

An Approximate Fisher Scoring Algorithm for Finite Mixtures of Multinomials

Andrew M. Raim, Minglei Liu, Nagaraj K. Neerchal and Jorge G. Morel

Abstract

Finite mixture distributions arise naturally in many applications including clustering and classification. Since they usually do not yield closed forms for maximum likelihood estimates (MLEs), numerical methods using the well known Fisher Scoring or Expectation-Maximization algorithms are considered. In this work, an approximation to the Fisher Information Matrix of an arbitrary mixture of multinomial distributions is introduced. This leads to an Approximate Fisher Scoring algorithm (AFSA), which turns out to be closely related to Expectation-Maximization, and is more robust to the choice of initial value than Fisher Scoring iterations. A combination of AFSA and the classical Fisher Scoring iterations provides the best of both computational efficiency and stable convergence properties.

Key Words: Multinomial; Finite mixture, Maximum likelihood, Fisher Information Matrix, Fisher Scoring.

1 Introduction

A finite mixture model arises when observations are believed to originate from one of several populations, but it is unknown to which population each observation belongs. Model identifiability of mixtures is an important issue with a considerable literature; see, for example, Robbins (1948) and Teicher (1960). The book by McLachlan and Peel (2000) is a good entry point to the literature on mixtures. A majority of the mixture literature deals with mixtures of normal distributions; however, Blischke (1962; 1964), Krolikowska (1976), and Kabir (1968) are a few early works which address mixtures of discrete distributions. The present work focuses on mixtures of multinomials, which have wide applications in clustering and classification problems, as well as modeling overdispersion data (Morel and Nagaraj, 2012).

It is well known that computation of maximum likelihood estimates (MLE) under mixture distributions is often analytically intractable, and therefore iterative numerical methods are needed. Classical iterative techniques such as Newton-Raphson and Fisher Scoring are two widely used methods. The more recent Expectation-Maximization (EM) algorithm discussed in Dempster, Laird and Rubin (1977) has become another standard technique to compute MLEs. EM is a framework for performing estimation in missing data problems. The idea

Andrew Raim is graduate student at Department of Mathematics and Statistics, University of Maryland, Baltimore County, Baltimore, MD 21250, USA (Email: araim1@umbc.edu). Minglei Liu is Senior Principal Biostatistician at Medtronic, Santa Rosa, CA, USA. Nagaraj Neerchal is Professor at Department of Mathematics and Statistics, University of Maryland, Baltimore County, Baltimore, MD 21250, USA (Email: nagaraj@umbc.edu). Jorge Morel is Principal Statistician at Procter & Gamble, Cincinnati, OH, USA.

is to solve a difficult incomplete data problem by repeatedly solving tractable complete-data problems. If the unknown population labels are treated as missing information, then estimation under the mixture distribution can be considered a missing data problem, and an EM algorithm can be used. Unlike Fisher Scoring (FSA), EM does not require computation of an expected Hessian in each iteration, which is a great advantage if this matrix is difficult to compute. Slow speed of convergence has been cited as a disadvantage of EM. Variations and improved versions of the EM algorithm have been widely used for obtaining MLEs for mixtures (Mclachlan and Peel, 2000, chapter 2).

Fisher Scoring iterations require the inverse of the Fisher Information Matrix (FIM). In the mixture setting, computing the FIM involves a complicated expectation which does not have an analytically tractable form. The matrix can be approximated numerically by Monte Carlo simulation for example, but this is computationally expensive, especially when repeated over many iterations. Morel and Nagaraj (1991; 1993) proposed a variant of Fisher Scoring using an approximate FIM in their study of a multinomial model with extra variation. This model, now referred to as the Random Clumped Multinomial (see Example 3.1 for details), is a special case of the finite mixture of multinomials. The approximate FIM was justified asymptotically, and was used to obtain MLEs for the model and to demonstrate their efficiency. In the present paper, we extend the approximate FIM idea to general finite mixtures of multinomials and hence formulate the Approximate Fisher Scoring Algorithm (AFSA) for this family of distributions. By using the approximate FIM in place of the true FIM, we obtain an algorithm which is closely related to EM. Both AFSA and EM have a slower convergence rate than Fisher Scoring once they are in the proximity of a maximum, but both are also much more robust than Fisher Scoring in finding such regions from an arbitrary initial value.

