

Abstract

We survey the Monte Carlo EM algorithm, focusing on its implementation and a few alternative methods.

- I have changed my iteration labels. $\hat{\theta}_{k-1}$ is now the maximizer of the current MCEM objective function, and $\hat{\theta}_{k-1}$ was used to construct that objective function. Watch for this when editing.
- Set standard terminology for the observed information matrix of the observed data likelihood. I think it's best to just call this the observed data information.

1 The EM Algorithm

We define three distributions which will be central to our study of missing data problems. Let Y be the observed data and X be the missing data. Note that X need not correspond to any actual real-world process, but may instead be a conceptual device which facilitates analysis of the data which were actually observed. We refer to the distribution of Y as the “observed data distribution”, and write f for its density (or mass function). We refer to the joint distribution of Y and X as the “complete data distribution”, and write f_c for its density. We refer to the conditional distribution of X given Y as the “missing data distribution”, and write f_m for its density. Note that the missing data distribution is not the marginal distribution of the missing data, but rather its conditional distribution given the observed data. We also write ℓ , ℓ_c and ℓ_m for the log-likelihoods based on the observed, complete and missing data distributions respectively. Similar with S , S_c , S_m for the scores (gradients of the matching ℓ 's) and I , I_c , I_m for the observed information matrices (negative Hessians of the matching ℓ 's). We emphasize that, in our notation, a subscript c denotes “complete” rather than “conditional”, and a subscript m denotes “missing” rather than “marginal”. We write $\theta \in \Theta \subseteq \mathbb{R}^p$ for the parameter of interest, and note that f , f_c and f_m are parameterized by the same θ (although in principle they need not all depend on every component of θ)

The EM algorithm is a method for analyzing incomplete data which was formalized by Dempster et al. (1977). See McLachlan and Krishnan (2008) for an excellent book-length overview of the EM algorithm. We begin by discussing a probabilistic framework within which the EM algorithm is often applied. We then present the EM algorithm in detail. Finally, we discuss some limitations of this method. Throughout, we illustrate our presentation with a toy problem based on linear regression with a single, unobserved, covariate.

The EM algorithm consists of iterating two steps. First is the expectation, or “E”, step, in which an objective function is constructed from the complete data likelihood. Second

is the maximization, or “M”, step, in which the previously computed objective function is maximized. These two steps are then alternated until some convergence criterion is met. Whatever value of θ the algorithm converges to is used as our parameter estimate. We now go into more detail on each of the two steps.

The E-step of the EM algorithm is where we construct the objective function which will be used to update our parameter estimate. This objective function is the conditional expectation of the complete data likelihood, given the observed data. If our complete data can be partitioned into an observed component, Y , and a missing component, X , then our objective function at iteration k is given by

$$Q(\theta|\theta_{k-1}) = \mathbb{E}_{\theta_{k-1}}[\ell_c(\theta; y, X)|Y = y] \quad (1)$$

where ℓ_c is the log-likelihood of the complete data model. Note that the conditional expectation uses our parameter estimate from the previous iteration.

The M-step of the EM algorithm consists of maximizing the objective function constructed in the previous E-step. That is, we define $\theta_k = \underset{\theta}{\operatorname{argmax}} Q(\theta|\theta_{k-1})$. Typically, this optimization must be performed numerically via, e.g., gradient ascent or Newton’s method. See Nocedal and Wright (2006) for details and other optimization algorithms.

We can combine the E- and M-steps of the EM algorithm into a single “update function”. We write $M(\theta_{k-1}) = \underset{\theta}{\operatorname{argmax}} Q(\theta|\theta_{k-1})$. The EM algorithm can thus be viewed as the iterative application of this update function, M .

1.1 Properties

Section intro...

1.1.1 Ascent Property and Generalized EM

An important feature of the EM algorithm is its so-called “ascent property”. This property says that an iteration of the EM algorithm never decreases the observed data likelihood. This is somewhat surprising, since EM updates are computed without ever evaluating the observed data likelihood.

Proposition 1.1 (Ascent Property of EM). *Let $\theta \in \Theta$, and $\theta' = M(\theta)$ be the EM update from θ . Then $\ell(\theta') \geq \ell(\theta)$.*

Proof. We begin by noting that the following decomposition holds for any value of x :

$$\ell(\theta; y) = \ell_c(\theta; y, x) - \ell_m(\theta; y, x) \quad (2)$$

Subtracting the values of both sides at θ from their values at θ' and taking conditional expectations, we get

$$\ell(\theta'; y) - \ell(\theta; y) = Q(\theta'|\theta) - Q(\theta|\theta) + \mathbb{E}_\theta[\ell_m(\theta; y, x) - \ell_m(\theta'; y, x)] \quad (3)$$

$$= Q(\theta'|\theta) - Q(\theta|\theta) + \text{KL}(\theta \rightarrow \theta') \quad (4)$$

where the last term in line 4 is the Kullback-Leibler (KL) divergence from the missing data distribution with $\theta = \theta$ to the same distribution with $\theta = \theta'$ (van der Vaart, 1998). Note that KL divergences are always non-negative, so we get

$$\ell(\theta'; y) - \ell(\theta; y) \geq Q(\theta'|\theta) - Q(\theta|\theta) \quad (5)$$

Finally, since θ' maximizes $Q(\cdot|\theta)$, we have $\ell(\theta'; y) - \ell(\theta; y) \geq 0$. \square

In our proof of the ascent property, we only required that $Q(\theta'|\theta) \geq Q(\theta|\theta)$, not that θ' maximize $Q(\cdot|\theta)$. This observation leads to the definition of the “Generalized EM Algorithm”, which replaces the M-step with setting θ_k to any point in Θ such that $Q(\theta_k|\theta_{k-1}) \geq Q(\theta_{k-1}|\theta_{k-1})$.

1.1.2 Recovering Observed Data Likelihood Quantities

Under regularity conditions McLachlan and Krishnan (see 2008), it is possible to compute both the score vector and the observed information matrix of the observed data likelihood using complete data quantities. These regularity conditions consist of being able to interchange the order of differentiation and integration for various functions.

Proposition 1.2. *Provided that differentiation and integration can be exchanged and that all given expectations are finite, the following identities hold:*

- (i) $S(\theta; y) = \mathbb{E}_\theta[S_c(\theta; y, X)|Y = y]$
- (ii) $I(\theta) = \mathcal{I}_c(\theta) - \mathcal{I}_m(\theta)$
where $\mathcal{I}_c(\theta) := -\mathbb{E}_\theta [\nabla^2 \ell_c(\theta; y, X)|Y = y]$ and $\mathcal{I}_m(\theta) := -\mathbb{E}_\theta [\nabla^2 \ell_m(\theta; y, X)|Y = y]$

Proof. We start with expression (i). Let Ω be the complete data sample space. Let \mathcal{Y} and \mathcal{X} be the observed and missing data sample spaces respectively. For every $y \in \mathcal{Y}$, let

$\mathcal{X}(y) = \{x \in \mathcal{X} : (y, x) \in \Omega\}$. Note that $f(y; \theta) = \int_{\mathcal{X}(y)} f_c(y, x; \theta) dx$.

$$\begin{aligned}
\mathbb{E}_\theta[S_c(\theta; y, X)|Y = y] &= \int_{\mathcal{X}(y)} \nabla \ell_c(\theta; y, x) f_m(y, x; \theta) dx \\
&= \int_{\mathcal{X}(y)} \frac{f_m(y, x; \theta)}{f_c(y, x; \theta)} \nabla f_c(\theta; y, x) dx \\
&= \int_{\mathcal{X}(y)} \frac{1}{f(y; \theta)} \nabla f_c(\theta; y, x) dx \\
&= \frac{1}{f(y; \theta)} \int_{\mathcal{X}(y)} \nabla f_c(\theta; y, x) dx \\
&= \frac{1}{f(y; \theta)} \nabla \int_{\mathcal{X}(y)} f_c(\theta; y, x) dx \\
&= \frac{1}{f(y; \theta)} \nabla f(y; \theta) \\
&= S(\theta; y)
\end{aligned}$$

Proceeding now to (ii), we decompose the observed data log-likelihood as

$$\ell(\theta; y) = \ell_c(\theta; y, x) - \ell_m(\theta; y, x)$$

Differentiating twice and taking conditional expectations of both sides yields the required result. \square

Note that the matrices \mathcal{I}_c and \mathcal{I}_m are not observed information matrices (negative Hessians), but conditional expectations of observed information matrices. An alternative to Proposition 1.2 part (ii) which involves only conditional expectations of complete data quantities is given in the following proposition.

Proposition 1.3 (Louis' Identity). *Let $\hat{\theta}$ be a critical point of the observed data log-likelihood. Assuming that differentiation and integration can be exchanged and that all given expectations are finite, we can write the observed information of the observed data distribution at $\hat{\theta}$ as*

$$I(\theta) = \mathcal{I}_c(\theta) - \mathbb{E}_\theta[S_c(\theta)S_c(\theta)^T|Y = y] + S(\theta)S(\theta)^T \quad (6)$$

In particular, if $\hat{\theta}$ is a critical point of the observed data log-likelihood, then

$$I(\hat{\theta}) = (\mathcal{I}_c(\theta) - \mathbb{E}_\theta[S_c(\theta)S_c(\theta)^T|Y = y])|_{\theta=\hat{\theta}} \quad (7)$$

Proof. We follow the derivation of Louis (1982). For brevity, we write $f(\theta)$ and $f_c(\theta)$ for $f(y; \theta)$ and $f(y, x; \theta)$ respectively. Consider the following two Hessians:

$$\nabla^2 \ell(\theta) = \nabla \left[\int_{\mathcal{X}(y)} \frac{\nabla f_c(\theta) dx}{f(\theta)} \right] \quad (8)$$

$$= \int_{\mathcal{X}(y)} \frac{\nabla^2 f_c(\theta)}{f(\theta)} dx - \frac{1}{f(\theta)^2} \left(\int_{\mathcal{X}(y)} \nabla f_c(\theta) dx \right) \left(\int_{\mathcal{X}(y)} \nabla f_c(\theta) dx \right)^T \quad (9)$$

$$= \mathbb{E}_\theta \left[\frac{\nabla^2 f_c(\theta)}{f_c(\theta)} \middle| Y = y \right] - \mathbb{E}_\theta \left[\frac{\nabla f_c(\theta)}{f_c(\theta)} \middle| Y = y \right] \mathbb{E}_\theta \left[\frac{\nabla f_c(\theta)}{f_c(\theta)} \middle| Y = y \right]^T \quad (10)$$

$$= \mathbb{E}_\theta \left[\frac{\nabla^2 f_c(\theta)}{f_c(\theta)} \middle| Y = y \right] - S(\theta; y) S(\theta; y)^T \quad (11)$$

$$\nabla^2 \ell_c(\theta) = \nabla \left(\frac{\nabla f_c(\theta)}{f_c(\theta)} \right) \quad (12)$$

$$= \frac{\nabla^2 f_c(\theta)}{f_c(\theta)} - S_c(\theta) S_c(\theta)^T \quad (13)$$

Combining lines 11 and 13, we get

$$\nabla^2 \ell(\theta) = \mathbb{E}_\theta[\nabla^2 \ell_c(\theta) | Y = y] + \mathbb{E}_\theta[S_c(\theta) S_c(\theta)^T | Y = y] - S(\theta; y) S(\theta; y)^T \quad (14)$$

Finally, evaluating line 14 at $\theta = \hat{\theta}$ makes the rightmost term vanish, thereby yielding the required expression. \square

Proposition 1.3 is known as Louis' standard error formula. Other decompositions for the observed information matrix of the observed data likelihood do exist; see, e.g., Oakes (1999); McLachlan and Krishnan (2008). However, the one due to Louis will be most useful to us later.

1.2 Example: Gene Frequency Estimation

Consider the problem of estimating allele frequencies based on observed phenotypes. Often, a single phenotype can be encoded by multiple genotypes with different configurations of dominant and recessive alleles. This is sometimes referred to as the problem of gene frequency estimation. Our analysis closely follows Example 2.4 from McLachlan and Krishnan (2008)¹

We investigate a simplified model for blood type which consists of only the ABO blood group. See Chapter 5 of Dean (2005), for a detailed introduction to blood types. There

¹Although these authors do give data in their example, there is little context around this data, and I have been unable to locate more information from their references. I have opted instead to use my own dataset.

Table 1: Observed frequency and theoretical probability of each blood type (Fujita et al., 1978)

Blood Type	O	A	B	AB
Random Variable	Y_1	Y_2	Y_3	Y_4
Observed Frequency	10	16	7	1
Probability	r^2	$p^2 + 2pr$	$q^2 + 2qr$	$2pq$

are three alleles for this gene: A, B and O. Allele O is recessive, while alleles A and B exhibit co-dominance. That is, genotypes AO and AA encode blood type A, genotypes BO and BB encode blood type B, genotype OO encodes blood type O, and genotype AB encodes blood type AB. Suppose that we seek to estimate the proportion of each allele within a population, based on a sample of individuals' phenotypes. Fujita et al. (1978) report blood types of 4,464,349 people in Japan collected between 1964 and 1975. This sample is so large that any standard errors are practically zero. To retain a reasonable level of uncertainty, we focus on a single administrative division, Oto, in Nara Prefecture. See Figure 1 for details.

