamma^t$ for some constant $C \geq 0$ and $\gamma \in (0, 1)$, then the mixing coefficients are *geometrically decaying*. We then have the following central limit theorem.

Theorem 1 (Robert et al., 1998) *If the Markov chain $\{\mathbf{u}_k\}$ is ergodic and strongly mixing with geometrically decaying mixing coefficients, and*

$$E_g |h(\mathbf{u})|^{2+\delta} < \infty$$

for some $\delta > 0$, then the normalized sum S_m converges weakly to the standard normal distribution.

We choose s functions $h_j(\mathbf{u}) = \frac{\partial}{\partial \psi^{(j)}} \ln f(\mathbf{u}, \mathbf{y} | \Psi)$, $j = 1, \dots, s$, to study MC error in estimation of each of the s components of Ψ . Let $\mathbf{h}(\mathbf{u}) = (h_1(\mathbf{u}), \dots, h_s(\mathbf{u}))'$. Under this choice, the expectation is

$$E_g[\mathbf{h}(\mathbf{u})] = Q^{(1)}(\Psi | \Psi') = \frac{\partial Q(\Psi | \Psi')}{\partial \Psi}.$$

Theorem 1 will be used to choose the Monte Carlo sample size m each iteration. For direct application of the theorem, we assume the components of the s -vector $\mathbf{h}(\mathbf{u})$ are independent.

In particular, we will construct a $(1 - \alpha)$ confidence region around $Q^{(1)}(\Psi \mid \hat{\Psi}^{(r)})$, namely, the set of intervals for each $j = 1, \dots, s$,

$$\hat{\mu}_{m;j} \pm z_{1-\alpha/2} \cdot \hat{v}_{m;j}. \quad (8)$$

Here $z_{1-\alpha}$ denotes the $(1 - \alpha)$ critical value of the standard normal distribution and $\hat{\mu}_{m;j}$ and $\hat{v}_{m;j}$ are the estimators of the expectation and variance of $h_j(\mathbf{u})$ from (5) and (6) respectively. Note these estimators for the mean and variance are computed by substituting the most recent MLE update $\hat{\Psi}^{(r+1)}$.

We will increase m following iteration $(r + 1)$ if for any $j = 1, \dots, s$ the value from the previous iteration,

$$Q_{m;j}^{(1)}(\hat{\Psi}^{(r)} \mid \Psi^{(r-1)}) = \sum_{k=1}^N w_{t_k} \frac{\partial}{\partial \psi(j)} \ln f(\mathbf{u}_{t_k}, \mathbf{y} \mid \Psi) \bigg/ \sum_{t=1}^m w_{t_k} \bigg|_{\Psi = \hat{\Psi}^{(r)}},$$

based on the subsample, lies in the confidence interval (8). The rationale is that if the confidence interval covers the r th value, then the Monte Carlo error swamps the current MLE estimate $\hat{\Psi}^{(r+1)}$. Consequently, more Monte Carlo samples are required to attain reasonable precision at the next iteration. We choose to increase m to $m + m/c$ where c is a positive constant, as suggested by Booth and Hobert (1999). In the application of Section 3, $c = 3$ and $\alpha = 0.25$.

A number of remarks are in order.

1. Most Markov chains induced by MCMC routines such as the Metropolis-Hastings algorithm and the Gibbs sampler applied in practice are ergodic (for example, see Robert and Casella, 1999; Chapters 6 and 7).
2. The α -mixing criterion of Theorem 1 may seem foreboding. However, every positive recurrent aperiodic Markov chain is α -mixing (Robert, 1994). Additionally, relationships exist between α -mixing and other types of mixing as well as minorization conditions (for example see Tierney, 1994).
3. The presumed independence between components of $\mathbf{h}(\mathbf{u})$ may be lifted through an

application of the Cramér-Wold device (Robert et al., 1998); but for simplicity we restrict ourselves to the rectangular confidence region in (8).

4. Note $Q_m^{(1)}(\Psi^{*(r+1)} | \hat{\Psi}^{(r)}) = \mathbf{0}$ for the true MLE $\Psi^{*(r+1)}$ at iteration $(r + 1)$. Therefore the mean value $Q_m^{(1)}(\hat{\Psi}^{(r+1)} | \hat{\Psi}^{(r)}) \approx \mathbf{0}$. We may thus reduce computations slightly by constructing the confidence interval (8) to be symmetric around zero.

2.3 Algorithm

1. Initialize $m, \Psi^{(0)}$.
2. Generate $\mathbf{u}_1, \dots, \mathbf{u}_m \sim g(\mathbf{u} | \mathbf{y}, \Psi^{(0)})$ via a Metropolis-Hastings algorithm.

At iteration $r + 1$

3. Compute the importance weights

$$w_t = \frac{L(\hat{\Psi}^{(r)} | \mathbf{u}_t, \mathbf{y})}{L(\Psi^{(0)} | \mathbf{u}_t, \mathbf{y})}.$$

4. *E-step*: Estimate $Q(\Psi | \hat{\Psi}^{(r)})$ by

$$Q_m(\Psi | \hat{\Psi}^{(r)}) = \sum_{t=1}^m w_t \ln f(\mathbf{u}_t, \mathbf{y} | \Psi) / \sum_{t=1}^m w_t$$

5. *M-step*: Maximize $Q_m(\Psi | \hat{\Psi}^{(r)})$ to obtain $\hat{\Psi}^{(r+1)}$.

6. *MC error estimation*:

- a. Compute for each $j = 1, \dots, s$

$$\hat{\mu}_{m;j} = \sum_{t=1}^m w_t \frac{\partial}{\partial \psi(j)} \ln f(\mathbf{u}_t, \mathbf{y} | \Psi) / \sum_{t=1}^m w_t \Big|_{\Psi = \hat{\Psi}^{(r+1)}}$$

- b. Compute for each $j = 1, \dots, s$

$$\hat{v}_{m;j} = \sum_{t=1}^m w_t \left[\frac{\partial}{\partial \psi(j)} \ln f(\mathbf{u}_t, \mathbf{y} | \Psi) \right]^2 / \sum_{t=1}^m w_t - \hat{\mu}_{m;j}^2 \Big|_{\Psi = \hat{\Psi}^{(r+1)}}$$

c. Obtain for each $j = 1, \dots, s$ a $(1-\alpha)$ confidence interval about $Q_j^{(1)}(\Psi | \hat{\Psi}^{(r)})$

$$\hat{\mu}_{m;j} \pm z_{1-\alpha/2} \cdot \hat{v}_{m;j}$$

where $z_{1-\alpha/2}$ is the $(1 - \alpha/2)$ cutoff of the standard normal distribution.

7. Obtain subsampling instants $t_k = x_1 + \dots + x_k$ where $x_k - 1 \sim \text{Poisson}(\nu_k)$, $k = 1, \dots, N$ and $N = \sup\{n : t_n \leq m\}$.