The rest of the paper is organized as follows. In section 2, a large cluster approximation for the Fisher Information Matrix is derived and some of its properties are presented. This approximate information matrix is easily computed and has an immediate application in Fisher Scoring, which is presented in section 3. Simulation studies are presented in section 4, illustrating convergence properties of the approximate information matrix and approximate Fisher Scoring. Concluding remarks are given in section 5.

2 An Approximate Fisher Information Matrix

Consider the multinomial sample space with m trials placed into k categories at random,

$$\Omega = \left\{ \mathbf{x} = (x_1, \dots, x_k) : x_j \in \{0, 1, \dots, m\}, \sum_{j=1}^k x_j = m \right\}.$$

The standard multinomial density is

$$f(\mathbf{x}; \mathbf{p}, m) = \frac{m!}{x_1! \dots x_k!} p_1^{x_1} \dots p_k^{x_k} \cdot I(\mathbf{x} \in \Omega),$$

where $I(\cdot)$ is the indicator function, and the parameter space is

$$\left\{ \mathbf{p} = (p_1, \dots, p_{k-1}) : 0 < p_j < 1, \sum_{j=1}^{k-1} p_j < 1 \right\} \subseteq \mathbb{R}^{k-1}.$$

If a random variable \mathbf{X} has distribution $f(\mathbf{x}; \mathbf{p}, m)$, we will write $\mathbf{X} \sim \text{Mult}_k(\mathbf{p}, m)$. Since $x_k = m - \sum_{j=1}^{k-1} x_j$ and $p_k = 1 - \sum_{j=1}^{k-1} p_j$, the k th category can be considered as redundant information. Following the sampling and overdispersion literature, we will refer to the number of trials m as the “cluster size” of a multinomial observation.

Now suppose there are s multinomial populations

$$\text{Mult}_k(\mathbf{p}_1, m), \dots, \text{Mult}_k(\mathbf{p}_s, m), \quad \mathbf{p}_\ell = (p_{\ell 1}, \dots, p_{\ell, k-1})$$

where the ℓ th population occurs with proportion π_ℓ for $\ell = 1, \dots, s$. If we draw \mathbf{X} from the mixed population, its probability density is a finite mixture of multinomials

$$f(\mathbf{x}; \boldsymbol{\theta}) = \sum_{\ell=1}^s \pi_\ell f(\mathbf{x}; \mathbf{p}_\ell, m), \quad \boldsymbol{\theta} = (\mathbf{p}_1, \dots, \mathbf{p}_s, \boldsymbol{\pi}), \quad (2.1)$$

and we will write $\mathbf{X} \sim \text{MultMix}_k(\boldsymbol{\theta}, m)$. The dimension of $\boldsymbol{\theta}$ is

$$q := s(k-1) + (s-1) = sk - 1,$$

disregarding the redundant parameters $p_{1k}, \dots, p_{sk}, \pi_s$. We will also make use of the following slightly-less-cumbersome notation for densities,

$$\begin{aligned} P(\mathbf{x}) &:= f(\mathbf{x}; \boldsymbol{\theta}, m) : && \text{the mixture,} \\ P_\ell(\mathbf{x}) &:= f(\mathbf{x}; \mathbf{p}_\ell, m) : && \text{the } \ell\text{th component of the mixture.} \end{aligned}$$

The setting of this paper will be an independent sample

$$\mathbf{X}_i \sim \text{MultMix}_k(\boldsymbol{\theta}, m_i), \quad i = 1, \dots, n$$

with cluster sizes not necessarily equal. The resulting likelihood is

$$L(\boldsymbol{\theta}) = \prod_{i=1}^n f(\mathbf{x}_i; \boldsymbol{\theta}) = \prod_{i=1}^n \left\{ \sum_{\ell=1}^s \pi_\ell \left[\frac{m_i!}{x_{i1}! \dots x_{ik}!} p_{\ell 1}^{x_{i1}} \dots p_{\ell k}^{x_{ik}} \cdot I(\mathbf{x}_i \in \Omega) \right] \right\}. \quad (2.2)$$

The inner summation prevents closed-form likelihood maximization, hence our goal will be to compute the MLE $\hat{\boldsymbol{\theta}}$ numerically. Some additional preliminaries are given in Appendix A.