Let Y_1 , Y_2 , Y_3 and Y_4 be the number of people with blood type O, A, B and AB respectively, and $Y = (Y_1, Y_2, Y_3, Y_4)$. Let r , p and q be the proportions of alleles O, A and B respectively within the population of interest. Since $r + p + q = 1$, we let $\theta = (p, q)$ be our target of inference. Pretending that the population size is fixed and that iid sampling was employed, Y follows a multinomial distribution with $n = 4,464,349$. Assuming homogeneous genetic mixing, the probability vector for Y is $\pi = (r^2, p^2 + 2pr, q^2 + 2qr, 2pq)$.

Maximizing the likelihood in this model involves solving the score equation, a system of two 3rd-degree polynomials in p and q . This can be done numerically, and gives estimates $p = 0.299$ and $q = 0.128$. These values match the ones given by Fujita et al. (1978). The information matrix and asymptotic covariance (inverse information matrix) are given by

$$I(\hat{\theta}) = \begin{bmatrix} 276 & 84.8 \\ 84.8 & 584 \end{bmatrix} \quad (15)$$

$$\hat{\Sigma}_{\text{MLE}} = \begin{bmatrix} 3.79 \cdot 10^{-3} & -5.49 \cdot 10^{-4} \\ -5.49 \cdot 10^{-4} & 1.79 \cdot 10^{-3} \end{bmatrix} \quad (16)$$

1.2.1 Complete Data

The problem of gene frequency estimation would be much simpler if we could observe individuals' genotypes. We consider augmenting the observed data Y by further classifying individuals by genotype. Let $X = (X_1, \dots, X_6)$ be the genotypes of the individuals represented in Table 1. See Table 2.

Note that we can write Y in terms of X . Specifically, $Y_1 = X_1$, $Y_2 = X_2 + X_3$,

Table 2: Terminology and probabilities for our augmented version of the dataset in Fujita et al. (1978). We also give the blood type coded for be each genotype.

Genotype	OO	AO	AA	BO	BB	AB
Random Variable	X_1	X_2	X_3	X_4	X_5	X_6
Probability	r^2	$2pr$	p^2	$2qr$	q^2	$2pq$
Blood Type	O	A	A	B	B	AB

$Y_3 = X_4 + X_5$ and $Y_4 = X_6$. This corresponds to summing components of X which correspond to the same blood type. The distribution of X is multinomial, with the same sample size as Y , and probability vector given in Table 2.

See Appendix A.2 for the complete data likelihood function and its derivatives.

1.2.2 EM Algorithm

The gene frequency estimation problem fits nicely into the EM algorithm framework. In this section, we present key quantities and results of our analysis. See Appendix A, especially part A.3, for more details.

The EM objective function at iteration $k + 1$ is

$$Q(\theta|\theta_k) \equiv \mathbb{E}_{\theta_k}(n_O|y) \log r + \mathbb{E}_{\theta_k}(n_A|y) \log p + \mathbb{E}_{\theta_k}(n_B|y) \log q \quad (17)$$

$$=: \nu_O^{(k)} \log r + \nu_A^{(k)} \log p + \nu_B^{(k)} \log q \quad (18)$$

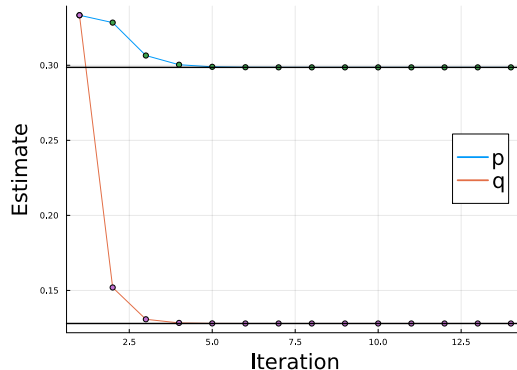
where $\nu_O^{(k)}$, $\nu_A^{(k)}$ and $\nu_B^{(k)}$ are the expected number of O, A and B alleles respectively given $Y = y$ and $\theta = \theta_k$. See Appendix A.3 for explicit formulas. Maximizing this objective function in p and q gives the following EM update function:

$$M(\theta_k) = \begin{pmatrix} p_{k+1} \\ q_{k+1} \end{pmatrix} \quad (19)$$

$$= \begin{pmatrix} \nu_A^{(k)} / 2n \\ \nu_B^{(k)} / 2n \end{pmatrix} \quad (20)$$

Starting with $\theta_0 = (1/3, 1/3)$ corresponding to equal proportions of the three alleles, Figure 1 gives trajectories of the EM estimates for p and q using the data in Figure 1. These estimates converge quite quickly to the maximizer of the observed data likelihood. Beyond computing the observed data MLE, we also need the standard error of this estimator. To this end, we compute the observed data information matrix using Louis' Method (see Proposition 1.3). The asymptotic covariance matrix of our MLE is then the inverse of this information matrix. Omitting details (see Appendices A.2 and A.3), both the observed data information matrix and asymptotic covariance match those obtained from the observed data likelihood.

Figure 1: Trajectory of EM estimates for p and q for the blood type example. Horizontal lines give the values of the MLE.



2 The Monte Carlo EM Algorithm

The Monte Carlo EM, or MCEM, algorithm was first proposed by Wei and Tanner (1990). See also Section 4.5 of Tanner (1996) for a textbook-level treatment of the basic idea with multiple examples. This method proceeds by replacing the conditional expectation in the E-step of the EM algorithm with a Monte Carlo average. More precisely, at each iteration we generate observations from the missing data distribution (i.e. the conditional distribution of the missing data given the observed data), and average the complete data likelihood over this Monte Carlo sample. Formally, at a given iteration of the MCEM algorithm, let X_1, \dots, X_M be a Monte Carlo sample from the law of $X|Y = y$ with θ set to the value from the previous iteration, say θ_0 . Write

$$\hat{Q}(\theta|\theta_0) = \sum_{i=1}^M w_i \ell_c(\theta; y, X_i) \quad (21)$$

where the w_i are sampling weights. We write $\hat{\mathbb{E}}\phi$ for the weighted average of a function ϕ based on an importance sample, so $\hat{Q}(\theta|\theta_0)$ can be re-written as $\hat{\mathbb{E}}\ell_c(\theta; y, X)$. Under iid sampling we simply get $w_i = M^{-1}$ for every i , but more intricate sampling schemes may have more complicated weights. In this section, we focus only on iid sampling. See Section 5 for discussions of some alternative sampling methods. The estimate of θ is then the maximizer of the MCEM objective function: $\hat{\theta} = \operatorname{argmax}_{\theta} \hat{Q}(\theta|\theta_0)$. Write $\hat{\theta}_{k-1}$ for the k th MCEM estimate.

Provided that the MCEM algorithm converges to a critical point of the observed data likelihood, we can use Proposition 1.3 to estimate the observed data information matrix. Specifically, after declaring convergence, we generate a new Monte Carlo sample and use

it to approximate all conditional expectations in Equation (7).

The MCEM algorithm has the advantage of circumventing the challenge of computing potentially intractable conditional expectations for the EM algorithm. However, this analytical simplification does come at the cost of introducing some new computational problems. In this section, we outline the main problems faced by the MCEM algorithm and present various solutions which have been proposed in the literature. We focus primarily on practical aspects of the MCEM algorithm; see Fort and Moulines (2003) for a thorough analysis of the convergence properties of MCEM, and Neath (2013) for a survey of this and other theoretical considerations.

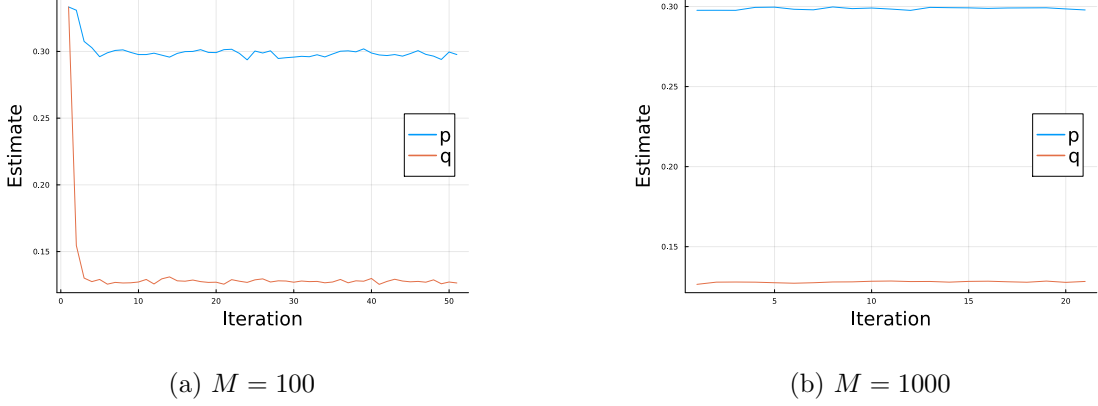
Two problems which have received considerable attention in the literature are how to choose the Monte Carlo sample size at each iteration, and when to terminate the MCEM algorithm. These were identified early by Wei and Tanner (1990), but did not receive systematic treatment until later. We here give a brief overview of different authors' approaches to solving these two problems. The rest of this section goes into more detail on each of the methods. Wei and Tanner (1990) suggest examining a plot of the parameter estimates across iterations, and either terminating or increasing the Monte Carlo size when the plot appears to stabilize. Chan and Ledolter (1995) use a pilot study to choose the Monte Carlo sample size, and terminate when a confidence interval for the improvement of the observed data log-likelihood between successive iterations contains zero. Booth and Hobert (1999) frame each MCEM iteration as an M-estimation problem targeting the deterministic EM update. They increase the Monte Carlo size if an asymptotic confidence interval for the EM update contains the previous iteration's parameter estimate, and terminate when multiple consecutive iterations' estimates have sufficiently small relative error. Caffo et al. (2005) build confidence bounds for the increment in the EM objective function at each iteration of the MCEM algorithm. They increase the Monte Carlo size until the lower bound is positive and terminate when the upper bound is sufficiently small.

In the rest of this section, we give more detail on each of the implementations introduced above. We also illustrate each method on the blood type dataset described in Section 1.2. The relevant conditional distribution and likelihood calculations are described in Appendix A.3.

2.1 Early Work (Wei and Tanner, 1990)

In their seminal work, Wei and Tanner (1990) propose the MCEM algorithm and present a simple implementation. They illustrate that the complete data gradient and Hessian are easily obtained at each iteration from the Monte Carlo sample and, following Louis (1982), give an estimator for the observed data information matrix. Regarding convergence, Wei and Tanner recommend plotting the parameter estimates across iterations and stopping when the estimates appear to stabilize around some constant. When this stabilization is detected, one can either declare convergence and stop, or increase the Monte Carlo size and continue iterating until the estimate trajectory again stabilizes.

Figure 2: Trajectory of MCEM estimates of p and q for the blood type example. The horizontal lines correspond to maximum likelihood estimates.



In order to apply the MCEM algorithm to estimate allele frequencies in the blood type problem, we must specify the number of iterations, K , and the Monte Carlo size for each iteration, M . Starting conservatively, we use $K = 50$ and $M = 100$. Figure 2a gives trajectories of the MCEM estimates of p and q . These estimates appear to converge quickly to a stationary mean, but there is still some uncertainty around this mean. As such, we run MCEM for another 20 iterations with $M = 1000$, starting with the final value from our first run. See Figure 2b. The trajectories from our second run are much more stable around their means. We use the final values from these trajectories as our estimates: $\hat{p} = 0.298$ and $\hat{q} = 0.128$. These values closely match the maximizer of the observed data likelihood.

2.2 Running a Pilot Study (Chan and Ledolter, 1995)

Building on the ideas of Wei and Tanner, Chan and Ledolter (1995) develop a method for both choosing the Monte Carlo size and deciding when to terminate the MCEM algorithm. The method of Chan and Ledolter includes numerous choices for which they do not give specific guidance. We describe the procedure in general terms, but details of any particular implementation will need to be explored in the context of the dataset being investigated.