8. If $Q_m^{(1)}(\hat{\Psi}^{(r)} | \Psi^{(r-1)})$ lies in the confidence region from step 6, then

a. Set $m_o = m$

b. Set $m = m_o + \lfloor m_o/c \rfloor$ for some $c > 0$.

c. Obtain $\mathbf{u}_{m_o+1}, \dots, \mathbf{u}_m \sim g(\mathbf{u} | \mathbf{y}, \Psi^{(0)})$ via a Metropolis-Hastings algorithm.

9. Compute for each $j = 1, \dots, s$

$$Q_{m;j}^{(1)}(\hat{\Psi}^{(r+1)} | \hat{\Psi}^{(r)}) = \sum_{k=1}^N w_{t_k} \frac{\partial}{\partial \psi^{(j)}} \ln f(\mathbf{u}_{t_k}, \mathbf{y} | \Psi) \bigg/ \sum_{t=1}^m w_{t_k} \bigg|_{\Psi = \hat{\Psi}^{(r+1)}}$$

10. Repeat steps three through nine until convergence.

As mentioned in Section 2.1, the importance weights are sensitive to the “distance” between the target and candidate distribution, namely, at iteration r , $g(\mathbf{u} | \mathbf{y}, \hat{\Psi}^{(r)})$ and $g(\mathbf{u} | \mathbf{y}, \Psi^{(0)})$. We can include the following burn-in in step one to alleviate this problem.

1. Initialize m , $\Psi^{(0)}$, and run burn-in

a. Set importance weights $w_t = 1$ for all $t = 1, \dots, m$.

At iteration b

b. Generate $\mathbf{u}_1, \dots, \mathbf{u}_m \sim g(\mathbf{u} | \mathbf{y}, \Psi^{(b)})$ via a Metropolis-Hastings algorithm.

c. Run E and M steps above with $r = b$.

d. Repeat steps (1b) and (1c) for B burn-in iterations.

e. Re-initialize $\Psi^{(0)} = \Psi^{(B)}$.

In the application of Section 3, we run the burn-in for one minute. Note that m is not changed during the burn-in process. This burn-in is essentially the McCulloch (1997) MCEM where random variates are generated each iteration.

3 Applications

We apply the algorithm in Section 2.3 to fit a simple logit-normal model from McCulloch (1997), a survival model from Fan et al. (1998) under censored data, and the salamander data of McCullagh and Nelder (1989).

3.1 Logit-Normal model

Suppose for $i = 1, 2, \dots, q$ and $j = 1, 2, \dots, n$

$$Y_{ij} \mid \mathbf{u} \sim \text{Bernoulli}(\pi_{ij}) \quad (9)$$

conditionally independent. Furthermore,

$$u_i \sim N(0, \sigma^2) \quad (10)$$

independent and identically distributed over i . The logit link function relates \mathbf{Y} to \mathbf{u} in that

$$\ln \left(\frac{\pi_{ij}}{1 - \pi_{ij}} \right) = \beta x_{ij} + u_i. \quad (11)$$

Thus the latent variables here are the random effects \mathbf{u} . McCulloch (1997) simulates data from this single random effects model with $q = 10$, $n = 15$, $\beta = 5$, $\sigma^2 = 0.5$, and $x_{ij} = i/15$. We use the data generated by Booth and Hobert (1999) presented in Table 1. The exact MLE computed from numerical integration for this data is $\hat{\beta} = 6.132$ and $\hat{\sigma}^2 = 1.766$.

Figures 1 and 2 display output from our algorithm with and without importance sampling. The algorithms are run on a 533 MHz DEC alpha with 128 MB RAM under initial values $\Psi^{(0)} = (\beta^{(0)}, \sigma^{(0)}) = (2, 1)$ and $m = 100$. Each iteration of the sample $\mathbf{u}_1^{(r)}, \dots, \mathbf{u}_m^{(r)}$

i/j	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
1	1	0	0	0	0	1	1	0	1	1	1	1	1	1	1
2	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1
3	0	1	0	1	1	1	1	1	1	1	1	1	1	1	1
4	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
5	0	1	1	1	1	1	1	1	1	1	0	1	1	1	1
6	0	0	0	1	0	1	1	1	0	1	1	1	1	1	1
7	0	1	0	0	1	1	1	1	1	1	1	1	1	1	1
8	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
9	1	0	0	1	1	0	1	1	1	1	1	1	1	1	1
10	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1

Table 1: Simulated data from Booth and Hobert (1999) for the logit-normal model.

is generated from a Metropolis-Hastings algorithm with $N(0, \sigma^2)$ candidate distribution. McCulloch (1997) presents the relevant likelihoods and acceptance probabilities.

Both MCEM approaches converge to a reasonable neighborhood around the true MLEs, but continue to show variation. As stated by McCulloch (1997), the number of replications required to obtain the MLE within three or four decimal places of accuracy would be very large. As expected, the importance sampling algorithm is faster performing 108 iterations including 16 burn-in in sixty minutes as opposed to 58 iterations run by MCEM without importance sampling. Generation of an MCMC sample each iteration of the algorithm is quite expensive.

The Markov chain induced by the Metropolis-Hastings algorithm with normal candidate distribution is strongly mixing with geometrically decaying mixing rates (see the Appendix). We may thus monitor the Monte Carlo sample size through the Central Limit Theorem of Theorem 1. The sample size m increases from 100 to 4197. Figure 3 displays the iterations at which the sample size is increased and by how much. Also included in Figure 3 is McCulloch's (1997) predetermined choice for increasing m . For comparison purposes, we

excluded the burn-in from the iteration count as m is not increased during the burn-in process. Note the predetermined choices of m proposed by McCulloch (1997) for running MCEM without importance sampling follows the sample size for MCEM under importance for the first 40 iterations. Following iteration 40, the predetermined m 's take a large jump to 5000, overshooting the sample size for MCEM without importance sampling. Thus, at early stages of the algorithm the McCulloch (1997) MCEM does not increase the MC sample size quickly enough. At later stages, the McCulloch (1997) MCEM runs at an unnecessarily high computational cost by generating more MCMC variates than are needed.

Our stopping rule at this point is based purely on time, namely stop after sixty minutes. We can apply the stopping rules suggested by Booth and Hobert (1999) used to diagnose convergence. A temporal stopping rule is used to allow for easy cost comparisons between the various MCEM algorithms.

3.2 Survival analysis

The MCEM algorithm provides a mechanism for fitting survival functions in the presence of censored data. For example, Wei and Tanner (1990) introduce the MCEM algorithm via an application to regression analysis with censored data and Sinha et al. (1994) analyze grouped survival data via the MCEM algorithm.

In this section, we consider an application of our MCEM algorithm to the twin data set (Duffy et al., 1990). This data consists of responses to a survey of twins in Australia concerning their past history of a number of diseases and operations. One of the operations of interest is appendectomy, in particular, the genetics of acute appendicitis. Of course, not all twins in the study had an appendectomy at the time of the survey. Hence, the data are censored in that for some subjects we know only the appendectomy has not occurred before some given time.