In general, as mentioned earlier, the Fisher Information Matrix (FIM) for mixtures involves a complicated expectation which does not have a tractable form. Since the multinomial

mixture has a finite sample space, it can be computed naively by using the definition of the expectation

$$\mathcal{I}(\boldsymbol{\theta}) = \sum_{\mathbf{x} \in \Omega} \left\{ \frac{\partial}{\partial \boldsymbol{\theta}} \log f(\mathbf{x}; \boldsymbol{\theta}) \right\} \left\{ \frac{\partial}{\partial \boldsymbol{\theta}} \log f(\mathbf{x}; \boldsymbol{\theta}) \right\}^T f(\mathbf{x}; \boldsymbol{\theta}), \quad (2.3)$$

given a particular value for $\boldsymbol{\theta}$. Although the number of terms $\binom{k+m-1}{m}$ in the summation is finite, it grows quickly with m and k , and this method becomes intractable as m and k increase. For example, when $m = 100$ and $k = 4$, the sample space Ω contains more than 178,000 elements. To avoid these potentially expensive computations, we extend the approximate FIM approach of Morel and Nagaraj (1991; 1993) to the general finite mixture of multinomials. The following theorem states our main result.

Theorem 2.1. *Suppose $\mathbf{X} \sim \text{MultMix}_k(\boldsymbol{\theta}, m)$ is a single observation from the mixed population. Denote the exact FIM with respect to \mathbf{X} as $\mathcal{I}(\boldsymbol{\theta})$. Then an approximation to the FIM with respect to \mathbf{X} is given by the $(sk - 1) \times (sk - 1)$ block-diagonal matrix*

$$\tilde{\mathcal{I}}(\boldsymbol{\theta}) := \text{Blockdiag}(\pi_1 \mathbf{F}_1, \dots, \pi_s \mathbf{F}_s, \mathbf{F}_\pi),$$

where for $\ell = 1, \dots, s$

$$\mathbf{F}_\ell = m [\mathbf{D}_\ell^{-1} + p_{\ell k}^{-1} \mathbf{1}\mathbf{1}^T] \quad \text{and} \quad \mathbf{D}_\ell = \text{diag}(p_{\ell 1}, \dots, p_{\ell, k-1})$$

are $(k - 1) \times (k - 1)$ matrices,

$$\mathbf{F}_\pi = \mathbf{D}_\pi^{-1} + \pi_s^{-1} \mathbf{1}\mathbf{1}^T \quad \text{and} \quad \mathbf{D}_\pi = \text{diag}(\pi_1, \dots, \pi_{s-1})$$

is a $(s - 1) \times (s - 1)$ matrix, and $\mathbf{1}$ denotes a vector of ones of the appropriate dimension. To emphasize the dependence of the FIM and the approximation on m , we will also write $\mathcal{I}_m(\boldsymbol{\theta})$ and $\tilde{\mathcal{I}}_m(\boldsymbol{\theta})$. If the vectors $\mathbf{p}_1, \dots, \mathbf{p}_s$ are distinct (i.e. $\mathbf{p}_a \neq \mathbf{p}_b$ for every pair of populations $a \neq b$), then $\mathcal{I}_m(\boldsymbol{\theta}) - \tilde{\mathcal{I}}_m(\boldsymbol{\theta})$ as $m \rightarrow \infty$.

A proof is given in Appendix B. Notice that the matrix \mathbf{F}_ℓ is exactly the FIM of $\text{Mult}_k(\mathbf{p}_\ell, m)$ for the ℓ th population, and \mathbf{F}_π is the FIM of $\text{Mult}_s(\boldsymbol{\pi}, 1)$ corresponding to the mixing probabilities $\boldsymbol{\pi}$; see Appendix A for details. The approximate FIM turns out to be equivalent to a complete data FIM, as shown in Proposition 2.2 below, which provides an interesting connection to EM. This matrix can be formulated for any finite mixture whose components have a well-defined FIM, and is not limited to the case of multinomials.