The algorithm presented by Chan and Ledolter is based on an identity which allows us to estimate observed data likelihood ratios by Monte Carlo averages of complete data likelihood ratios. More precisely, we write

$$\frac{\mathcal{L}(\theta_1; y)}{\mathcal{L}(\theta_2; y)} = \mathbb{E}_{\theta_2} \left[\frac{\mathcal{L}_c(\theta_1; y, X)}{\mathcal{L}_c(\theta_2; y, X)} \middle| Y = y \right] \quad (22)$$

See Chan and Ledolter (1995) for a derivation. To apply Equation 22 to the MCEM algorithm, we replace the conditional expectation on the right-hand side with a Monte

Carlo average from the corresponding conditional distribution. This adjustment allows us to estimate log-likelihood ratios from the observed data distribution, without ever directly evaluating the observed data likelihood.

The method of Chan and Ledolter can be divided into two parts. The first part is a pilot study, in which we compute a standard error for our log-likelihood ratio estimator near the MLE, and determine what Monte Carlo size is required to get this standard error below a pre-specified threshold. In the second part, we increase the Monte Carlo size appropriately, and continue iterating until a confidence interval for the true observed data log-likelihood ratio contains zero. This two-part procedure reflects the suggestion of Wei and Tanner (1990) to run MCEM until it appears to stabilize, then increase the Monte Carlo sample size to get a more precise estimate.

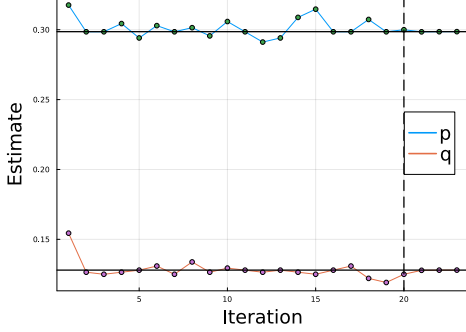
The pilot study portion of Chan and Ledolter’s method proceeds with a fixed, “moderately large” Monte Carlo size. In addition to tracking the parameter estimates across iterations, we also record the estimated log-likelihood ratio of the current estimate relative to the starting point of the algorithm. This ratio is computed by keeping a running cumulative sum of all one-step log-likelihood ratios.

Note that the observed data log-likelihood ratio is estimated by a Monte Carlo average. In order to avoid bias, we do not use the same Monte Carlo sample to estimate this ratio as we used to compute the current parameter update. Instead, we generate the Monte Carlo sample which will be used to compute the next iteration’s parameter update, and use this new sample to estimate the log-likelihood ratio². Since Monte Carlo variability is independent across the two iterations’ samples (something, something, conditional independence), we remove any possible bias from our estimate of the log-likelihood ratio. We are then free to use this new Monte Carlo sample to compute the next iteration’s parameter update.

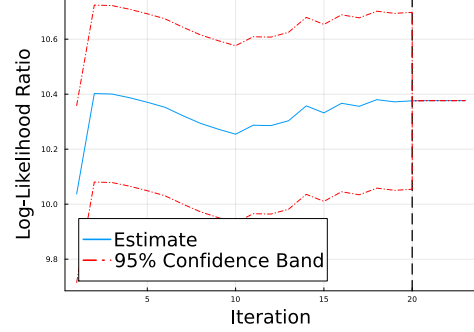
Upon completion of the pilot study (i.e., after a pre-specified number of iterations), we select the estimate with maximal estimated observed data likelihood. We also select a few nearby estimates (Chan and Ledolter suggest the 10 which follow the maximizer but give no justification for this number), and perform a small number of single-iteration MCEM runs from each of the chosen estimates. We compute the estimated observed data log-likelihood ratio for each of these MCEM updates, and get the variance of these updates for each parameter estimate. We then pool variances across parameter estimates into a standard error for our log-likelihood ratio estimator. Finally, we use this standard error, along with its known scaling rate with the Monte Carlo sample size³, to determine what size is needed to get our standard error below the pre-specified threshold. Note that, if a more complicated sampling scheme is being used (e.g. importance sampling, see Section 5.1), then the formula given by Chan and Ledolter for the new Monte Carlo size may not

²Due to our Monte Carlo sample using the new parameter estimate rather than the old, we must actually estimate the reciprocal of the likelihood ratio which we want, then multiply its log by -1 . This is reflected in the formulas of Chan and Ledolter but not discussed explicitly.

Figure 3: MCEM estimates of p and q for the blood type example, based on the method of Chan and Ledolter (1995). The vertical dashed line shows the end of the pilot study. The horizontal lines correspond to maximum likelihood estimates.



(a) Estimate trajectories



(b) Estimated log-likelihood ratio trajectory

apply.

With this newly determined Monte Carlo size, we return to the optimal parameter estimate from our pilot study and continue iterating the MCEM algorithm. At each step now, we also construct a confidence interval for the true observed data log-likelihood ratio corresponding to the current parameter update (i.e., between two consecutive parameter values, not from the starting point to the current estimate). We terminate the algorithm when such a confidence interval contains zero. This corresponds to no evidence of an improvement in the observed data likelihood.

Applying Chan and Ledolter's method to our blood-type dataset, we get the parameter estimate trajectory shown in Figure 3a. Figure 3b gives the trajectory of estimated observed data log-likelihood ratios relative to the starting point of the algorithm ($\hat{p}_0 = \hat{q}_0 = 1/3$), along with pointwise 95% Wald-type confidence bands using the standard error computed after the pilot study and scaled appropriately for the Monte Carlo sample size. We require that the standard error of our estimated log-likelihood increment be at most 10^{-6} . Note that this standard error only applies near the maximizer of the observed data MLE.

2.3 Uncertainty Quantification for the Parameter Estimate (Booth and Hobert, 1999)

Booth and Hobert (1999) use a somewhat different approach from either Wei and Tanner (1990) or Chan and Ledolter (1995) to understand the behaviour of the MCEM algorithm.

³Provided that the parameter estimates used to compute our standard error are close enough to the MLE, Chan and Ledolter show that the standard error scales like $1/M$ rather than the usual $1/\sqrt{M}$, where M is the Monte Carlo size.

The method of Booth and Hobert is based on quantifying Monte Carlo uncertainty of the MCEM update as an approximation to the update which would be made by the deterministic EM algorithm from the same starting point. They then recommend starting the MCEM algorithm with a small Monte Carlo size, and adding more observations only when the parameter estimates are no longer changing discernibly across iterations. More formally, Booth and Hobert suggest building a confidence interval for the EM update based on the Monte Carlo variability of the MCEM update at each iteration. If this interval contains the previous iteration's parameter estimate, then the parameter updates are too small relative to the amount of Monte Carlo variability and more samples are required. These authors similarly recommend assessing convergence by checking for small relative error in the parameter updates. To account for the possibility of Monte Carlo variability leading to two consecutive estimates being similar before the algorithm has 'converged', they suggest waiting until the relative error is small for three consecutive iterations.

The confidence interval used to quantify Monte Carlo uncertainty within an iteration is obtained by framing each step of the MCEM algorithm as the solution of an M-estimation problem. This allows us to inherit the desirable properties of M-estimators; specifically, asymptotic normality (see, e.g. van der Vaart, 1998). Following the usual M-estimator construction and assuming that the relevant regularity conditions hold, we are able to estimate the asymptotic variance of the MCEM parameter estimator at each iteration. Note that this standard error is based on the Monte Carlo variability within an iteration; it does not measure sampling variability due to the observed data.

More formally, write $\tilde{\theta}_k$ for the EM update based on $\hat{\theta}_{k-1}$. Note that $\hat{\theta}_{k-1}$ is held fixed. Analysis of an MCEM update is done conditional on the previous iteration's estimate. Unless stated otherwise, all expectations are taken with $\theta = \hat{\theta}_{k-1}$. Assuming sufficient smoothness and moment conditions, we get the following expression for the MCEM update:

$$\sqrt{M_k}(\hat{\theta}_k - \tilde{\theta}_k) = -\sqrt{M_k} \left[\nabla^2 Q(\tilde{\theta}_k | \hat{\theta}_{k-1}) \right]^{-1} \left[\nabla \hat{Q}(\tilde{\theta}_k | \hat{\theta}_{k-1}) \right] + o_p(1) \quad (23)$$

as $M_k \rightarrow \infty$, where M_k is the Monte Carlo size used to compute $\hat{\theta}_k$ and ∇ denotes differentiation with respect to the left argument of Q or \hat{Q} . Note that the first expression on the right-hand side is the inverse Hessian of the EM objective function (fixed) while the second is the gradient of the MCEM objective function (an average). Thus, $\hat{\theta}_k$ is asymptotically normal with asymptotic variance

$$\left[\nabla^2 Q(\tilde{\theta}_k | \hat{\theta}_{k-1}) \right]^{-1} \mathbb{V} \left[S_c(\tilde{\theta}_k) | Y = y \right] \left[\nabla^2 Q(\tilde{\theta}_k | \hat{\theta}_{k-1}) \right]^{-1} \quad (24)$$

$$\approx \left[\nabla^2 \hat{Q}(\hat{\theta}_k | \hat{\theta}_{k-1}) \right]^{-1} \hat{\mathbb{E}} \left[S_c(\hat{\theta}_k) S_c(\hat{\theta}_k)^T | Y = y \right] \left[\nabla^2 \hat{Q}(\hat{\theta}_k | \hat{\theta}_{k-1}) \right]^{-1} \quad (25)$$

where S_c is the complete data score vector, and $\hat{\mathbb{E}}$ is the Monte Carlo average over the missing data with $\hat{\theta}_k$ held fixed. Note that there is no first moment term in the conditional variance of S_c because $\hat{\theta}_k$ is a maximizer of $\hat{\mathbb{E}}[\ell_c | Y = y]$.

Based on the above discussion, we can build an asymptotic confidence interval for $\tilde{\theta}_k$, the EM update based on the MCEM estimate from iteration k . Booth and Hobert recommend checking whether this interval contains $\hat{\theta}_{k-1}$ and, if so, increasing the Monte Carlo size for the next iteration. Specifically, they suggest starting the next iteration with $M_{k+1} = M_k(1 + 1/r)$, with $r = 3, 4$ or 5 working well in their examples.

To assess convergence of the MCEM algorithm, Booth and Hobert present two criteria. The first is a familiar measure of relative error in parameter estimates between consecutive iterations:

$$\max_j \left(\frac{|\hat{\theta}_{k,j} - \hat{\theta}_{k-1,j}|}{|\hat{\theta}_{k-1,j}| + \delta_1} \right) < \delta_2 \quad (26)$$

where δ_1 and δ_2 are small positive constants, and the subscript j ranges over components of θ . Booth and Hobert suggest using $\delta_1 = 10^{-3}$ and δ_2 between $2 \cdot 10^{-3}$ and $5 \cdot 10^{-3}$. See p. 436 of (Sea06) and the references therein for a discussion of the form of Equation 26.

Alternatively, since Booth and Hobert apply their method to the analysis of generalized linear mixed models, where pathologies may arise due to parameter estimates being too close to a boundary, they propose a second stopping rule:

$$\max_j \left(\frac{|\hat{\theta}_{k,j} - \hat{\theta}_{k-1,j}|}{\sqrt{\nabla \hat{\theta}_{k-1,j}} + \delta'_1} \right) < \delta'_2 \quad (27)$$

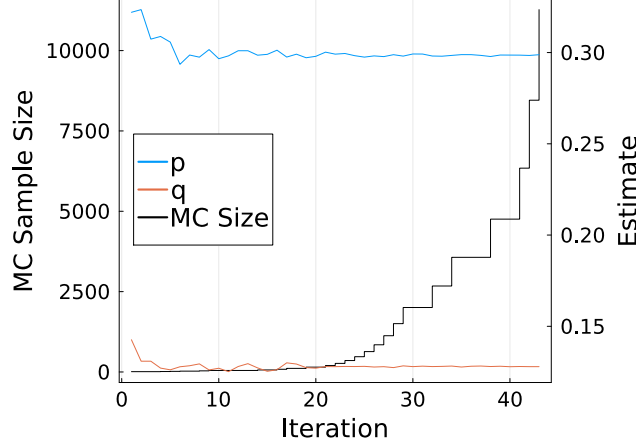
The purpose of condition (27) is to detect when estimated variance components are very close to zero and the numerical precision needed to satisfy condition (26) requires a prohibitive amount of computation.

We apply the method of Booth and Hobert (1999) to our blood type example, with the settings recommended in their paper. Specifically, they suggest setting $\alpha = 0.25$, $k = 3$ (alternatively, 4 or 5), $\delta_1 = 0.001$, and $\delta_2 = 0.002$ (or as high as 0.005). We also start with a Monte Carlo size of 10. Figure 4 gives trajectories of the MCEM estimates, as well as the Monte Carlo size used to obtain each of these estimates. Note how the trajectory stabilizes as the MC size increases. The final estimate from this method is $\hat{p} = 0.299$, $\hat{q} = 0.128$; again, very close to the MLE.