The censored data can be viewed as incomplete or missing data. Let $T \geq 0$ denote a random variable representing the age at which the subject undergoes an appendectomy. Denote the probability of event occurrence or “survival” probability by $G(t; \Psi) = P(T \geq$

t), characterized by some unknown parameter vector Ψ , with associated density $g(t; \Psi)$. Assume the study consists of n uncensored individuals and k censored subjects, combined in the $(n + k)$ -vector $\mathbf{T} = (T_1, \dots, T_n, T_{n+1}, \dots, T_{n+k})'$. Note the variables $(T_{n+1}, \dots, T_{n+k})$ are not observable. Instead we observe the k censoring times c_1, \dots, c_k .

The log-likelihood, $l(\Psi; T)$ consists of a contribution from the uncensored and censored subjects,

$$l(\Psi; \mathbf{T}) = \sum_{u=1}^n \ln[g(T_u; \Psi)] + \sum_{i=1}^k \ln[G(c_i; \Psi)] \quad (12)$$

assuming independent subjects (Cox and Oakes, 1984: Chapter 3). However, the log-likelihood $l_0(\theta; T)$ of the complete data \mathbf{T} that would be observable in the absence of censoring may have a simpler form as it depends solely on the survival density $g(t; \theta)$

$$l_0(\Psi; \mathbf{T}) = \sum_{i=1}^{n+k} \ln[g(T_i; \Psi)].$$

The EM algorithm utilizes the complete data likelihood by maximizing at iteration r the expectation

$$Q(\Psi | \hat{\Psi}^{(r)}) = E[l_0(\hat{\Psi}^{(r)}; \mathbf{T}) | T_1, \dots, T_n, \mathbf{v}; \Psi]$$

where \mathbf{v} is an $(n + k)$ -vector containing the censoring status of each subject. Here, the unknown observations for the censored subjects, $\mathbf{u} = (T_{n+1}, \dots, T_{n+k})'$ are the latent variables. In many instances, the E -step is more difficult to compute in closed form, than numerical maximization in the M -step (Cox and Oakes, 1984; Chapter 11). Consequently, the MCEM algorithm suggests itself as a useful algorithm in the survival function estimation setting.

For purposes of illustration and comparison of the MCEM algorithms presented in Section 2, we fit the survival function on a subset of the twin data set: one of the twins from each of the 327 pairs, 50 of which did not have an appendectomy at the time of the survey (15% censoring). We will fit the survival function

$$G(t) = \exp \left[-\lambda t + \lambda \xi e^{-\delta a} \left(\frac{1 - e^{-\delta t}}{\delta} \right) \right] \quad (13)$$

introduced by Fan et al. (1998). This distribution models the occurrence of shocks to the system as well as aging. In particular, the occurrence of acute appendicitis is a consequence of

stresses to the human body. As the subject ages, the body is increasingly more susceptible to damage by these stresses. The function (13) models these processes through a Poisson shock process under exponential aging with shock process rate $\lambda > 0$, shock damage $\xi \in (0, 1)$, aging parameter $\delta > 0$, and initial age $a \geq 0$ (for more details and derivation, see Fan et al., 1998).

Assume the initial age of the subject is zero (birth) and shock damage $\xi = 1$. The survival model is characterized by two parameters. The MLE computed via numerical maximization from the likelihood function (12) is $\hat{\lambda} = 0.11$, and $\hat{\delta} = 0.038$. Figures 4 and 5 display convergence of the MCEM algorithm with and without importance sampling under initial values of $(\lambda^{(0)}, \delta^{(0)}) = (0.3, 0.5)$ and $m = 50$. The final values are $\hat{\lambda} = 0.11$ and $\hat{\delta} = 0.037$.

At each iteration, the sample $\mathbf{u}_1^{(r)}, \dots, \mathbf{u}_m^{(r)}$ is generated from a Metropolis-Hastings algorithm with truncated normal instrumental distribution (truncated at the observed censoring times). The mean and standard deviation of the candidate are taken as the numerically approximated mean and standard deviation of the time T given $\hat{\Psi}^{(r)}$.

A truncated normal candidate distribution is suggested from the density $g(t; \Psi)$. In particular, expanding the exponential terms of the density in a Taylor series and setting $\xi = 1$ produces

$$g(t; \Psi) \propto \lambda \delta t (1 - \delta t) \cdot \exp\left(\frac{-\lambda \delta t^2}{2}\right) \cdot I(t > c)$$

where c is the censoring time and $I(A)$ is the indicator function over the set A . The truncated normal distribution provides a candidate with heavier tails than this density, allowing for acceptance rates greater than 60%.

The exponential distribution also suggests itself as a potential instrumental distribution. In particular, the survival function (13) is motivated as an extension of the exponential survival function in the presence of aging (Fan et al., 1998). However, an exponential candidate exhibits a lighter tail resulting in prohibitively slow convergence of the Metropolis-Hastings chain. Nonetheless, the exponential form of the survival function motivated the implementation of a Metropolis-Hastings scheme to obtain a sample from this complicated distribution.

Both algorithms, with and without importance sampling, perform admirably, quickly oscillating in a small neighborhood of the MLE. In fact, the Matlab maximization routines involve an equal amount of computation time to attain the numerical MLE $\hat{\Psi}$ as the MCEM. The MCEM routine with importance sampling performed 58 iterations including 29 burn-in, comparable to the 22 iterations run by MCEM without importance sampling.

The Markov chain induced by the Metropolis-Hastings algorithm with truncated normal candidate distribution is strongly mixing with geometrically decaying mixing coefficients (see the Appendix). Thus we utilize Theorem 1 to adjust the Monte Carlo sample size in both algorithms. Figure 6 displays the change in m over the MCEM algorithms. The MCEM sample size increases from 50 to 6704 under importance sampling and from 50 to 2828 when drawing a Metropolis-Hastings sample each iteration. The Metropolis-Hastings algorithm is increasingly computer intensive as the number of censored subjects k increases. Hence, importance sampling may be particularly useful in such circumstances.

3.3 Salamander data

As a final illustration of the MCEM algorithm developed in this paper, we consider the salamander data of McCullagh and Nelder (1989), Section 14.5. The data consists of mating success between two species of salamanders, rough-butt (R) and whiteside (W). Twenty males and twenty females, ten of each of the two species were mated resulting in four types of crosses. The design we follow is from Table 14.3 of McCullagh and Nelder (1989) consisting of $n = 120$ matings from the first of three experiments. The observed outcome is binary denoted by W_i taking a value one if the i th mating is successful and zero otherwise.

We apply the MCEM algorithm to perform a hierarchical probit regression. The algorithm utilizes the Gibbs sampler to estimate the expectation in the E-step. The model fitting is thus analogous to the work of McCulloch (1994), Chan and Kuk (1997), and van Dyk (1998).