Proposition 2.2. *The matrix $\tilde{\mathcal{I}}(\boldsymbol{\theta})$ is equivalent to the FIM of (\mathbf{X}, Z) , where*

$$Z = \begin{cases} 1 & \text{with probability } \pi_1 \\ \vdots & \\ s & \text{with probability } \pi_s, \end{cases} \quad \text{and} \quad (\mathbf{X} \mid Z = \ell) \sim \text{Mult}_k(\mathbf{p}_\ell, m). \quad (2.4)$$

Proof of Proposition 2.2. Here Z represents the population from which \mathbf{X} was drawn. The complete data likelihood is then

$$L(\boldsymbol{\theta} \mid \mathbf{x}, z) = \prod_{\ell=1}^s \left[\pi_{\ell} f(\mathbf{x} \mid \mathbf{p}_{\ell}, m) \right]^{I(z=\ell)}.$$

This likelihood leads to the score vectors

$$\begin{aligned} \frac{\partial}{\partial \mathbf{p}_a} \log L(\boldsymbol{\theta}) &= \Delta_a \left[\mathbf{D}_a^{-1} \mathbf{x}_{-k} - \frac{x_k}{p_{ak}} \mathbf{1} \right], \\ \frac{\partial}{\partial \boldsymbol{\pi}} \log L(\boldsymbol{\theta}) &= \mathbf{D}_{\pi}^{-1} \boldsymbol{\Delta}_{-s} - \frac{\Delta_s}{\pi_s} \mathbf{1}, \end{aligned}$$

where $\boldsymbol{\Delta} = (\Delta_1, \dots, \Delta_s)$ so that $\Delta_{\ell} = I(Z = \ell) \sim \text{Bernoulli}(\pi_{\ell})$, and $\boldsymbol{\Delta}_{-s}$ denotes the vector $(\Delta_1, \dots, \Delta_{s-1})$. Taking second derivatives yields

$$\begin{aligned} \frac{\partial^2}{\partial \mathbf{p}_a \partial \mathbf{p}_a^T} \log L(\boldsymbol{\theta}) &= -\Delta_a \left[\mathbf{D}_a^{-2} \mathbf{x}_{-k} + \frac{x_k}{p_{ak}^2} \mathbf{1} \mathbf{1}^T \right], \\ \frac{\partial^2}{\partial \mathbf{p}_a \partial \mathbf{p}_b^T} \log L(\boldsymbol{\theta}) &= 0, \quad \text{for } a \neq b, \\ \frac{\partial^2}{\partial \mathbf{p}_a \partial \boldsymbol{\pi}^T} \log L(\boldsymbol{\theta}) &= 0, \\ \frac{\partial^2}{\partial \boldsymbol{\pi} \partial \boldsymbol{\pi}^T} \log L(\boldsymbol{\theta}) &= - \left[\mathbf{D}_{\pi}^{-2} \boldsymbol{\Delta}_{-s} + \frac{\Delta_s}{\pi_s^2} \mathbf{1} \mathbf{1}^T \right]. \end{aligned}$$

Now take the expected value of the negative of each of these terms, jointly with respect to (\mathbf{X}, Z) , to obtain the blocks of $\tilde{\mathcal{I}}(\boldsymbol{\theta})$. \square

Corollary 2.3. *Suppose $\mathbf{X}_i \sim \text{MultMix}(\boldsymbol{\theta}, m_i)$, $i = 1, \dots, n$, is an independent sample from the mixed population with varying cluster sizes, and $M = m_1 + \dots + m_n$. Then the approximate FIM with respect to $(\mathbf{X}_1, \dots, \mathbf{X}_n)$ is given by*

$$\begin{aligned} \tilde{\mathcal{I}}(\boldsymbol{\theta}) &= \text{Blockdiag}(\pi_1 \mathbf{F}_1, \dots, \pi_s \mathbf{F}_s, \mathbf{F}_{\pi}), \\ \mathbf{F}_{\ell} &= M \left[\mathbf{D}_{\ell}^{-1} + p_{\ell k}^{-1} \mathbf{1} \mathbf{1}^T \right], \quad \ell = 1, \dots, s \\ \mathbf{F}_{\pi} &= n \left[\mathbf{D}_{\pi}^{-1} + \pi_s^{-1} \mathbf{1} \mathbf{1}^T \right]. \end{aligned}$$