2.4 Uncertainty Quantification for the Objective Function (Caffo et al., 2005)

The approach of Caffo et al. (2005) is similar in spirit to that of Booth and Hobert (1999). Both sets of authors seek to quantify Monte Carlo uncertainty in the MCEM algorithm as an approximation to the EM algorithm. The difference is that where Booth and Hobert measure uncertainty in the parameter estimate, Caffo et al. focus on uncertainty in the

Figure 4: Trajectory of estimates for p and q , as well as Monte Carlo sample sizes, from the method of Booth and Hobert.



objective function. Specifically, Caffo et al. base their analysis on asymptotic normality of the MCEM increment:

Proposition 2.1. *Let $\Delta\hat{Q}(\hat{\theta}_k|\hat{\theta}_{k-1}) = \hat{Q}(\hat{\theta}_{k-1}|\hat{\theta}_{k-1}) - \hat{Q}(\hat{\theta}_k|\hat{\theta}_{k-1})$. Define $\Delta Q(\hat{\theta}_k|\hat{\theta}_{k-1})$ similarly. Let M_k be the Monte Carlo size at iteration k . Then, **assuming some regularity conditions**,*

$$\sqrt{M_k} \left[\Delta\hat{Q}(\hat{\theta}_k|\hat{\theta}_{k-1}) - \Delta Q(\hat{\theta}_k|\hat{\theta}_{k-1}) \right] \rightsquigarrow N(0, \Sigma_k) \quad (28)$$

as $M_k \rightarrow \infty$, where Σ_k is an asymptotic covariance matrix.

See Caffo et al. for a proof. Provided that we are able to estimate Σ_k , Proposition 2.1 allows us to build asymptotic confidence intervals for the EM increment, ΔQ . Recall that in Section 1.1.1, we defined the Generalized EM algorithm by requiring that $\Delta Q \geq 0$, and showed that this requirement guarantees the ascent property. While the stochastic nature of the MCEM algorithm makes it impossible to guarantee that the EM increment is positive, we are able to use Proposition 2.1 to construct asymptotic confidence bounds for ΔQ . Provided that we can estimate Σ_k , we can then be reasonably confident that $\Delta Q > 0$.

Estimating the asymptotic variance under iid or rejection sampling is fairly straightforward. Importance sampling however, is somewhat more complicated; particularly when a normalizing constant must be estimated. Caffo et al. give a formula for importance sampling based on the Delta Method. They also give some guidance for calculating standard errors based on Markov chain Monte Carlo sampling, which we do not go into here.

See Section 5.1 and 5.2 for more details on importance sampling and Markov chain Monte Carlo respectively.

We now return to the key MCEM problems of choosing the Monte Carlo size and when to terminate. For the former, Caffo et al. advise constructing a lower confidence limit for the EM increment, ΔQ . If this limit is positive, then we proceed to the next iteration. If not, then we augment the Monte Carlo sample at the current iteration (with, say, M_k/r , with r some small positive integer as in Booth and Hobert, 1999), and compute a new confidence bound. At the next iteration, Caffo et al. advise using a starting Monte Carlo sample which is at least as large as the final sample from the previous iteration. In fact, a larger sample may be required based on extrapolating the MC variability from the previous iteration. The paper gives a formula for this updated size based on a normal approximation to increments in the MCEM objective function.

Caffo et al. base their termination criterion on stopping when there is evidence that the algorithm is no longer yielding sufficient improvement in the EM objective function. Specifically, they start by choosing a tolerance, $\tau > 0$, then calculate an upper confidence limit for the EM increment at each iteration. If this upper confidence limit is below τ , then we declare that there is little room for improvement left in the EM objective, and terminate our algorithm.

Still need to apply this method to the blood type dataset.

3 Alternatives to the MCEM Algorithm

In this section, we outline some alternatives to the MCEM algorithm for maximizing the likelihood of an incomplete dataset. Examples include stochastic approximation (Lai, 2003), the Monte Carlo maximum likelihood method (Gelfand and Carlin, 1993; Geyer, 1994), and variational inference (Blei et al., 2017; Tsikas et al., 2008).

I haven't yet applied any of these methods on the blood type dataset. I wasn't planning on doing so, but I actually don't think it will be too hard. I've already got most of the infrastructure code written and most of these methods are not very complicated. The possible exception is variational inference, but I've never really tried to implement that one. Maybe it won't be too bad.

3.1 Stochastic Approximation

Stochastic approximation (SA) is a method originally proposed by Robbins and Monro (1951) for finding roots of functions which can only be evaluated with noise. This method was expanded upon rapidly by, for example, Kiefer and Wolfowitz (1952) into a method for derivative-free optimization, and by Dvoretzky (1956) with a systematic theoretical framework. Since the mid-20th century, SA methods have grown into a thriving research area. See, e.g., Kushner and Yin (1997) or Borkar (2022) for textbook-length treatments, and Lai (2003) for a survey paper which focuses on applications to statistics.

The core idea of SA is to iteratively update our estimate of the root, θ_* , of some unobservable function, ϕ , based on the value of a noisy realization of that function at the current estimate. Specifically, if θ_k is our estimate at iteration k and $\hat{\phi} \approx \phi$, then our estimate at iteration $k+1$ is $\theta_{k+1} = \theta_k - \alpha_k \hat{\phi}(\theta_k)$, where α_k is a sequence which goes to zero at a particular rate. Since our sequence of weights goes to zero with k , the update terms become negligible in the limit and our estimate of θ_* stabilizes. The precise requirement for these weights is that $\sum_{k=1}^{\infty} \alpha_k = \infty$ and $\sum_{k=1}^{\infty} \alpha_k^2 < \infty$. A common choice is $\alpha_k = k^{-1}$. Numerous authors have studied convergence of the SA method in probability, almost surely and in \mathcal{L}^1 under various regularity conditions. See Lai (2003) for an excellent review of the history of stochastic approximation convergence theory.

There are a few specific applications of the general stochastic approximation which merit special attention. First, suppose that we want to optimize a function, f , which we can only evaluate approximately. Assume further that we are able to approximately evaluate the gradient of f , ∇f . Setting $\phi = \nabla f$ and running the SA algorithm with $\hat{\phi} \approx \nabla f$ gives us a critical point for f . We can therefore use SA to optimize functions which cannot be evaluated exactly. This version of SA sees considerable use in the Machine Learning community under the name Stochastic Gradient Descent (Bottou, 2010).

Another specific instance of the SA method is when computing fixed points. For example, suppose that we want to find a point, θ_* , such that $f(\theta_*) = \theta_*$. In this case, we set $\phi(\theta) = f(\theta) - \theta$ and run SA. The update equation in this case becomes a convex combination of the previous iteration's estimate and the current function value: $\theta_{k+1} = (1 - \alpha_k)\theta_k + \alpha_k \hat{f}(\theta_k)$.

We now discuss application of stochastic approximation to the missing data problem. We refer to such an application as a stochastic approximation EM (SAEM) algorithm. Both cases of the SA method described above find use in constructing particular SAEM algorithms. An advantage of the SAEM algorithm over MCEM is that in SAEM we choose the Monte Carlo size once at the beginning and leave it fixed for every iteration. We can think of the method as automatically increasing the MC size since the estimate at each iteration is a weighted sum of all the estimates which came before it. A disadvantage of SAEM is that it requires us to select the sequence $\{\alpha_k\}$, commonly referred to as the “step size”. Choosing α_k too large will mean the algorithm takes a long time to stabilize, while choosing α_k too small causes the algorithm to stabilize before it reaches its limiting value (and will therefore take a long time to reach this limit). Jank (2006) gives some guidance on choosing this step size based on the goal of balancing bias with variance. This work also presents a convergence diagnostic based on the ideas of Caffo et al. (2005) which allows for a posteriori assessment of whether the step size was too small.

A particular implementation of the SAEM method by Gu and Li (1998) suggests setting ϕ to the observed data score, $\phi(\theta) = S(\theta; y)$. This function can be estimated using Proposition 1.2 (i) and the same Monte Carlo approach we use with MCEM. In the case of a vector-valued parameter, Gu and Li also recommend pre-multiplying $\hat{\phi}$ by a matrix which converges to the inverse of the observed data information matrix. Such a sequence

of matrices can be constructed using Louis’ Identity (Proposition 1.3) and Monte Carlo. Note that we use the same sample of missing data to update our estimates of the parameter and the observed data information matrix at each iteration. Similar work by Gu and Kong (1998) extend the same SAEM construction to accommodate Markov chain Monte Carlo sampling. Gu and Zhu (2001) discuss how to apply the above methodology to the analysis of spatial models, and incorporate a second stage to the method in which estimates are averaged across iterations. This averaging has been found to improve performance of SAEM, and of stochastic approximation more generally (Delyon et al., 1999; Polyak and Juditsky, 1992).

Jank (2006) presents an SAEM implementation in which $\phi(\theta)$ is the EM update function (see Section 1) minus θ (see also Delyon et al., 1999). Thus, the limit of such an algorithm is a fixed-point of the EM update. This ties convergence of Jank’s method more closely to the convergence theory of the original EM algorithm (Wu, 1983). In fact, Jank leverages this connection to exploit the relationship between the EM algorithm’s convergence rate and certain information matrices (Meng and Rubin, 1994). This connection allows Jank to estimate the bias and variance of the SAEM estimates, which are then used to calibrate the step size, α_k .

A subtly different line of research on the SAEM method has been developed by a group at the National Institute for Research in Digital Science and Technology (INRIA) in France (see, e.g., Celeux et al., 1995, for a comparison of some of their methods). The goal here is to augment the EM algorithm, rather than to facilitate the application of EM-type methods when ordinary EM is intractable. More precisely, methods from this group introduce a stochastic perturbation to the EM algorithm, with the goal of escaping fixed points which are locally, but not globally, optimal. Early work centered around a Stochastic EM (SEM) algorithm, which is equivalent to the MCEM algorithm with a Monte Carlo size of one (Celeux and Diebolt, 1987; Celeux et al., 1995; see also Nielsen, 2000). Later, they also propose a method which they refer to as SAEM (although it does not quite fit into our framework), in which each iteration consists of first computing both the EM and MCEM updates from the previous iteration’s estimate, then combining these two updates in a convex combination as the estimate for the current iteration. Here, as with the SEM algorithm, the MCEM update is computed with a Monte Carlo size of one (Celeux and Diebolt, 1992; Celeux et al., 1995).

3.2 Monte Carlo Maximum Likelihood

The Monte Carlo maximum likelihood (MCML) method was developed to handle settings when the likelihood function cannot be evaluated exactly. The idea is to approximate the whole likelihood surface (up to an additive constant) using a single Monte Carlo sample. This approximate likelihood is then maximized, either numerically or analytically. The core idea of MCML is to choose a fixed reference parameter value, θ_* , and estimate likelihood ratios relative to θ_* . The likelihood ratio is written as an expectation with respect to

the fixed reference value, and this expectation is approximated by Monte Carlo. A key feature of this methodology is that a single Monte Carlo sample can be re-used to estimate the likelihood ratio at any number of target parameter values. Finally, we maximize our estimated likelihood ratio as a proxy of the unknown likelihood function.

The most basic form of Monte Carlo maximum likelihood (Geyer, 1991) applies when we only know the likelihood up to a normalizing constant, say $f(y; \theta) = h(y; \theta)/c(\theta)$, where h is known but c is not. We note that $\int h(y; \theta) dy = c(\theta)$, and write

$$\log \frac{f(y; \theta)}{f(y; \theta_*)} = \log \frac{h(y; \theta)}{h(y; \theta_*)} - \log \frac{c(\theta)}{c(\theta_*)} \quad (29)$$

$$= \log \frac{h(y; \theta)}{h(y; \theta_*)} - \log \int \frac{f(y; \theta)}{c(\theta_*)} dy \quad (30)$$

$$= \log \frac{h(y; \theta)}{h(y; \theta_*)} - \log \int \frac{h(y; \theta)}{h(y; \theta_*)} f(y; \theta_*) dy \quad (31)$$

$$= \log \frac{h(y; \theta)}{h(y; \theta_*)} - \log \mathbb{E}_{\theta_*} \frac{h(y; \theta)}{h(y; \theta_*)} \quad (32)$$

$$\approx \log \frac{h(y; \theta)}{h(y; \theta_*)} - \log \frac{1}{M} \sum_{i=1}^M \frac{h(y_i; \theta)}{h(y_i; \theta_*)} \quad (33)$$

where the y_i are sampled iid from $f(y; \theta_*)$. The second term in line 33 may need to be modified if an alternative sampling scheme is used.