The probit-normal model is a threshold model, assuming the underlying response is an unobserved continuous random n -vector \mathbf{Y} which governs mating success. The observed W_i

states whether Y_i exceeds the threshold of zero or not. The effect of each salamander on the outcome is assumed random while the effect of the species cross is assumed fix. In this way, the mixed model is written

$$\begin{aligned} \mathbf{Y} &= \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}_f\mathbf{U}_f + \mathbf{Z}_m\mathbf{U}_m + \boldsymbol{\epsilon} \\ \mathbf{U}_f &\sim N_{20}(\mathbf{0}, \sigma_f^2\mathbf{I}) \\ \mathbf{U}_m &\sim N_{20}(\mathbf{0}, \sigma_m^2\mathbf{I}) \\ \boldsymbol{\epsilon} &\sim N_n(\mathbf{0}, \mathbf{I}) \end{aligned} \tag{14}$$

where \mathbf{U}_f and \mathbf{U}_m are normal random q -vectors denoting the effect of the female and male salamanders respectively, $\boldsymbol{\beta} = (\beta_{RR}, \beta_{RW}, \beta_{WR}, \beta_{WW})^T$ denotes the fixed effects of each of the four species crossings, \mathbf{X} is the design matrix for the type of cross, \mathbf{Z}_f and \mathbf{Z}_m denotes the design matrix for the female and male salamanders respectively, and \mathbf{I} denotes an identity matrix of the appropriate size, $\mathbf{0}$ denotes a vector of zeros of the appropriate dimension, and T denotes vector transpose. The random effects are assumed independent of each other and mutually independent of the normal random errors $\boldsymbol{\epsilon}$. Let the parameter vector be denoted by $\boldsymbol{\Psi} = (\boldsymbol{\beta}^T, \sigma_f^2, \sigma_m^2)^T$.

Under model (14), the complete data is $(\mathbf{Y}, \mathbf{U}_f, \mathbf{U}_m)$, as the latent variable \mathbf{Y} is unobserved in addition to the random effects. However, through results about the multivariate normal density, we can show, at iteration r of the EM algorithm, calculation of $Q(\boldsymbol{\Psi} | \hat{\boldsymbol{\Psi}}^{(r)})$ is dependent solely on \mathbf{Y} (see McCulloch, 1994 for details). Thus the MC estimation in the E-step requires a sample from the latent variables \mathbf{Y} .

Recall the observations are binary responses $\mathbf{W} = (W_1, \dots, W_n)^T$. Thus the distribution from which we must sample is the truncated multivariate normal conditional distribution of \mathbf{Y} given \mathbf{W} truncated at zero. Note the full conditional distributions of $Y_i | Y_1, \dots, Y_{i-1}, Y_{i+1}, \dots, Y_n$ for $i = 1, \dots, n$ is univariate truncated normal distributions, truncated at zero. As the univariate truncated normal is easy to sample, a Gibbs sampler may be implemented to obtain a sample $\mathbf{Y}^{(1)}, \dots, \mathbf{Y}^{(m)}$ from the truncated multivariate normal distribution of interest (Robert, 1995b). McCulloch (1994) presents the Gibbs routine for sampling from the truncated multivariate normal distribution in detail.

At iteration i of the Gibbs sampler, Y_i is generated from the truncated univariate normal distribution given the most recent value of the other $n - 1$ components. Thus, at each iteration of the Gibbs chain, we must compute n conditional expectations as well as generate a sample from the truncated normal distribution. This task is costly, particularly for later iterations of the MCEM algorithm in which m may be of the order of tens of thousands. Hence we implement the MCEM algorithm with importance sampling as discussed in Section 2.3. The algorithm is initiated at $\Psi = (0, 0, 0, 0, 0.2, 0.2)$ and $m = 10$.

Figure 7 displays convergence of the ML estimates for the parameter Ψ . The algorithm was run for 60 minutes during which it performed 46 EM iterations including 16 burn-in. The final estimates are $\hat{\beta}_{RR} = 0.81$, $\hat{\beta}_{RW} = 0.54$, $\hat{\beta}_{WR} = -0.96$, $\hat{\beta}_{WW} = 0.73$, $\hat{\sigma}_f^2 = 0.62$, $\hat{\sigma}_m^2 = 0.088$, similar to the values found by McCulloch (1994).

Robert (1995b) states the Gibbs sampler utilized to generate the truncated multivariate normal sample induces a geometrically ergodic Markov chain. Consequently, the chain is strongly mixing with geometrically decaying mixing rates (Chan and Geyer, 1994) putting Theorem 1 into use. The MC sample size m increases from 10 to 12857. The large increase in MC sample size suggests the recommendation of Chan and Kuk (1997) to generate a Gibbs sample of 1000 iterates each EM iteration is not appropriate. This fixed sampling routine oversamples at early iterations of the MCEM algorithm and does not draw enough samples in later iterations to obtain precise EM estimates.

The amount by which the MC sample should be increased in our algorithm, of course, is flexible. Our algorithm tells the user when to increase the MC sample size, but the actual value taken by m each EM iteration is user determined. As van Dyk (1998) states, trial and error is necessary for choosing appropriate MC sample sizes. Nonetheless, a linear increase of m seems to be the preferred choice in the literature (see for example McCulloch, 1994 and van Dyk, 1998). Though our choice of increasing m by 4/3 appears to allow a reasonable tradeoff of computational cost with EM precision, users may choose to increase or decrease this value as deemed appropriate for the application of interest.

4 Appendix

The Metropolis-Hastings algorithms utilized in Sections 3.1 and 3.2 are both independence chains in that the candidate distribution does not depend on the most recent sample point (see Robert and Casella, 1999, Chapter 6 for more details). We will use the following theorem concerning convergence of independent Metropolis-Hastings algorithms, a combination of a result from Chan and Geyer (1994) and Theorem 2.1 of Mengersen and Tweedie (1996).

Theorem 2 *The Markov chain induced by the independent Metropolis-Hastings algorithm is a strongly mixing chain with geometrically decaying mixing coefficients if there exists some $M > 0$ such that for all $\mathbf{u} \in \text{supp}(\pi)$, $q(\mathbf{u}) > 0$ and*

$$\frac{\pi(\mathbf{u})}{q(\mathbf{u})} \leq M$$

where π is the stationary distribution of the chain and q is the candidate distribution.

Mengersen and Tweedie (1996) actually show under the conditions of Theorem 2, the chain is *uniformly ergodic*, a much stronger conclusion than α -mixing. Chan and Geyer (1994) show uniformly ergodic chains are strongly mixing with geometrically decaying mixing rate. We will apply this theorem to the logit-normal and survival analysis examples in turn.

Recall the candidate distribution $q(\mathbf{u})$ in the logit-normal example of Section 3.1 is multivariate normal $N(\mathbf{0}, \sigma^2 \mathbf{I})$ where \mathbf{I} is an $s \times s$ identity matrix with $s = 10$. Hence generations from the candidate distribution are independent of variates generated earlier by the Metropolis-Hastings sampler.