Proof of Corollary 2.3. Let $\tilde{\mathcal{I}}_i(\boldsymbol{\theta})$ represent the approximate FIM with respect to observation \mathbf{X}_i . The result is obtained using

$$\tilde{\mathcal{I}}(\boldsymbol{\theta}) = \tilde{\mathcal{I}}_1(\boldsymbol{\theta}) + \dots + \tilde{\mathcal{I}}_n(\boldsymbol{\theta}),$$

corresponding to the additive property of exact FIMs for independent samples. The additive property can be justified by noting that each $\tilde{\mathcal{I}}_i(\boldsymbol{\theta})$ is a true (complete data) FIM, by Proposition 2.2. \square

Since $\tilde{\mathcal{I}}(\boldsymbol{\theta})$ is a block diagonal matrix, some useful expressions can be obtained in closed form.

Corollary 2.4. *Let $\tilde{\mathcal{I}}(\boldsymbol{\theta})$ represent the FIM with respect to an independent sample $\mathbf{X}_i \sim \text{MultMix}(\boldsymbol{\theta}, m_i)$, $i = 1, \dots, n$. Then:*

(a) *The inverse of $\tilde{\mathcal{I}}(\boldsymbol{\theta})$ is given by*

$$\begin{aligned} \tilde{\mathcal{I}}^{-1}(\boldsymbol{\theta}) &= \text{Blockdiag}(\pi_1^{-1} \mathbf{F}_1^{-1}, \dots, \pi_s^{-1} \mathbf{F}_s^{-1}, \mathbf{F}_\pi^{-1}), \\ \mathbf{F}_\ell^{-1} &= M^{-1} \{ \mathbf{D}_\ell - \mathbf{p}_\ell \mathbf{p}_\ell^T \}, \quad \ell = 1, \dots, s \\ \mathbf{F}_\pi^{-1} &= n^{-1} \{ \mathbf{D}_\pi - \boldsymbol{\pi} \boldsymbol{\pi}^T \}. \end{aligned} \tag{2.5}$$

(b) *The trace of $\tilde{\mathcal{I}}(\boldsymbol{\theta})$ is given by*

$$\text{tr}(\tilde{\mathcal{I}}(\boldsymbol{\theta})) = \sum_{\ell=1}^s \sum_{j=1}^{k-1} M \pi_\ell \{ p_{\ell j}^{-1} + p_{\ell k}^{-1} \} + \sum_{\ell=1}^{s-1} n \{ \pi_\ell^{-1} + \pi_s^{-1} \}.$$

(c) *The determinant of $\tilde{\mathcal{I}}(\boldsymbol{\theta})$ is given by*

$$\det(\tilde{\mathcal{I}}(\boldsymbol{\theta})) = \left(\prod_{\ell=1}^s p_{\ell k}^{-1} \prod_{j=1}^{k-1} M \pi_\ell p_{\ell j}^{-1} \right) \left(\pi_s^{-1} \prod_{\ell=1}^{s-1} n \pi_\ell^{-1} \right).$$

Proof of Corollary 2.4 (a). Since $\tilde{\mathcal{I}}(\boldsymbol{\theta})$ is block diagonal, its inverse can be obtained by inverting the blocks, which can immediately be seen to be (2.5). To find the expressions for the individual blocks, we can apply the Sherman-Morrison formula (see for example Rao 1965, chapter 1)

$$(\mathbf{C} + \mathbf{u} \mathbf{v}^T)^{-1} = \mathbf{C}^{-1} - \frac{\mathbf{C}^{-1} \mathbf{u} \mathbf{v}^T \mathbf{C}^{-1}}{1 + \mathbf{v}^T \mathbf{C}^{-1} \mathbf{u}}.$$

For the case of \mathbf{F}_π^{-1} , for example, take $\mathbf{C} = \mathbf{D}_\pi^{-1}$, $\mathbf{u} = \pi_s^{-1/2} \mathbf{1}$, and $\mathbf{v} = \pi_s^{-1/2} \mathbf{1}^T$ and use the expressions in Corollary 2.3. \square