An alternative formulation of the MCML method given by Gelfand and Carlin (1993) is more relevant for use with missing data (see also, Geyer, 1994). A similar derivation to the one given above shows that, in the presence of missing data, X ,

$$\log \frac{f(y; \theta)}{f(y; \theta_*)} = \log \mathbb{E}_{\theta_*} \left[\frac{h(y, X; \theta)}{h(y, X; \theta_*)} \middle| Y = y \right] - \log \mathbb{E}_{\theta_*} \left[\frac{h(Y, X; \theta)}{h(Y, X; \theta_*)} \right] \quad (34)$$

$$\approx \log \left[\frac{1}{M_1} \sum_{i=1}^{M_1} \frac{h(y, X_i; \theta)}{h(y, X_i; \theta_*)} \right] - \log \left[\frac{1}{M_2} \sum_{j=1}^{M_2} \frac{h(Y_j, X_j; \theta)}{h(Y_j, X_j; \theta_*)} \right] \quad (35)$$

where the X_i are generated from the conditional distribution of $X|Y = y$, the (Y_j, X_j) pairs are generated from the joint distribution of Y and X , and $h(y, x; \theta)$ is proportional to this joint distribution. Note that two Monte Carlo samples are required to evaluate 35.

Note that the MCML procedure closely resembles importance sampling, with $f(y; \theta_*)$, $f_c(y, x; \theta_*)$ or $f_m(y, x; \theta_*)$ being the proposal distribution (see Section 5.1 for a discussion of importance sampling). Jank and Booth (2003) take this idea further, and suggest directly estimating $f(y; \theta)$ from the complete data density, $f_c(y, x; \theta)$, using importance sampling. That is, they write $f(y; \theta) = \int f_c(y, x; \theta) dx = \int [f_c(y, x; \theta)/g(x)] g(x) dx$, where g is an arbitrary density function. The final integral is then approximated by an average of samples drawn iid from g . While the method of Jank and Booth does present an interesting

direction in which to generalize MCML, unfortunately, it is only applicable if f_c is known exactly.

3.3 Variational Methods

This section was written by memory (i.e. without looking anything up). The idea is to more efficiently get words on the page, then I can fix things later.

Variational inference is a set of methods for approximating intractable densities (Citation Needed). The literature on variational inference is vast. See Blei et al. (2017) for a recent survey paper. Gelman et al. (2013, Section 13.7) give a textbook-level overview focusing on applications to Bayesian inference. Numerous authors have discussed the connection between variational inference and the EM algorithm; see, e.g., Neal and Hinton (1998); Tsikas et al. (2008); Blei et al. (2017).

The central idea of variational inference is to frame the target density as the exact solution of a functional optimization problem, where the decision variable typically ranges over densities. The domain of functions is then restricted to some set which is easier to work with. Finally, optimization is performed over this restricted class of functions. The result of this optimization is then our approximation to the target density. **This discussion doesn't make it clear whether our estimator is the optimizer or the optimal objective function.**

More generally, the optimization problem described above may be just one in a sequence of problems which may form part of an iterative algorithm. Herein lies the connection to the EM algorithm. In the M-step of EM, we compute expectations with respect to a particular conditional distribution. To embed this procedure in the variational inference framework, we define an optimization problem whose solution is the same conditional distribution.

We now give some details. First, let q be some plausible marginal probability density for the missing data, X . We can write the observed data log-likelihood as:

$$\ell(\theta; y) = \log f_c(y, x; \theta) - \log f_m(y, x; \theta) \quad (36)$$

$$= \log \left[\frac{f_c(y, x; \theta)}{q(x)} \right] - \log \left[\frac{f_m(y, x; \theta)}{q(x)} \right] \quad (37)$$

$$= \mathbb{Q} \log \left[\frac{f_c(y, X; \theta)}{q(X)} \right] - \mathbb{Q} \log \left[\frac{f_m(y, X; \theta)}{q(X)} \right] \quad (38)$$

$$=: F(q, \theta) + \text{KL}(q \rightarrow f_m(\theta)) \quad (39)$$

$$\geq F(q, \theta) \quad (40)$$

where $\text{KL}(q \rightarrow f_m(\theta))$ is the KL divergence from q to $f_m(\theta)$. Line (38) holds because the left-hand side does not depend on x , and the last line follows from non-negativity of the KL divergence. We refer to the first term in line (39) as the “evidence lower-bound” (ELBO), as well as the “variational free energy”, depending on the field of application (Citation Needed). The latter name comes from an analogue in physics (Citation Needed).

The former name comes from applications to Bayesian inference, where the observed data log-likelihood can be seen as an evidence if we take Y to be the data and X to be the parameters.

The EM algorithm can be re-framed as alternately maximizing F with respect to q and θ . To see why, first recall that the KL divergence between two distributions is non-negative, and is zero if and only if the two distributions are equal. Starting at a fixed parameter value, θ_0 , maximizing F with respect to q is accomplished by setting $q(x) = f_m(y, x; \theta_0)$, as this minimizes the KL divergence in (39) and the left-hand side is constant in q . The resulting ELBO is equal to $Q(\theta|\theta_0) + \xi$, where Q is the EM objective function and ξ is constant in θ . Maximizing the ELBO in θ is therefore equivalent to maximizing the EM objective function. This maximization gives a new parameter value, which is used as input to the next iteration. See Theorem 1 of Neal and Hinton (1998) for a more formal proof.

While the EM algorithm is obtained by maximizing q over an unrestricted class of densities, other procedures can be formulated by restricting this class. One popular example is the “mean-field” approximation, which involves optimizing over the class of densities which factor over their arguments. I.e., the class of functions $\mathcal{Q}_{MF} := \{\text{Densities } q : q(X_1, \dots, X_p) = \prod_{j=1}^p q_j(X_j)\}$.

A major advantage of the mean-field approximation is that an iterative algorithm exists for finding the density, q , which maximizes the ELBO. This algorithm performs coordinate ascent, and the coordinate updates are closely related to computation of the full conditional distributions in Gibbs sampling (Citation Needed). Write $q^{(k)} = \prod q_j^{(k)}$ for the current value of q , and $\mathbb{Q}_{-j}^{(k)}$ for expectation with respect to all the missing variables except j , under $q^{(k)}$. The update formula is

$$q_j^{(k+1)} \propto \exp \left[\mathbb{Q}_{-j}^{(k)} \ell_c(y, x_j, X_{-j}) \right] \quad (41)$$

where X_{-j} is all the missing variables other than X_j . See Section 2.4 of Blei et al. (2017) for a derivation of (41). The overall algorithm consists of repeatedly cycling through updating each coordinate’s distribution until some convergence criterion is met.

Note that, so far, our discussion of how to compute the mean-field approximate density for X has not addressed θ . To apply mean-field variational inference to EM-type problems, we substitute the mean-field density into the ELBO and maximize over θ . This new value of θ is fed back into (41), giving us a different complete data likelihood function and, hence, a new optimal density.

4 Comparisons Between Methods

Several authors have performed comparisons between methods discussed in Sections 2 and 3. In this section, we discuss these comparisons and their findings.

McCulloch (1997) use a simulation study to compares the MCEM and MCML methods, along with a Monte Carlo version of the Newton-Raphson algorithm (MCNR). Their

MCCEM implementation starts with fixed Monte Carlo size, then increases this size at iterations 20 and 40. It is not clear how these jump points were selected, nor how convergence was assessed beyond examining plots. Their MCCEM implementation thus most closely resembles that of Wei and Tanner (1990). A similar schedule of Monte Carlo sizes and termination was used for MCNR, whereas MCML used a much larger Monte Carlo size (the sample size for MCML was not increased). McCulloch also considers MCCEM and MCNR followed by MCML (i.e. using MCCEM or MCNR to choose the reference parameter for MCML). This comparison is made using a logit-normal mixed-effects model with one random effect and one fixed effect. McCulloch found that MCCEM and MCNR perform better than MCML alone, but that following either MCCEM or MCNR with MCML was even better. They did not find that following-up with multiple iterations of MCML was preferable to a single run of MCML.

In addition to presenting their SAEM method, Gu and Li (1998) compare this method with MCCEM. They give a comparative analysis on a dataset of motorette failure times (see Rain et al., 2016, for a diagram and explanation of what a motorette is). The model being fit is a linear regression with right-censoring. Their MCCEM implementation “is from Tanner (1993)”, which I do not have access to. However, based on the third edition of this book (Tanner, 1996), the implementation in Gu and Li is the same as that of Wei and Tanner (1990). It is not clear what Monte Carlo size they use to start, and it appears that this size is ever augmented. Their SAEM implementation is as previously described in their paper, with an MC size of 1, step size at iteration k (i.e., α_k) of $1/k$, and pre-multiplying matrix chosen adaptively based on the current iteration as given by Equation (13) of their paper. Gu and Li find that SAEM converges much more quickly than MCCEM. In fact, based on their Figure 1, it is not clear that MCCEM is converging to the MLE at all. These authors also give a heuristic argument that SAEM should converge much more quickly than MCCEM based on the number of MC samples used at each iteration and the number of iterations required to converge to the MLE. However, this argument is based on a fixed Monte Carlo size at each iteration and is thus not directly relevant to the MCCEM implementations discussed in Section 2.

Booth et al. (2001) compare the MCCEM, SAEM and MCML (referred to as stochastic maximum likelihood, or SML) methods. They perform their comparison on a simple one-way mixed-effects linear model, with known variance component and error variance. Their MCCEM implementation matches that of Booth and Hobert (1999), while their SAEM implementation is that of Jank (2006) (or Deylon et al., 1999). Note that although the presentation of SAEM in Booth et al. appears different from ours, it is not hard to show that the two are equivalent. Their MCML implementation uses the missing data distribution with a fixed value for the unknown parameter as reference distribution. The comparison between these methods is done partly analytically and partly by simulation. The mean squared errors (MSEs) of MCCEM and SAEM for estimating the unknown parameter can be obtained analytically. Booth et al. thus compare these two methods directly and find that MCCEM performs better when the problem is harder (i.e. larger mean and variance

component). The MSE of MCML however must be approximated by simulation. This MSE is approximated using 500 replicates, and a 95% Wald-type confidence interval is constructed. The upper and lower bounds of this confidence interval are then compared to the analytical MSE of MCEM (the authors do not compare MCML with SAEM). The result is that MCML is competitive with MCEM, and that MCML even performs better for some parameter settings (typically, when the variance component is small). Booth et al. point out that the comparison between MCEM and MCML is not entirely fair here, because taking the variance component as known makes simulation for MCML unrealistically easy. In more serious problems, difficulty in choosing a proposal distribution for MCML will likely lead to worse performance.

Booth et al. (2001) also extend the heuristic argument in Gu and Li (1998) to account for Monte Carlo sizes changing with MCEM iteration. They argue that, in the scalar case, MCEM should converge more quickly than SAEM when the “fraction of missing information” is larger than $\exp(-1)$. The fraction of missing information is defined as $\mathcal{I}_c(\hat{\theta})^{-1}\mathcal{I}_m(\hat{\theta})$, where $\hat{\theta}$ is the observed data MLE (see Proposition 1.2 for definitions of \mathcal{I}_m and \mathcal{I}_c). This quantity is closely related to the convergence rate of the EM algorithm (Meng and Rubin, 1994; McLachlan and Krishnan, 2008).

Jank and Booth (2003) extend the work of Booth et al. (2001), specifically the comparison between MCEM and MCML. Jank and Booth focus on analytical comparisons; as such, they use a fixed Monte Carlo size across iterations to make their calculations more tractable. They also use an unrealistic proposal distribution for MCML which requires that we know the observed data MLE. The authors derive the asymptotic variance of MCEM and MCML, and investigate their asymptotic relative efficiency (ARE). It turns out that the ARE depends directly on the eigenvalues of the fraction of missing information defined above. In particular, the relative efficiency of MCEM relative to MCML goes to ∞ as the fraction of missing information goes to 1. That is, as a problem gets harder in the sense that less information is available in the observed data, we expect MCEM to perform better relative to MCML. We also expect MCML to perform worse on real problems, since the above analysis is based on an inaccessible proposal distribution. Jank and Booth illustrate their analytical calculations on the one-way mixed-effects linear model from Booth et al. (2001) and a logistic-normal generalized linear mixed-effects model. The latter consists of a simulation study which compares, among other things, the effect of the choice of proposal distribution on MCML. They found that the average estimates from MCEM and MCML (i.e. averaged over simulation replicates) are fairly consistent, but that variability of the MCML estimates is much higher. In particular, entries in the empirical covariance matrix for MCML grow rapidly as the reference parameter value for the proposal distribution moves away from the MLE.

4.1 Synthesis

Many of the comparisons we have discussed (McCulloch, 1997; Gu and Li, 1998; Jank and Booth, 2003) use unsophisticated implementations of MCEM. In particular, as far as I can tell, Gu and Li (1998) use a poor implementation which never increases the Monte Carlo size, so it is unsurprising that they found MCEM performs poorly. Jank and Booth (2003) also uses a fixed Monte Carlo size, but theirs is sufficiently large that we can expect the behaviour of MCEM to be close to that of the deterministic EM algorithm. Similarly, McCulloch (1997) increase the Monte Carlo size at fixed iterations, and the final size is quite large (5000 for the final ten iterations). Only Booth and Hobert (1999) use an adaptive rule for choosing the Monte Carlo size.