The invariant distribution $\pi(\mathbf{u})$ of the Markov chain induced by the Metropolis-Hastings algorithm is the conditional density

$$\begin{aligned} f(\mathbf{u} | \mathbf{Y}, \beta, \sigma^2) &= \frac{f(\mathbf{Y} | \mathbf{u}, \beta, \sigma^2) \cdot f(\mathbf{u} | \beta, \sigma^2)}{f(\mathbf{Y} | \beta, \sigma^2)} \\ &\propto \prod_{j=1}^s \prod_{i=1}^n \frac{\exp\{y_{ij}(\beta x_{ij} + u_j)\}}{1 + \exp\{y_{ij}(\beta x_{ij} + u_j)\}} \cdot \frac{\exp\{-u_j^2/(2\sigma^2)\}}{\sqrt{2\pi\sigma^2}} \end{aligned}$$

where the variables are all as defined in Section 3.1.

Theorem 2 is applicable as

$$\begin{aligned} \frac{f(\mathbf{u} \mid \mathbf{Y}, \beta, \sigma^2)}{q(\mathbf{u})} &\propto \prod_{j=1}^s \prod_{i=1}^n \frac{\exp\{y_{ij}(\beta x_{ij} + u_j)\}}{1 + \exp\{y_{ij}(\beta x_{ij} + u_j)\}} \cdot \frac{\exp\{-u_j^2/(2\sigma^2)\}/\sqrt{2\pi\sigma^2}}{\exp\{\sum_{j=1}^s u_j^2/(2\sigma^2)\}} \\ &\leq \left(\frac{1}{2\pi\sigma^2}\right)^n \end{aligned}$$

where the proportionality constant in the first equation is

$$\frac{(2\pi\sigma^2)^{-q/2}}{f(\mathbf{Y} \mid \beta, \sigma^2)}.$$

The Metropolis-Hastings algorithm in the lifetime distribution example of Section 3.2 utilizes a truncated normal distribution with mean μ , variance σ^2 , and left truncation value being the censoring time for each individual. Thus

$$q(\mathbf{u}) = \prod_{i=1}^k \frac{\exp[-(u_i - \mu)^2/(2\sigma^2)] \cdot I(u_i \geq c_i)}{\sqrt{2\pi\sigma^2} [1 - \Phi(-\mu/\sigma)]}$$

where \mathbf{u} is the k -vector of failure times for the k censored subjects, \mathbf{c} is the k -vector of censoring times, Φ is the normal cumulative distribution function, and $I(A)$ is the indicator function over the set A . The mean, variance, and truncation time, at any iteration in the algorithm, are independent of previous variates generated by the sampler.

Recall the MCEM algorithm requires MC samples only for the censored individuals. The invariant distribution is then the conditional probability of survival given the subject has survived beyond the time of censoring. The associated density $\pi(\mathbf{u})$ from (13) is then

$$g(\mathbf{u} \mid \mathbf{c}) = \frac{\prod_{i=1}^k \lambda(1 - \exp\{-\delta u_i\}) \cdot \exp\{-\lambda u_i + (\lambda/\delta)[1 - \exp(-\delta u_i)]\} \cdot I(u_i \geq c_i)}{-\lambda c_i + (\lambda/\delta)[1 - \exp(-\delta c_i)]}$$

where all variables are as defined in Section 3.2.

Note

$$\begin{aligned} \frac{q(\mathbf{u})}{\pi(\mathbf{u})} &\leq K_1 \prod_{j=1}^k \frac{\exp[-(u_j^2 - 2(\mu + \lambda\sigma^2)u_j + \mu^2)/(2\sigma^2)]}{K_2(j)(1 - \exp\{-\delta c_i\})} \\ &\propto \prod_{j=1}^k \frac{[\exp[-(u_j - \mu - \lambda\sigma^2)^2 - (\mu + \lambda\sigma^2)^2 + \mu^2]/(2\sigma^2)]}{1 - \exp\{-\delta c_i\}} \\ &\leq \frac{\exp[-(1 - (\mu + \lambda\sigma^2)^2 + \mu^2)/(2\sigma^2)]}{1 - \exp\{-\delta c_i\}} \end{aligned}$$

where $K_1 = [\sqrt{2\pi}\sigma(1 - \Phi(-\mu/\sigma))]^{-k}$ and $K_2(j) = \exp\{-\lambda c_j + (\lambda/\delta)[1 - \exp(-\delta c_j)]\}$. Thus Theorem 2 holds and the independent Metropolis-Hastings chain is strongly mixing with geometrically decaying mixing coefficients.

References

- [1] Booth, J. G. and Hobert, J. P. (1999). Maximizing Generalized Linear Mixed Model Likelihoods with an Automated Monte Carlo EM Algorithm. To appear in *Journal of the Royal Statistical Society, B*.
- [2] Casella, G. and Robert, C. P. (1996). Post-processing Accept-Reject Samples: Recycling and Rescaling. To appear in the *Journal of Computational and Graphical Statistics*.
- [3] Chan, J. S. K. and Kuk, A. Y. C. (1997). Maximum Likelihood Estimation for Probit-linear Mixed Models with Correlated Random Effects. *Biometrics* 53, 86-97.
- [4] Chan, K. S. and Geyer, C. (1994). Discussion of "Markov Chains for Exploring Posterior Distributions." *The Annals of Statistics* 22 (4), 1747-1758.
- [5] Chan, K. S. and Ledolter, J. (1995). Monte Carlo EM Estimation for Time Series Models Involving Counts. *Journal of the American Statistical Association* 90, 242-252.
- [6] Cox, D. R. and Oakes, D. (1984). *Analysis of Survival Data*. Chapman and Hall, New York.
- [7] Duffy, D. L., Martin, N. G., Mathews, J. D. (1990). Appendectomy in Australian Twins. *American Journal of Human Genetics* 47, 590-592.
- [8] Fan, J. J., Ghurye, S. G., Levine, R. A. (1998). Multicomponent Lifetime Distributions in the Presence of Aging. Technical report #335, Division of Statistics, University of California, Davis.
- [9] Geyer, C. J. (1992). Practical Markov Chain Monte Carlo. *Statistical Science* 7 (4), 473-483.

- [10] Geyer, C. J. (1991). Reweighting Monte Carlo Mixtures. Technical Report 568, University of Minnesota, School of Statistics.
- [11] Kipnis, C. and Varadhan, S. R. S. (1986). Central Limit Theorem for Additive Functionals of Reversible Markov Processes and Applications to Simple Exclusions. *Communications in Mathematical Physics* 104, 1-19.
- [12] Liu, J. S. (1996). Metropolized Independent Sampling with Comparisons to Rejection Sampling and Importance Sampling. *Statistics and Computing* 6, 113-119.
- [13] McCullagh, P. and Nelder, J. A. (1989). *Generalized Linear Models*, Second Edition. Chapman and Hall, New York.
- [14] McCulloch, C. E. (1997). Maximum Likelihood Algorithms for Generalized Linear Mixed Models. *Journal of the American Statistical Association* 92, 162-170.
- [15] McCulloch, C. E. (1994). Maximum Likelihood Variance Components Estimation for Binary Data. *Journal of the American Statistical Association* 89, 330-335.
- [16] Mengersen, K. L. and Tweedie, R. L. (1996). Rates of Convergence of the Hastings and Metropolis Algorithms. *The Annals of Statistics* 24, 101-121.
- [17] Mykland, P., Tierney, L., and Yu, B. (1995). Regeneration in Markov Chain Samplers. *Journal of the American Statistical Association* 90, 233 - 241.
- [18] Robert, C. P. (1995a). Convergence Control Methods for Markov Chain Monte Carlo Algorithms. *Statistical Science* 10 (3), 231-253.
- [19] Robert, C. P. (1995b). Simulation of Truncated Normal Variables. *Statistics and Computing* 5, 121-125.
- [20] Robert, C. P. (1994). Discussion of "Markov Chains for Exploring Posterior Distributions." *The Annals of Statistics* 22, 1742-1747.