Proof of Corollary 2.4 (b). Since the trace of a block diagonal matrix is the sum of the traces of its blocks, we have

$$\text{tr}(\tilde{\mathcal{I}}(\boldsymbol{\theta})) = \pi_1 \text{tr}(\mathbf{F}_1) + \dots + \pi_s \text{tr}(\mathbf{F}_s) + \text{tr}(\mathbf{F}_\pi). \tag{2.6}$$

The individual traces can be obtained as

$$\text{tr}(\mathbf{F}_\ell) = \text{tr} [M(\mathbf{D}_\ell^{-1} + p_{\ell k}^{-1} \mathbf{1} \mathbf{1}^T)] = \sum_{j=1}^{k-1} M \{ p_{\ell j}^{-1} + p_{\ell k}^{-1} \},$$

a summation over the diagonal elements. Similarly for the block corresponding to $\boldsymbol{\pi}$,

$$\text{tr}(\mathbf{F}_\pi) = \text{tr} \left[n \left(\mathbf{D}_\pi^{-1} + \pi_s^{-1} \mathbf{1}\mathbf{1}^T \right) \right] = \sum_{\ell=1}^{s-1} n \{ \pi_\ell^{-1} + \pi_s^{-1} \}.$$

The result is obtained by replacing these expressions into (2.6). \square

Proof of Corollary 2.4 (c). Since $\tilde{\mathcal{I}}(\boldsymbol{\theta})$ has a block diagonal structure,

$$\begin{aligned} \det \tilde{\mathcal{I}}(\boldsymbol{\theta}) &= \det \{ \mathbf{F}_\pi \} \times \prod_{\ell=1}^s \det \{ \pi_\ell \mathbf{F}_\ell \} \\ &= \left(n^{s-1} \det \{ \mathbf{D}_\pi^{-1} + \pi_s^{-1} \mathbf{1}\mathbf{1}^T \} \right) \left(\prod_{\ell=1}^s \pi_\ell^{k-1} M^{k-1} \det \{ \mathbf{D}_\ell^{-1} + \mathbf{p}_{\ell k}^{-1} \mathbf{1}\mathbf{1}^T \} \right) \end{aligned} \quad (2.7)$$

Recall the property (see for example Rao 1965, chapter 1) that for \mathbf{M} non-singular, we have

$$\det(\mathbf{M} + \mathbf{u}\mathbf{u}^T) = \begin{vmatrix} \mathbf{M} & -\mathbf{u} \\ \mathbf{u}^T & 1 \end{vmatrix} = \det(\mathbf{M}) (1 + \mathbf{u}^T \mathbf{M}^{-1} \mathbf{u}).$$

This yields, for instance

$$\begin{aligned} \det \{ \mathbf{D}_\pi^{-1} + \pi_s^{-1} \mathbf{1}\mathbf{1}^T \} &= \det \{ \mathbf{D}_\pi^{-1} \} (1 + \pi_s^{-1} \mathbf{1}^T \mathbf{D}_\pi \mathbf{1}) \\ &= \left[1 + \frac{1 - \pi_s}{\pi_s} \right] \prod_{\ell=1}^{s-1} \pi_\ell^{-1} = \pi_s^{-1} \prod_{\ell=1}^{s-1} \pi_\ell^{-1}. \end{aligned}$$

The result can be obtained by substituting the simplified determinants into (2.7). \square

The determinant and trace of the FIM are not utilized in the computation of MLEs, but are used in the computation of many statistics in subsequent analysis. In such applications, it may be preferable to have a closed form for these expressions. As one example, consider the Consistent Akaike Information Criterion with Fisher Information (CAICF) formulated in (Bozdogan, 1987). The CAICF is an information-theoretic criterion for model selection, and is a function of the log-determinant of the FIM.

It can also be shown that $\mathcal{I}_m^{-1}(\boldsymbol{\theta}) - \tilde{\mathcal{I}}_m^{-1}(\boldsymbol{\theta}) \rightarrow \mathbf{0}$ as $m \rightarrow \infty$, which we now state as a theorem. A proof is given in Appendix B. This result is perhaps more immediately relevant than Theorem 2.1 for our Fisher Scoring application presented in the following section.