Despite the above concerns about Monte Carlo size, MCEM tends to perform quite well in simulations and analytical comparisons. Only Gu and Li (1998) found that MCEM performed substantially worse than another method, although McCulloch, 1997 did find that following MCEM with MCML tends to improve performance. These findings suggest that the MCEM algorithm should be one of the first methods considered when approaching a missing data problem in which the calculations required to implement EM are intractable.

A limitation of the above comparisons is that all use iid sampling. The concerns raised by Booth et al. (2001) and Jank and Booth (2003) about choosing an efficient proposal distribution for MCML also apply to MCEM unless it is possible to sample directly from the missing data distribution. Indeed, see Ruth and Lockhart (2023) for an analysis in which great care is taken to select an appropriate proposal distribution for MCEM.

5 Simulation

An obstacle to implementing the MCEM algorithm which was not addressed in Section 2, and which is also relevant to many of the methods described in Section 3, is how to generate the necessary Monte Carlo samples. It is in general a hard problem to simulate from arbitrary conditional distributions. In this section, we discuss a few methods for simulating the required observations at each step of the MCEM algorithm. All the topics that we cover here have their own vast bodies of literature, which we cannot hope to cover in their entirety. We instead give only a brief overview, focusing on aspects which are particularly relevant to use with the MCEM algorithm, and direct the reader to other, more focused, reviews.

We start by discussing importance sampling, which is a method for using a sample from one distribution to compute moments under another distribution. Next, we cover Markov chain Monte Carlo sampling, which consists of constructing and sampling from a Markov chain whose stationary distribution matches the distribution from which we want to simulate. Finally, we briefly touch on rejection sampling, a way to simulate exactly from a target distribution at the expense of increased computation time, and sequential Monte Carlo, which generates a sequence of samples such that they converge to the target

distribution.

5.1 Importance Sampling

Broadly speaking, importance sampling is a framework for approximating intractable expectations. A typical use case is when we want to evaluate the expected value of some function, $\mathbb{E}h$, and this expectation is not just analytically intractable, but the distribution \mathbb{F} is impossible (or impractical) to sample from. The latter restriction prevents us from using ordinary Monte Carlo integration. Instead, we re-write $\mathbb{E}h$ as $\mathbb{E}[h \cdot (f/g)]$ for some new distribution, \mathbb{G} , where f and g are the densities of \mathbb{F} and \mathbb{G} respectively. Provided that \mathbb{G} is easy to sample from, we can estimate this alternative expression for our target expectation via Monte Carlo integration with samples drawn from \mathbb{G} . We call \mathbb{F} the target distribution and \mathbb{G} the proposal distribution.

A classic reference on importance sampling and other Monte Carlo methods is the book by Robert and Casella (2004); particularly Chapters 3 and 4. Chapter 8 of the book by Chopin and Papaspiliopoulos (2020) gives a more current overview of importance sampling, with a focus on its application to Sequential Monte Carlo. Elvira and Martino (2022) give a survey of modern methods for extending the importance sampling framework, with an emphasis on two main approaches: multiple importance sampling and adaptive importance sampling. See Elvira et al. (2019) for a review of multiple importance sampling methods, and Bugallo et al. (2017) for more on adaptive importance sampling. Agapiou et al. (2017) give a survey paper level treatment of some more theoretical considerations for importance sampling.

In the rest of this section, we describe a few simple modifications which can ease implementation and improve performance when using importance sampling with the MCEM algorithm.

When using importance sampling with the MCEM algorithm, our target distribution is the missing data distribution (i.e. the conditional distribution of the missing data given the observed data). In some settings, this distribution may be difficult to describe exactly. However, the integrand is the (log-)likelihood for the complete data distribution, so it is reasonable to expect that we can evaluate this complete data density. From the definition of conditional probability, the missing data density is proportional to the complete data density, provided that we treat the latter as a function of the missing data and hold the observed data fixed. The proportionality constant here is the observed data density, so it is unlikely that we will be able to normalize the missing data density exactly (otherwise, we could just work directly with the observed data likelihood).

A simple modification of importance sampling, which circumvents the need for normalized densities, is to compute the importance weights as usual, then normalize them to sum to one. This is referred to as “auto-normalized”, or “self-normalized”, importance sampling (see, e.g., Elvira and Martino, 2022). The reason self-normalized importance sampling works is that the unknown normalizing constant cancels in the numerator and

denominator of our normalized weights. In fact, our proposal distribution can also be un-normalized, and the ratio of the two normalizing constants cancels when we normalize our weights.

There are, however, disadvantages of self-normalized importance sampling compared to the exact importance sampling which we can do when both target and proposal densities are known exactly. One important limitation is that our estimator of $\mathbb{E}h$ is no longer unbiased, and is instead only asymptotically unbiased. In fact, as a ratio estimator, the standard error of a self-normalized importance sampling estimator can only be determined asymptotically. This is fine in problems where it is easy to sample from our proposal distribution, \mathbb{G} , but in high dimensional problems for example, even simulating from \mathbb{G} may be costly, and more care must be taken with the discrepancy between asymptotic results and finite-sample behaviour.

It is well-known that the performance of an importance sampling estimator depends on how closely the proposal distribution matches the target (Robert and Casella, 2004). One thing that can go wrong is if the importance weight (i.e. the likelihood ratio between these two distributions) does not have sufficiently many finite moments (Agapiou et al., 2017). A simple modification to our importance weights which guarantees infinitely many finite moments is to specify a threshold value, and truncate any weights which fall above the threshold (i.e., set weights which fall above the threshold equal to the threshold value). This “truncated importance sampling” method is proposed and analyzed by Ionides (2008). Two strategies are proposed in this work for selecting the threshold: the first is to simply use the square root of the Monte Carlo size, \sqrt{M} , while the latter is based on unbiased risk estimation and gives a value better tailored to the specific problem. Note that these recommendations are based on exact importance sampling. If self-normalization is used, the general recommendation is instead to truncate at \sqrt{M} times the mean of the un-normalized weights. Truncated importance sampling represents an example of the bias-variance trade-off. Truncating weights at some threshold reduces the variance of our importance sampling estimator, but also introduces bias. This trade-off highlights the importance of selecting an appropriate threshold value: too large and the variance reduction will be negligible, but too small and the bias will be unacceptable.

Vehtari et al. (2022) propose an alternative method for handling large importance weights, called Pareto Smoothed Importance Sampling (PSIS). The idea of this method is to fit a Generalized Pareto Distribution to the largest importance weights, then replace those large weights with quantiles of the fitted distribution. One of the parameter’s fitted values also serves as a useful diagnostic for how closely our proposal distribution matches the target. An advantage of PSIS over the truncated importance sampling method of Ionides (2008) is that the bias here is smaller, although it is not clear in general which method has better mean-squared error (see Vehtari et al., 2022, for extensive numerical comparisons).

An approach proposed by Levine and Casella (2001) seeks to use importance sampling to save computing time when running the MCEM algorithm. This computational efficiency

is especially important in their context where Markov chain Monte Carlo sampling is used (see 5.2), so generating a single sample takes quite some time. Their idea is to generate a single sample from the missing data distribution with some reference value for θ , θ_{ref} . This sample is then re-used at every MCEM iteration, with conditional expectations under the current parameter value computed by taking importance ratios with respect to θ_{ref} . This raises the question of whether a single proposal distribution can be adequate for every target distribution along the MCEM trajectory. To address this concern, Levine and Casella suggest running MCEM for a few iterations with fresh importance samples and only drawing the sample which will be used for their method after a sufficiently long “burn-in” period (they suggest running for one minute).

5.2 Markov Chain Monte Carlo

The core idea of Markov chain Monte Carlo (MCMC) sampling is to construct a Markov chain from which we can simulate, and which has stationary distribution equal to the target distribution. Popular methods to construct such a Markov chain are Gibbs sampling and the Metropolis-Hastings algorithm. See Gelman et al. (2013) or Robert and Casella (2004) for excellent textbook-length overviews. A popular implementation of MCMC sampling is the **Stan** programming language (Stan Development Team, 2022), and its **R** interface, **RStan** (Stan Development Team, 2023).

Both Gibbs sampling and Metropolis-Hastings start with a random variable, $X = (X_1, \dots, X_d)$, which we wish to simulate. Let f be the density of X . These methods proceed by iteratively simulating draws of the vector X from some distribution based on the draw from the previous iteration.

Gibbs sampling is based on successively sampling each component of X conditional on all the other components. The order in which this conditioning is performed is a bit subtle, however. When generating X_i , we condition on the values of X_1, \dots, X_{i-1} from the current iteration and the values of X_{i+1}, \dots, X_d from the previous iteration. Once we reach the end of X , we start a new iteration. More generally, we can group the components of X , and simulate an entire group conditional on the others.

The Metropolis-Hastings algorithm closely resembles rejection sampling. We start each iteration by generating a candidate value of X from some proposal distribution (this distribution may depend on the previous iteration’s value of X). Write $J(x|x_0)$ for the proposal density, where x_0 is the value of X from the previous iteration. Next, we define an acceptance probability, $r := f(x)J(x_0|x)/f(x_0)J(x|x_0)$, and accept the proposed value of X with probability $r \wedge 1$. With probability $(1 - r) \vee 0$, we reject the proposed x and instead set the current iteration’s value to x_0 . We then proceed to the next iteration. Note that rejecting still adds an observation to our Monte Carlo sample, this value just happens to be identical to the one proceeding it.

Limit theory for estimators based on MCMC sampling are more complicated than the similar theory for estimators based on iid or importance sampling (Geyer, 1991). This

additional complexity stems from the dependence between draws from an MCMC sampler. While convergence for iid and importance sampling can be established using the Law of Large Numbers and the classical Central Limit Theorem, similar results for MCMC sampling make use of the Ergodic Theorem and Markov chain Central Limit Theorem. A challenge in implementing the Markov chain Central Limit Theorem is that the asymptotic variance depends on pairwise covariances between points in the chain with arbitrarily large lags. Estimation of this asymptotic variance is therefore challenging in practice. See, e.g., Chapters 6 and 7 of Robert and Casella (2004).

Levine and Casella (2001) propose a method to simplify the analysis of estimates based on MCMC sampling. Their approach consists of subsampling the original chain using Poisson spacings with increasing mean. The result is that a target function averaged over the subsampled chain satisfies a classical Central Limit Theorem (i.e. one with no covariance terms in the asymptotic variance). We can estimate the mean and standard error of our subsample estimator using the entire chain and construct confidence intervals for the mean of the target function. Levine and Casella apply this strategy for constructing confidence intervals along the same lines as Booth and Hobert (1999). Here, we construct a confidence interval for the gradient in the EM objective function, and increase the Monte Carlo size if the confidence interval for the current iteration contains the subsample estimate from the previous iteration. This work is modified to provide uncertainty quantification for the MCEM update instead of the EM score, and also to incorporate MCMC sampling, by Levine and Fan (2004). Furthermore, Levine and Fan give a principled argument for how much to increase the Monte Carlo sample size by when doing so is deemed necessary.

5.3 Other Sampling Methods

While importance sampling and MCMC are the most discussed sampling schemes in the context of the MCEM algorithm, some others do exist. Rejection sampling and sequential Monte Carlo are two such alternatives. Rejection sampling closely resembles importance sampling, except instead of weighting each proposal by the likelihood ratio, we either accept or reject the proposed observation with probability proportional to the likelihood ratio (Chopin and Papaspiliopoulos, 2020). Typically, rejection sampling is continued until the number of accepted proposals reaches a desired sample size. These accepted points are then treated as an iid sample from the target distribution. Although the output of rejection sampling sounds ideal (much of the difficulty with importance sampling comes from having to account for simulated points not having been drawn from the target distribution), the cost comes in computation time. The number of draws from the proposal required to get a fixed number of accepted draws is random, and can be quite high if the proposal distribution does not closely match the target. Indeed, if we instead fix the number of draws from the proposal distribution, importance sampling can be shown to have lower variance than the associated rejection sampling scheme (see Section 8.8 of Chopin and Papaspiliopoulos, 2020).

Sequential Monte Carlo (SMC), is a form of adaptive sampling in which a sequence of samples is generated such that the distribution of these samples converges to some target. The update from one sample to the next seeks to balance preserving the effective sample size with maintaining computational efficiency. See Del Moral et al. (2006) for a survey paper, or Chopin and Papaspiliopoulos (2020) for a book-length overview of SMC. Moffa and Kuipers (2014) use SMC in an MCEM analysis to sample from a truncated multivariate normal distribution with difficult truncation domain.

Appendix A Likelihood for Gene Frequency Estimation

In this appendix, we present details for the analysis of our example of estimating gene frequency. See Section 1.2 for formulation of the model and definition of notation.