- [21] Robert, C. P. and Casella, G. (1999). *Monte Carlo Statistical Methods*. Springer-Verlag, New York.
- [22] Robert, C. P., Rydén, T., and Titterington, D. M. (1998). Convergence controls for MCMC algorithms, with applications to hidden Markov chains. Technical report from MCMC preprint service <http://www.stats.bris.ac.uk/MCMC>.
- [23] Sinha, D., Tanner, M. A., and Hall, W. J. (1994). Maximization of the Marginal Likelihood of Grouped Survival Data. *Biometrika* 81, 53-60.
- [24] Tanner, M. A. (1993). *Tools for Statistical Inference, second edition*. Springer-Verlag, New York.
- [25] Tierney, L. (1994). Markov Chains for Exploring Posterior Distributions. *The Annals of Statistics* 22 (4), 1701-1762.
- [26] Wei, G. C. G. and Tanner, M. A. (1990). A Monte Carlo Implementation of the EM Algorithm and the Poor Man's Data Augmentation Algorithms. *Journal of the American Statistical Association* 85, 699-704.
- [27] Wu, C. F. J. (1983). On the convergence properties of the EM algorithm. *Annals of Statistics* 11, 95-103.
- [28] van Dyk, D. (1998). Nesting EM Algorithms for Computational Efficiency. To appear in *Statistica Sinica*.

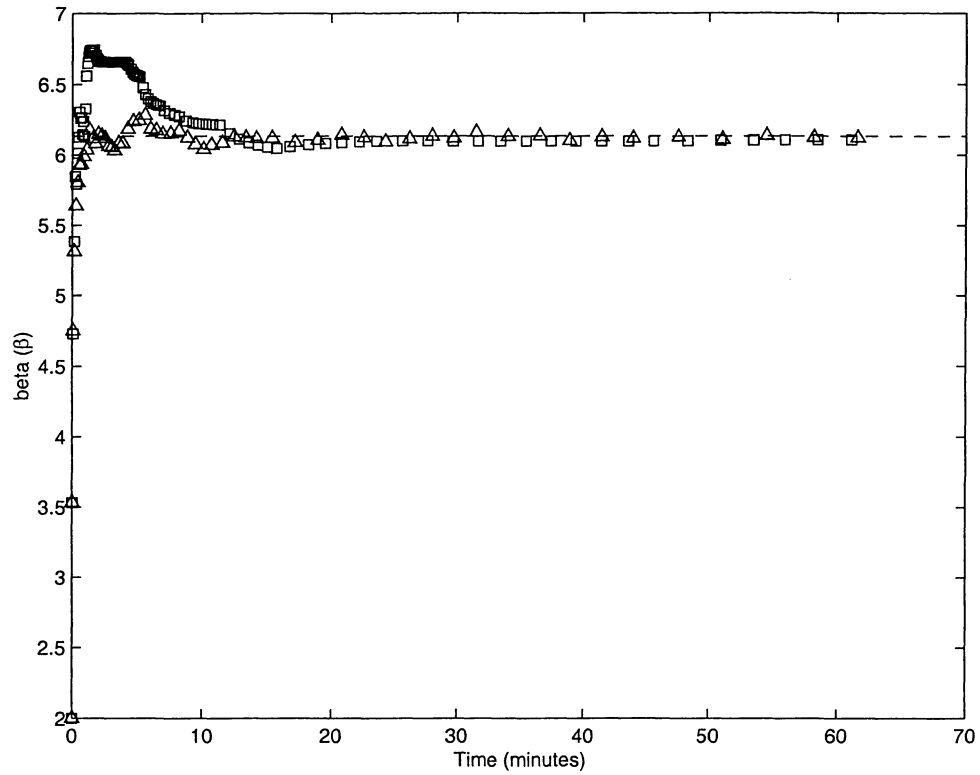


Figure 1: MLE of β at each iteration of MCEM algorithms with importance sampling (\square) and without importance sampling (\triangle). The iteration is measured in minutes. Dotted line specifies the true MLE $\hat{\beta} = 6.132$.

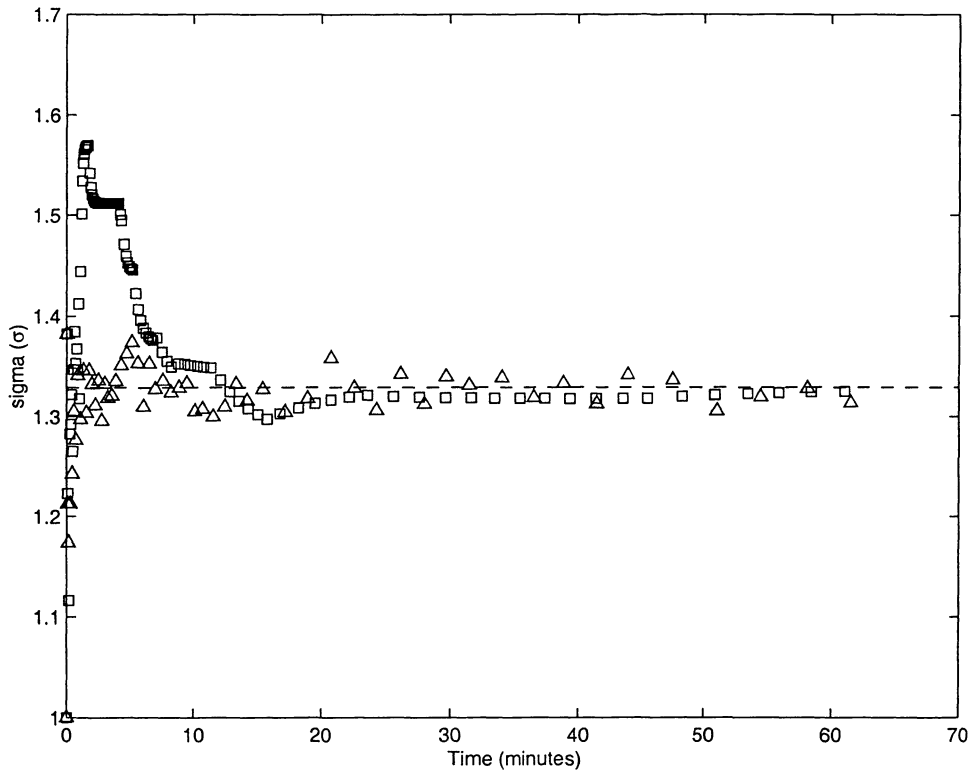


Figure 2: MLE of σ at each iteration of MCEM algorithms with importance sampling (\square) and without importance sampling (\triangle). The iteration is measured in minutes. Dotted line specifies the true MLE $\hat{\sigma} = 1.329$.