Theorem 2.5. *Let $\mathcal{I}_m(\boldsymbol{\theta})$ and $\tilde{\mathcal{I}}_m(\boldsymbol{\theta})$ be defined as in Theorem 2.1 (namely the FIM and approximate FIM with respect to a single observation with cluster size m). Then $\mathcal{I}_m^{-1}(\boldsymbol{\theta}) - \tilde{\mathcal{I}}_m^{-1}(\boldsymbol{\theta}) \rightarrow \mathbf{0}$ as $m \rightarrow \infty$.*

In the next section, we use the approximate FIM obtained in Theorem 2.1 to define an approximate Fisher Scoring algorithm and investigate its properties.

3 Approximate Fisher Scoring Algorithm

Consider an independent sample with varying cluster sizes

$$\mathbf{X}_i \sim \text{MultMix}_k(\boldsymbol{\theta}, m_i), \quad i = 1, \dots, n.$$

Let $\boldsymbol{\theta}^{(0)}$ be an initial guess for $\boldsymbol{\theta}$, and $S(\boldsymbol{\theta})$ be the score vector with respect to the sample (see Appendix A). Then by independence

$$S(\boldsymbol{\theta}) = \sum_{i=1}^n S(\boldsymbol{\theta}; \mathbf{x}_i),$$

where $S(\boldsymbol{\theta}; \mathbf{x}_i)$ is the score vector with respect to the i th observation. The Fisher Scoring Algorithm is given by computing the iterations

$$\boldsymbol{\theta}^{(g+1)} = \boldsymbol{\theta}^{(g)} + \mathcal{I}^{-1}(\boldsymbol{\theta}^{(g)}) S(\boldsymbol{\theta}^{(g)}), \quad g = 1, 2, \dots \quad (3.1)$$

until the convergence criteria

$$|\log L(\boldsymbol{\theta}^{(g+1)}) - \log L(\boldsymbol{\theta}^{(g)})| < \varepsilon$$

is met, for some given tolerance $\varepsilon > 0$. In practice, a line search may be used for every iteration after determining a search direction, but such modifications will not be considered here. Note that (3.1) uses the exact FIM which may not be easily computable. We propose to substitute the approximation $\tilde{\mathcal{I}}(\boldsymbol{\theta})$ for $\mathcal{I}(\boldsymbol{\theta})$, and will refer to the resulting method as the Approximate Fisher Scoring Algorithm (AFSA). The expressions for $\tilde{\mathcal{I}}(\boldsymbol{\theta})$ and its inverse are available in closed form, as seen in Corollaries 2.3 and 2.4.

AFSA can be applied to finite mixture of multinomial models which are not explicitly in the form of (2.2). We now give two examples which use AFSA to compute MLEs for such models. The first is the Random Clumped model for overdispersed multinomial data. The second is an arbitrary mixture of multinomials with links from parameters to covariates.

Example 3.1. In section 1 we have mentioned the Random Clumped Multinomial (RCM), a distribution that addresses overdispersion due to “clumped” sampling in the multinomial framework. RCM represents an interesting model for exploring computational methods. Recently, Zhou and Lange (2010) have used it as an illustrative example for the minorization-maximization principle. Raim et al (2012) have explored parallel computing in maximum likelihood estimation using large RCM models as a test problem. It turns out that RCM conforms to the finite mixture of multinomials representation (2.1), and can therefore be fitted by the AFSA algorithm. Once the mixture representation is established, the score vector and approximate FIM can be formulated by the use of transformations; see for example section 2.6 of Lehmann and Casella (1998). Hence, we can obtain the algorithm presented in Morel and Nagaraj (1993) and Neerchal and Morel (1998) as an AFSA-type algorithm.