A.1 Observed Data Likelihood, Score and Information

Let π_i be the probability of blood type i . The observed data likelihood for our model can be written as follows:

$$\ell(\theta; y) = \log \binom{n}{y} + \sum y_i \log \pi_i(\theta) \quad (42)$$

$$\equiv \sum y_i \log \pi_i \quad (43)$$

$$\equiv 2y_1 \log r + y_2 \log(p^2 + 2pr) + y_3 \log(q^2 + 2qr) + y_4 \log pq \quad (44)$$

where we use \equiv to denote equality up to additive constants which do not depend on θ .

Differentiating ℓ wrt θ and recalling that $r = 1 - p - q$, so $\partial_p r = \partial_q r = -1$, we get the following expression for the observed data score, S .

$$S(\theta; y) = \begin{pmatrix} \partial_p \ell(\theta; y) \\ \partial_q \ell(\theta; y) \end{pmatrix}, \text{ where} \quad (45)$$

$$\partial_p \ell(\theta; y) = -\frac{2y_1}{r} + \frac{2ry_2}{p^2 + 2pr} - \frac{2qy_3}{q^2 + 2qr} + \frac{y_4}{p} \quad (46)$$

$$\partial_q \ell(\theta; Y) = -\frac{2y_1}{r} - \frac{2py_2}{p^2 + 2pr} + \frac{2ry_3}{q^2 + 2qr} + \frac{y_4}{q} \quad (47)$$

Solving the score equation, $S(\theta) = 0$, thus reduces to solving a system of two polynomials in p and q . Since p and q are proportions, we reject any roots outside the unit simplex.

Differentiating ℓ again and multiplying by -1 gives the observed data information matrix, I . To simplify notation, let $p_y = p^2 + 2pr$ and $q_y = q^2 + 2qr$.

$$I(\theta; y) = - \begin{bmatrix} \partial_p^2 \ell(\theta; y) & \partial_{p,q} \ell(\theta; y) \\ \partial_{p,q} \ell(\theta; y) & \partial_q^2 \ell(\theta; y) \end{bmatrix}, \text{ where} \quad (48)$$

$$\partial_p^2 \ell(\theta; y) = \frac{2y_1}{r^2} + \frac{2y_2(p_y + 2r^2)}{p_y^2} + \frac{4y_3q^2}{q_y^2} + \frac{y_4}{p^2} \quad (49)$$

$$\partial_{p,q} \ell(\theta; y) = \frac{2y_1}{r^2} + \frac{2y_2p^2}{p_y^2} + \frac{2y_3q^2}{q_y^2} \quad (50)$$

$$\partial_q^2 \ell(\theta; y) = \frac{y_1}{r^2} + \frac{4y_2p^2}{p_y^2} + \frac{2y_3(q_y + 2r)}{q_y^2} + \frac{y_4}{q^2} \quad (51)$$

The asymptotic standard error of our MLE is I^{-1} , evaluated at the estimate.

A.2 Complete Data Likelihood, Score and Information

The complete data distribution for our model can be written as follows. Write ρ_i for the probability of genotype i . See Table 2 for the values of these probabilities.

$$\ell_c(\theta; y, x) = \log \binom{n}{x} + \sum x_i \log \rho_i(\theta) \quad (52)$$

$$\equiv \sum y_i \log \rho_i \quad (53)$$

$$\equiv 2x_1 \log r + x_2 \log pr + 2x_3 \log p + x_4 \log qr + 2x_5 \log q + x_6 \log pq \quad (54)$$

$$= (2x_1 + x_2 + x_4) \log r + (x_2 + 2x_3 + x_6) \log p + (x_4 + 2x_5 + x_6) \log q \quad (55)$$

$$= n_O \log r + n_A \log p + n_B \log q \quad (56)$$

where n_O , n_A and n_B are the number of times allele O, A and B arise respectively in the sampled genotypes. Note that ℓ_c depends on y only through x , so we suppress y from our notation for complete data quantities. The complete data score function is

$$S_c(\theta; x) = \begin{pmatrix} \partial_p \ell_c(\theta; x) \\ \partial_q \ell_c(\theta; x) \end{pmatrix}, \text{ where} \quad (57)$$

$$\partial_p \ell_c(\theta; x) = \frac{x_2 + 2x_3 + x_6}{p} - \frac{2x_1 + x_2 + x_4}{r} = \frac{n_A}{p} - \frac{n_O}{r} \quad (58)$$

$$\partial_q \ell_c(\theta; x) = \frac{x_4 + 2x_5 + x_6}{q} - \frac{2x_1 + x_2 + x_4}{r} = \frac{n_B}{q} - \frac{n_O}{r} \quad (59)$$

Notice that the score is linear in x . To make this relationship explicit, we write $S_c(\theta; x) = \mathcal{S}(\theta)x$, where $\mathcal{S}(\theta) \in \mathbb{R}^{2 \times 6}$ consists of the coefficients on x in (58) and (59). We will make use of this linearity later.

Next, we give the information matrix for the complete data.

$$I_c(\theta; x) = - \begin{bmatrix} \partial_p^2 \ell_c(\theta; x) & \partial_{p,q} \ell_c(\theta; x) \\ \partial_{p,q} \ell_c(\theta; x) & \partial_q^2 \ell_c(\theta; x) \end{bmatrix}, \text{ where} \quad (60)$$

$$\partial_p^2 \ell_c(\theta; x) = \frac{x_2 + 2x_3 + x_6}{p^2} + \frac{2x_1 + x_2 + x_4}{r^2} = \frac{n_A}{p^2} + \frac{n_O}{r^2} \quad (61)$$

$$\partial_{p,q} \ell_c(\theta; x) = \frac{2x_1 + x_2 + x_4}{r^2} = \frac{n_O}{r^2} \quad (62)$$

$$\partial_q^2 \ell_c(\theta; x) = \frac{x_4 + 2x_5 + x_6}{q^2} + \frac{2x_1 + x_2 + x_4}{r^2} = \frac{n_B}{q^2} + \frac{n_O}{r^2} \quad (63)$$

A.3 Missing Data Distribution

Many quantities which arise in the EM and MCEM algorithms depend on the missing data distribution (i.e. the conditional distribution of X given $Y = y$). This distribution is best described componentwise in X . First, note that $X_1 = y_1$ and $X_6 = y_4$. Next, we have that

$X_2 + X_3 = y_2$ and $X_4 + X_5 = y_3$. Thus, we can write $X_2|Y = y \sim \text{Bin}(y_2, 2pr/(p^2 + 2pr))$ and $X_4|Y = y \sim \text{Bin}(y_3, 2qr/(q^2 + 2qr))$. Finally, we recover X_3 and X_5 by subtracting X_2 from y_2 and X_4 from y_3 respectively.

We make frequent use of the first few conditional moments of X , so they are listed here for convenience. Let $\alpha_1 = 2pr/(p^2 + 2pr)$ be the probability parameter for the binomial distribution of X_2 given Y , and $\alpha_2 = 1 - \alpha_1$. Similarly, let $\beta_1 = 2qr/(q^2 + 2qr)$ correspond to X_4 and $\beta_2 = 1 - \beta_1$.

$$\mathbb{E}(X|Y = y) = (y_1, y_2\alpha_1, y_2\alpha_2, y_3\beta_1, y_3\beta_2, y_4)^T \quad (64)$$

$$=: \mu_m \quad (65)$$

$$\mathbb{V}(X|Y = y) = \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & y_2\alpha_1\alpha_2 & -y_2\alpha_1\alpha_2 & 0 & 0 & 0 \\ 0 & -y_2\alpha_1\alpha_2 & y_2\alpha_1\alpha_2 & 0 & 0 & 0 \\ 0 & 0 & 0 & y_3\beta_1\beta_2 & -y_3\beta_1\beta_2 & 0 \\ 0 & 0 & 0 & -y_3\beta_1\beta_2 & y_3\beta_1\beta_2 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} \quad (66)$$

$$=: \Sigma_m \quad (67)$$

$$\mathbb{E}(XX^T|Y = y) = \Sigma_m + \mu_m\mu_m^T \quad (68)$$

Conditional expectations of the number of alleles of each kind will be of particular interest.

$$\nu_O := \mathbb{E}(n_O|y) \quad (69)$$

$$= 2y_1 + \frac{y_2 pr}{p^2 + 2pr} + \frac{y_3 qr}{q^2 + 2qr} \quad (70)$$

$$= 2y_1 + y_2 \left(\frac{\rho_2}{\rho_2 + \rho_3} \right) + y_3 \left(\frac{\rho_4}{\rho_4 + \rho_5} \right) \quad \left(= 2y_1 + y_2 \left(\frac{\rho_2}{\pi_2} \right) + y_3 \left(\frac{\rho_4}{\pi_3} \right) \right) \quad (71)$$

$$\nu_A := \mathbb{E}(n_A|y) \quad (72)$$

$$= \frac{2y_2 pr}{p^2 + 2pr} + \frac{2y_2 p^2}{p^2 + 2pr} + y_4 \quad (73)$$

$$= y_2 \left(\frac{\rho_2}{\rho_2 + \rho_3} + 2 \frac{\rho_3}{\rho_2 + \rho_3} \right) + y_4 \quad \left(= y_2 \left(\frac{\rho_2}{\pi_2} + 2 \frac{\rho_3}{\pi_2} \right) + y_4 \right) \quad (74)$$

$$= y_2 \left(1 + \frac{p^2}{p^2 + 2pr} \right) + y_4 \quad (75)$$

$$\nu_B := \mathbb{E}(n_B|y) \quad (76)$$

$$= \frac{2y_3 qr}{q^2 + 2qr} + \frac{2y_3 q^2}{q^2 + 2qr} + y_4 \quad (77)$$

$$= y_3 \left(\frac{\rho_4}{\rho_4 + \rho_5} + 2 \frac{\rho_5}{\rho_4 + \rho_5} \right) + y_4 \quad \left(= y_3 \left(\frac{\rho_4}{\pi_3} + 2 \frac{\rho_5}{\pi_3} \right) + y_4 \right) \quad (78)$$

$$= y_3 \left(1 + \frac{q^2}{q^2 + 2qr} \right) + y_4 \quad (79)$$

A.4 EM Algorithm

In order to apply the EM algorithm, we must construct and optimize the EM objective function. That is, we must compute $Q(\theta|\theta_0) = \mathbb{E}_{\theta_0} [\ell_c(\theta; y, X) | Y = y]$. The conditional distribution of X given $Y = y$ is best described componentwise in X . First, note that X_1 and X_6 occur in Y , so we have $X_1 = y_1$ and $X_6 = y_4$. Next, we have that $X_2 + X_3 = y_2$ and $X_4 + X_5 = y_3$. Thus, we can write $X_2 \sim \text{Bin}(y_2, 2pr/(p^2 + 2pr))$ and $X_4 \sim \text{Bin}(y_3, 2qr/(q^2 + 2qr))$. Finally, we recover X_3 and X_5 by subtracting X_2 from y_2 and X_4 from y_3 respectively.

Before writing out the EM objective function, we first compute conditional expectations

of the sample allele counts.

$$\mathbb{E}(n_O|y) = 2y_1 + \frac{y_2 pr}{p^2 + 2pr} + \frac{y_3 qr}{q^2 + 2qr} \quad (80)$$

$$= 2y_1 + y_2 \left(\frac{\rho_2}{\rho_2 + \rho_3} \right) + y_3 \left(\frac{\rho_4}{\rho_4 + \rho_5} \right) \quad \left(= 2y_1 + y_2 \left(\frac{\rho_2}{\pi_2} \right) + y_3 \left(\frac{\rho_4}{\pi_3} \right) \right) \quad (81)$$

$$=: \nu_O \quad (82)$$

$$\mathbb{E}(n_A|y) = \frac{2y_2 pr}{p^2 + 2pr} + \frac{2y_2 p^2}{p^2 + 2pr} + y_4 \quad (83)$$

$$= y_2 \left(\frac{\rho_2}{\rho_2 + \rho_3} + 2 \frac{\rho_3}{\rho_2 + \rho_3} \right) + y_4 \quad \left(= y_2 \left(\frac{\rho_2}{\pi_2} + 2 \frac{\rho_3}{\pi_2} \right) + y_4 \right) \quad (84)$$

$$= y_2 \left(1 + \frac{p^2}{p^2 + 2pr} \right) + y_4 \quad (85)$$

$$=: \nu_A \quad (86)$$

$$\mathbb{E}(n_B|y) = \frac{2y_3 qr}{q^2 + 2qr} + \frac{2y_3 q^2}{q^2 + 2qr} + y_4 \quad (87)$$

$$= y_3 \left(\frac{\rho_4}{\rho_4 + \rho_5} + 2 \frac{\rho_5}{\rho_4 + \rho_5} \right) + y_4 \quad \left(= y_3 \left(\frac{\rho_4}{\pi_3} + 2 \frac{\rho_5}{\pi_3} \right) + y_4 \right) \quad (88)$$

$$= y_3 \left(1 + \frac{q^2}{q^2 + 2qr} \right) + y_4 \quad (89)$$

$$=: \nu_B \quad (90)$$

The EM objective function can be written as

$$Q(\theta|\theta_0) := \mathbb{E}_{\theta_0}[\ell_c(\theta; X)|Y = y] \quad (91)$$

$$\equiv \nu_O^{(0)} \log r + \nu_A^{(0)} \log p + \nu_B^{(0)} \log q \quad (92)$$

where a superscript zero denotes that the quantity is computed by taking an expectation under θ_0 . Maximizing Q analytically with respect to θ gives the following expression for the EM update:

$$M(\theta_{k-1}) = \begin{pmatrix} \hat{p}_k \\ \hat{q}_k \end{pmatrix} \quad (93)$$

$$= \begin{pmatrix} \nu_A^{(k-1)}/2n \\ \nu_B^{(k-1)}/2n \end{pmatrix} \quad (94)$$

where a superscript $k-1$ denotes that the quantity is computed by taking an expectation under θ_{k-1} . Note that the equation for a critical point of the EM algorithm, $M(\theta) = \theta$, has identical roots to the observed data score equation, $S(\theta) = 0$ **(confirm this)**.