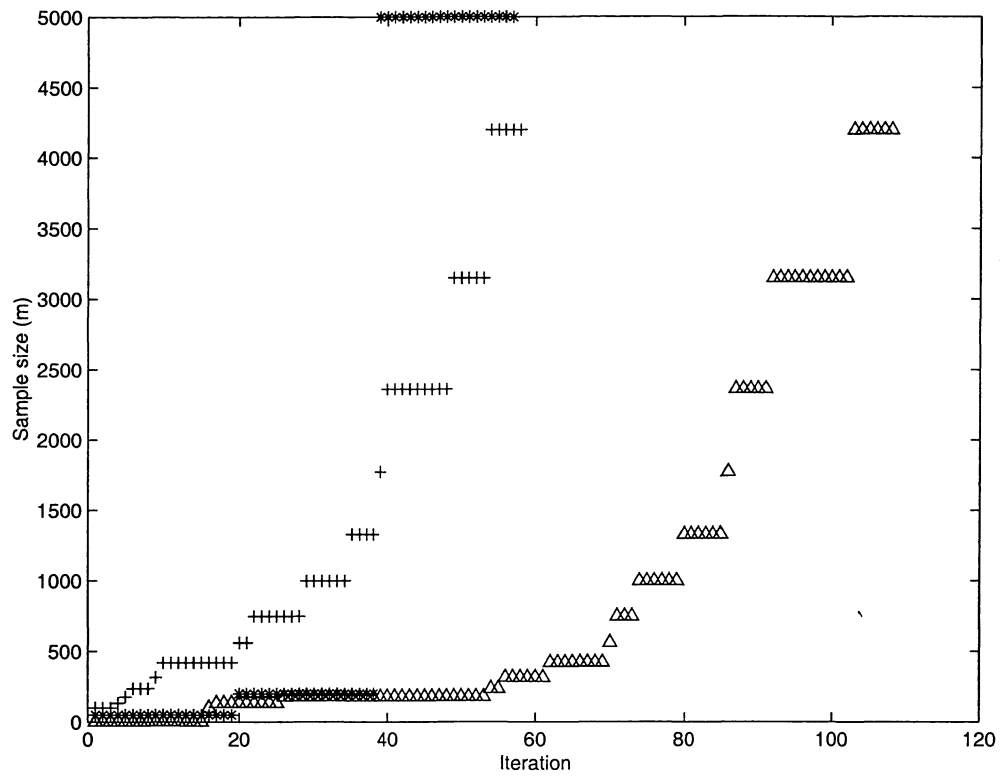


Figure 3: Increases in Monte Carlo sample size m for three MCEM algorithms: MCEM with importance sampling (Δ), without importance sampling (+), and the predetermined values from McCulloch, 1997 (*).

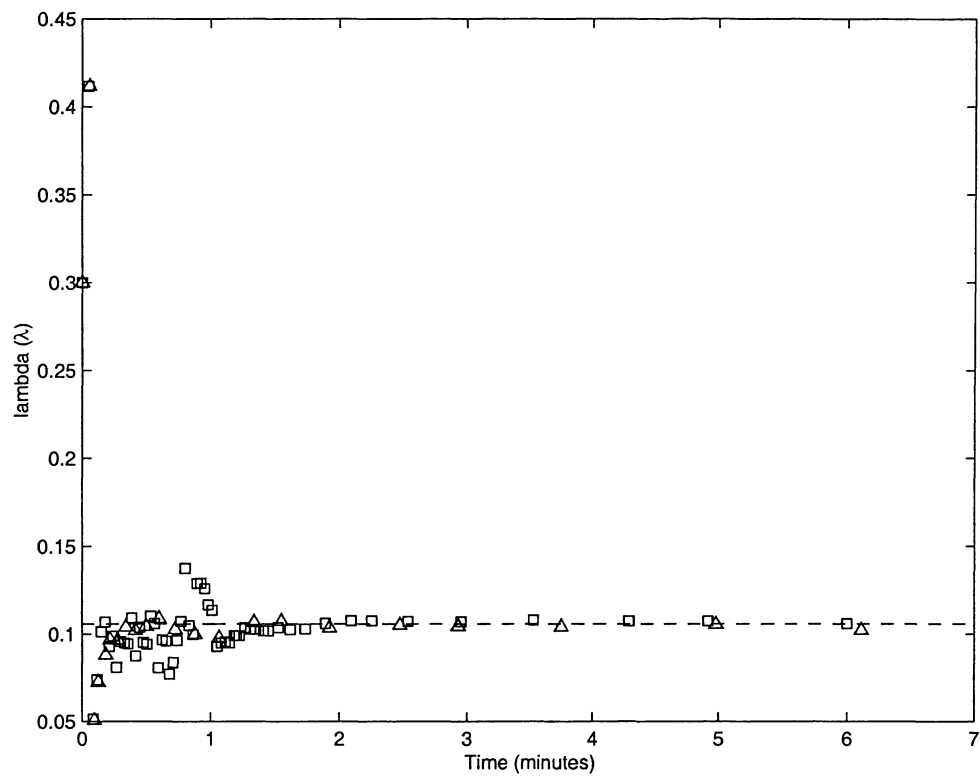


Figure 4: MLE of λ at each iteration of MCEM algorithms with (\square) and without importance sampling (\triangle). Dotted line specifies the true MLE $\hat{\lambda} = 0.11$.

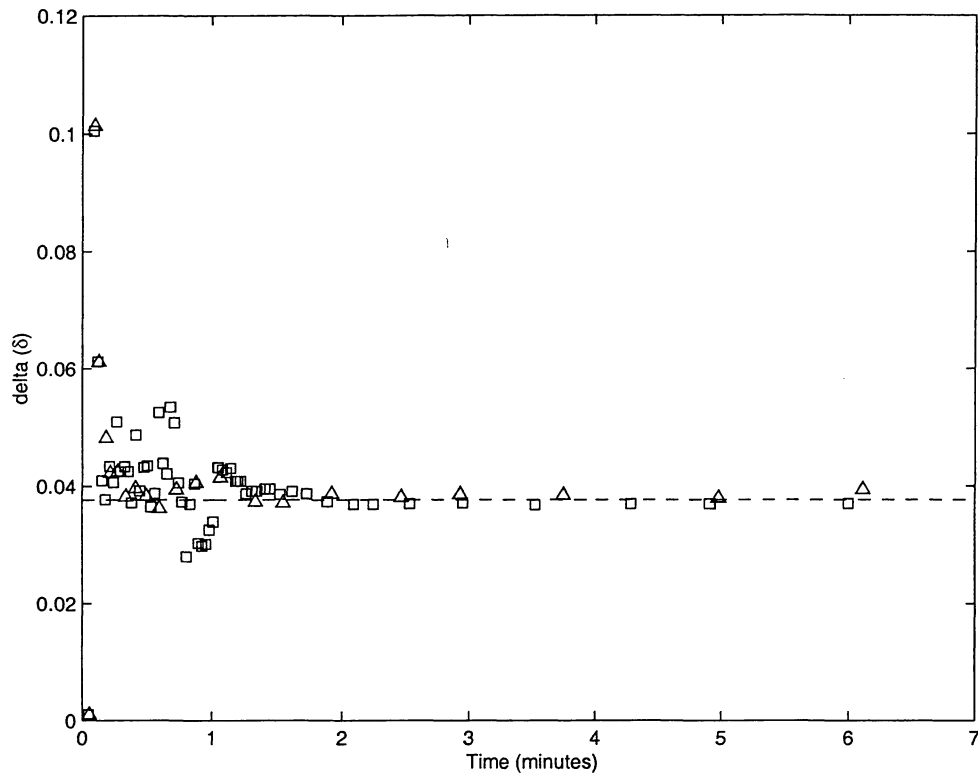


Figure 5: MLE of δ at each iteration of MCEM algorithms with (\square) and without importance sampling (\triangle). Dotted line specifies the true MLE $\hat{\delta} = 0.038$. The initial value of $\delta^{(0)} = 0.5$ is not plotted.

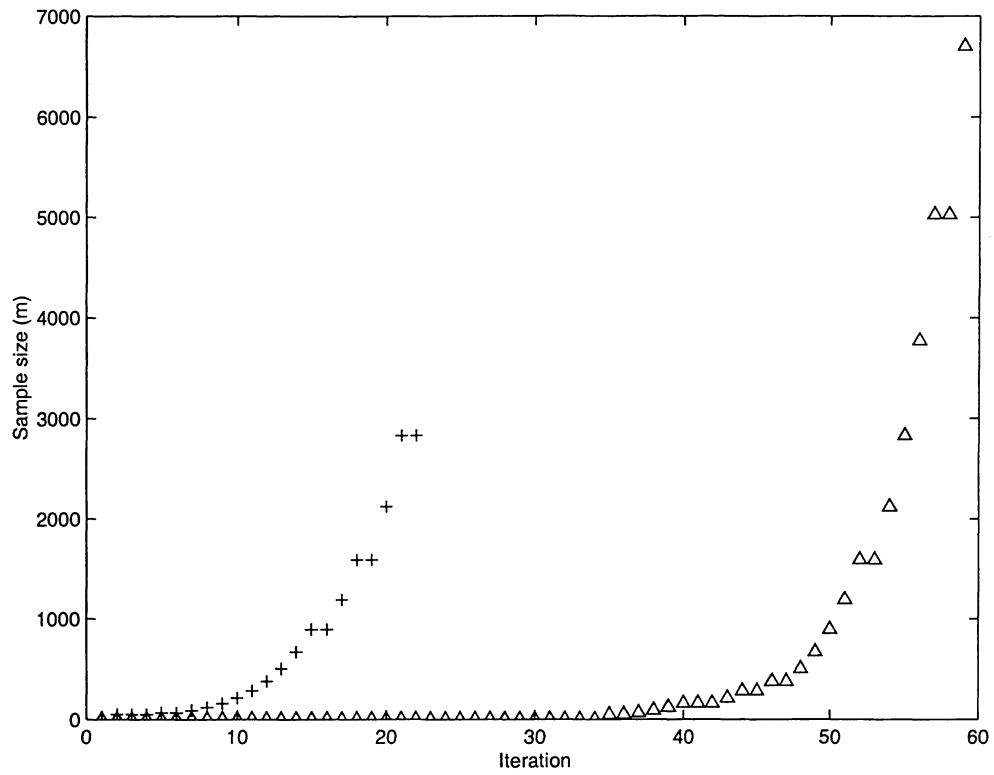


Figure 6: Increases in Monte Carlo sample size m for MCEM with (Δ) and without (+) importance sampling.

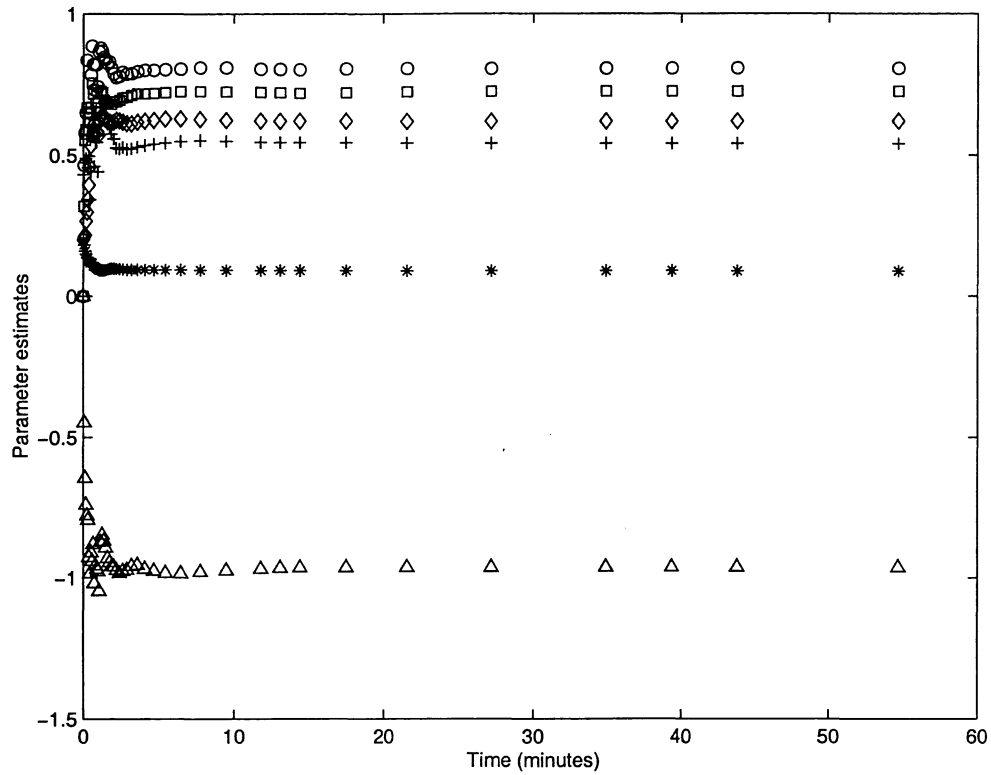


Figure 7: MLE of β_{RR} (○), β_{RW} (+), β_{WR} (△), β_{WW} (□), σ_f^2 (◇), and σ_m^2 (*) at each iteration of MCEM algorithm with importance sampling.