A.4.1 Asymptotic Standard Error

Recall that the EM algorithm computes the MLE, which has asymptotic covariance matrix equal to the inverse Fisher information matrix evaluated at the true parameter value. In practice, we estimate this covariance with the inverse of the observed information matrix evaluated at the MLE. Using Proposition 1.3, we can calculate the observed information matrix using conditional expectations of quantities derived from the complete data likelihood.

To this end, we need to evaluate the conditional expectations in expression (7) of Proposition 1.3. It is convenient for us to write $S_c(\theta) =: \mathcal{S}(\theta)X$ (see Appendix A.2). Then

$$I_c(\hat{\theta}) = \begin{bmatrix} \frac{\nu_A}{p^2} + \frac{\nu_Q}{r^2} & \frac{\nu_Q}{r^2} \\ \frac{\nu_Q}{r^2} & \frac{\nu_B}{q^2} + \frac{\nu_Q}{r^2} \end{bmatrix}, \text{ and} \quad (95)$$

$$\mathbb{E}_{\hat{\theta}}[S_c(\hat{\theta})S_c(\hat{\theta})^T|Y=y] = \mathcal{S}(\hat{\theta})\mathbb{E}_{\hat{\theta}}[XX^T|Y=y]\mathcal{S}(\hat{\theta}) \quad (96)$$

$$= \mathcal{S}(\hat{\theta})(\Sigma_m + \mu_M\mu_M^T)\mathcal{S}(\hat{\theta}) \quad (97)$$

$$(98)$$

While it is possible to expand the above expressions, they quickly become too long to easily interpret. We instead leave these as computational formulas and use them as a guide for writing `R` or `Julia` code.

A.5 Monte Carlo EM

Check for “Citation Needed” before publishing.

References

- S. Agapiou, O. Papaspiliopoulos, D. Sanz-Alonso, and A. M. Stuart. Importance sampling: Intrinsic dimension and computational cost. *Statistical Science*, 32(3), 2017.
- David M. Blei, Alp Kucukelbir, and Jon D. McAuliffe. Variational inference: a review for statisticians. *Journal of the American Statistical Association*, 112(518), 2017.
- James G. Booth and James P. Hobert. Maximizing generalized linear mixed model likelihoods with an automated monte carlo em algorithm. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, 61(1), 1999.
- James G. Booth, James P. Hobert, and Wolfgang Jank. A survey of Monte Carlo algorithms for maximizing the likelihood of a two-stage hierarchical model. *Statistical Modelling*, 1(4), 2001.
- Vivek S. Borkar. *Stochastic Approximation: A dynamical systems viewpoint*. Springer, 2nd edition, 2022.
- Léon Bottou. Large-scale machine learning with stochastic gradient descent. In *Proceedings of COMPSTAT’2010: 19th International Conference on Computational Statistics Paris France, August 22-27, 2010 Keynote, Invited and Contributed Papers*, 2010.
- Mónica F. Bugallo, Víctor Elvira, Luca Martino, David Luengo, Joaquín Míguez, and Petar M. Djurić. Adaptive importance sampling: The past, the present, and the future. *IEEE Signal Processing Magazine*, 34(4), 2017.
- Brian S. Caffo, Wolfgang Jank, and Galin L. Jones. Ascent-based Monte Carlo expectation-maximization. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, 67(2), 2005.
- Gilles Celeux and Jean Diebolt. The EM and SEM algorithms for mixtures: statistical and numerical aspects. Technical report, INRIA, 1987.
- Gilles Celeux and Jean Diebolt. A stochastic approximation type EM algorithm for the mixture problem. *Stochastics and Stochastic Reports*, 41(1-2), 1992.
- Gilles Celeux, Didier Chauveau, and Jean Diebolt. On stochastic versions of the EM algorithm. Technical report, INRIA, 1995.
- K. S. Chan and Johannes Ledolter. Monte Carlo EM estimation for time series models involving counts. *Journal of the American Statistical Association*, 90(429), 1995.

- Nicolas Chopin and Omiros Papaspiliopoulos. *An Introduction to Sequential Monte Carlo*. Springer, 2020.
- Citation Needed.
- Laura Dean. Blood groups and red cell antigens [internet]. Technical report, National Center for Biotechnology Information (US), 2005.
- Pierre Del Moral, Arnaud Doucet, and Ajay Jasra. Sequential Monte Carlo samplers. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 68(3), 2006.
- Bernard Delyon, Marc Lavielle, and Eric Moulines. Convergence of a stochastic approximation version of the EM algorithm. *The Annals of Statistics*, 27(1), 1999.
- A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society. Series B (Methodological)*, 39(1), 1977.
- Bernard Delyon, Marc Lavielle, and Eric Moulines. Convergence of a stochastic approximation version of the EM algorithm. *The Annals of Statistics*, 27(1), 1999.
- Aryeh Dvoretzky. On stochastic approximation. In Jerzy Neyman, editor, *Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Contributions to the Theory of Statistics*, 1956.
- Víctor Elvira, Luca Martino, David Luengo, and Mónica Bugallo. Generalized multiple importance sampling. *Statistical Science*, 34(1), 2019.
- Víctor Elvira and Luca Martino. Advances in importance sampling. *arXiv:2102.05407v3*, 2022.
- Gersende Fort and Eric Moulines. Convergence of the Monte Carlo expectation maximization for curved exponential families. *The Annals of Statistics*, 31(4), 2003.
- Yoshiko Fujita, Masako Tanimura, and Katumi Tanaka. The distribution of the ABO blood groups in Japan. *Japanese Journal of Human Genetics*, 23, 1978.
- Alan E. Gelfand and Bradley P. Carlin. Maximum-likelihood estimation for constrained-or missing-data models. *The Canadian Journal of Statistics*, 21(3), 1993.
- Andrew Gelman, John B. Carlin, Hal S. Stern, David B. Dunson, Aki Vehtari, and Donald B. Rubin. *Bayesian Data Analysis*. CRC Press, third edition, 2013.
- Charles J. Geyer. Markov chain Monte Carlo maximum likelihood. *Interface Foundation of North America*, 1991.

- Charles J. Geyer. On the convergence of Monte Carlo maximum likelihood calculations. *Journal of the Royal Statistical Society. Series B (Methodological)*, 56(1), 1994.
- Ming Gao Gu and Fan Hui Kong. A stochastic approximation algorithm with Markov chain Monte-Carlo method for incomplete data estimation problems. *Proceedings of the National Academy of Sciences*, 95(13), 1998.
- Ming Gao Gu and Shaolin Li. A stochastic approximation algorithm for maximum-likelihood estimation with incomplete data. *The Canadian Journal of Statistics*, 26(4), 1998.
- Ming Gao Gu and Hong-Tu Zhu. Maximum likelihood estimation for spatial models by Markov Chain Monte Carlo stochastic approximation. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, 63(2), 2001.
- Edward L. Ionides. Truncated importance sampling. *Journal of Computational and Graphical Statistics*, 17(2), 2008.
- Wolfgang Jank. Implementing and diagnosing the stochastic approximation EM algorithm. *Journal of Computational and Graphical Statistics*, 15(4), 2006.
- Wolfgang Jank and James Booth. Efficiency of Monte Carlo EM and simulated likelihood in two-stage hierarchical models. *Journal of Computational and Graphical Statistics*, 12(1), 2003.
- J. Kiefer and J. Wolfowitz. Stochastic estimation of the maximum of a regression function. *The Annals of Mathematical Statistics*, 23(3), 1952.
- Harold J. Kushner and G. George Yin. *Stochastic Approximation Algorithms and Applications*. Springer, 1997.
- Tze Leung Lai. Stochastic approximation: invited paper. *The Annals of Statistics*, 31(2), 2003.
- Richard A. Levine and George Casella. Implementations of the Monte Carlo EM algorithm. *Journal of Computational and Graphical Statistics*, 10(3), 2001.
- Richard A. Levine and Juanjuan Fan. An automated (Markov chain) Monte Carlo EM algorithm. *Journal of Statistical Computation and Simulation*, 74(5), 2004.
- Thomas A. Louis. Finding the observed information matrix when using the EM algorithm. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, 44(2), 1982.
- Charles E. McCulloch. Maximum likelihood algorithms for generalized linear mixed models. *Journal of the American Statistical Association*, 92(437), 1997.

- Geoffrey J. McLachlan and Thriyambakam Krishnan. *The EM algorithm and extensions*. Wiley, 2nd edition, 2008.
- Xiao-Li Meng and Donald B. Rubin. On the global and componentwise rates of convergence of the EM algorithm. *Linear Algebra and its Applications*, 199, 1994.
- Giuseppe Moffa and Jack Kuipers. Sequential Monte Carlo EM for multivariate probit models. *Computational Statistics and Data Analysis*, 72, 2014.
- Radford M. Neal and Geoffrey E. Hinton. *A view of the EM algorithm that justifies incremental, sparse, and other variants*, pages 355–368. Springer, 1998.
- Ronald C. Neath. On convergence properties of the Monte Carlo EM algorithm. *Advances in Modern Statistical Theory and Applications*, 10, 2013.
- Søren Feodor Nielsen. The stochastic EM algorithm: estimation and asymptotic results. *Bernoulli*, 6(3), 2000.
- Jorge Nocedal and Stephen J. Wright. *Numerical Optimization*. Springer, 2nd edition, 2006.
- David Oakes. Direct calculation of the information matrix via the EM algorithm. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, 61(2), 1999.
- B. T. Polyak and A. B. Juditsky. Acceleration of stochastic approximation by averaging. *SIAM Journal on Control and Optimization*, 30(4), 1992.
- P. Rain, F. Loubeau, A. Durieux, F. Le Strat, and F. Fresnet. Using motorettes for the experimental and numerical determinations of the pdiv in an electric motor. In *2016 IEEE International Conference on Dielectrics (ICD)*, volume 2, pages 967–970, 2016.
- Herbert Robbins and Sutton Monroe. A stochastic approximation method. *The Annals of Mathematical Statistics*, 22(3), 1951.
- Christian P. Robert and George Casella. *Monte Carlo statistical methods*. Springer, second edition, 2004.
- William Ruth and Richard Lockhart. *Lthc analysis*, 2023.
- Stan Development Team. Stan modelling language users guide and reference manual, 2022. URL <http://mc-stan.org/>. Version 2.32.
- Stan Development Team. RStan: the R interface to Stan, 2023. URL <https://mc-stan.org/>. R package version 2.21.8.
- Martin A. Tanner. *Tools for Statistical Inference*. Springer, 2nd edition, 1993.

- Martin A. Tanner. *Tools for Statistical Inference*. Springer, 3rd edition, 1996.
- Dimitris G. Tsikas, Aristidis C. Likas, and Nikolaos P. Galatsanos. The variational approximation for Bayesian inference. *IEEE Signal Processing Magazine*, 25(6), 2008.
- A. W. van der Vaart. *Asymptotic Statistics*. Cambridge University Press, 1998.
- Aki Vehtari, Daniel Simpson, Andrew Gelman, Yuling Yao, and Jonah Gabry. Pareto smoothed importance sampling, 2022.
- Greg C. G. Wei and Martin A. Tanner. A Monte Carlo implementation of the EM algorithm and the poor man’s data augmentation algorithms. *Journal of the American Statistical Association*, 85(411), 1990.
- C. F. Jeff Wu. On the convergence properties of the EM algorithm. *The Annals of Statistics*, 11(1), 1983.

Index

$\hat{\mathbb{E}}$, 8

Blood Type, 5

Complete Data Distribution, 1

Gene Frequency Estimation, 5

Generalized EM Algorithm, 3

Gibbs Sampling, 27

Metropolis-Hastings Algorithm, 27

Missing Data Distribution, 1

Observed Data Distribution, 1