Consider a cluster of m trials, where each trial results in one of k possible outcomes with probabilities (π_1, \dots, π_k) . Suppose a default category is also selected at random, so that each

trial either results in this default outcome with probability ρ , or an independent choice with probability $1 - \rho$. Intuitively, if $\rho \rightarrow 0$, RCM approaches a standard multinomial distribution. Using this idea, an RCM random variable can be obtained from the following procedure. Let $\mathbf{Y}_0, \mathbf{Y}_1, \dots, \mathbf{Y}_m \stackrel{\text{iid}}{\sim} \text{Mult}_k(\boldsymbol{\pi}, 1)$ and $\mathbf{U}_1, \dots, \mathbf{U}_m \stackrel{\text{iid}}{\sim} U(0, 1)$ be independent samples, then

$$\begin{aligned} \mathbf{X} &= \mathbf{Y}_0 \sum_{i=1}^m I(\mathbf{U}_i \leq \rho) + \sum_{i=1}^m \mathbf{Y}_i I(\mathbf{U}_i > \rho) \\ &= \mathbf{Y}_0 N + (\mathbf{Z} \mid N) \end{aligned} \tag{3.2}$$

follows the distribution $\text{RCM}_k(\boldsymbol{\pi}, \rho)$. The representation (3.2) emphasizes that $N \sim \text{Binomial}(m, \rho)$, $(\mathbf{Z} \mid N) \sim \text{Mult}_k(\boldsymbol{\pi}, m - N)$, and $\mathbf{Y}_0 \sim \text{Mult}_k(\boldsymbol{\pi}, 1)$, where N and \mathbf{Y}_0 are independent.

RCM is also a special case of the finite mixture of multinomials, so that

$$\begin{aligned} \mathbf{X} \sim f(\mathbf{x}; \boldsymbol{\pi}, \rho) &= \sum_{\ell=1}^k \pi_\ell f(\mathbf{x}; \mathbf{p}_\ell, m), \\ \mathbf{p}_\ell &= (1 - \rho)\boldsymbol{\pi} + \rho \mathbf{e}_\ell, \quad \text{for } \ell = 1, \dots, k - 1 \\ \mathbf{p}_k &= (1 - \rho)\boldsymbol{\pi}, \end{aligned}$$

where $f(\mathbf{x}; \mathbf{p}, m)$ is our usual notation for the density of $\text{Mult}_k(\mathbf{p}, m)$. This mixture representation can be derived using moment generating functions, as shown in (Morel and Nagaraj, 1993). Notice that in this mixture $s = k$ so that the number of mixture components matches the number of categories. There are also only k distinct parameters rather than $sk - 1$ as in the general mixture.

The approximate FIM for the RCM model can be obtained by transformation, starting with the expression for the general mixture. Consider transforming the k dimensional $\boldsymbol{\eta} = (\boldsymbol{\pi}, \rho)$ to the $q = sk - 1 = (k + 1)(k - 1)$ dimensional $\boldsymbol{\theta} = (\mathbf{p}_1, \dots, \mathbf{p}_s, \boldsymbol{\pi})$ so that

$$\boldsymbol{\theta}(\boldsymbol{\eta}) = \begin{pmatrix} (1 - \rho)\boldsymbol{\pi} & + & \rho \mathbf{e}_1 \\ & & \vdots \\ (1 - \rho)\boldsymbol{\pi} & + & \rho \mathbf{e}_{k-1} \\ (1 - \rho)\boldsymbol{\pi} \\ \boldsymbol{\pi} \end{pmatrix}.$$

The $q \times k$ Jacobian of this transformation is

$$\frac{\partial \boldsymbol{\theta}}{\partial \boldsymbol{\eta}} = \left(\left(\frac{\partial \theta_i}{\partial \eta_j} \right) \right) = \begin{pmatrix} (1 - \rho)\mathbf{I}_{k-1} & -\boldsymbol{\pi} + \mathbf{e}_1 \\ \vdots & \vdots \\ (1 - \rho)\mathbf{I}_{k-1} & -\boldsymbol{\pi} + \mathbf{e}_{k-1} \\ (1 - \rho)\mathbf{I}_{k-1} & -\boldsymbol{\pi} \\ \mathbf{I}_{k-1} & \mathbf{0} \end{pmatrix}.$$

Using the relations

$$S(\boldsymbol{\eta}) = \frac{\partial}{\partial \boldsymbol{\eta}} \log f(\mathbf{x}; \boldsymbol{\theta}) = \left(\frac{\partial \boldsymbol{\theta}}{\partial \boldsymbol{\eta}} \right)^T \frac{\partial}{\partial \boldsymbol{\theta}} \log f(\mathbf{x}; \boldsymbol{\theta}),$$

$$\mathcal{I}(\boldsymbol{\eta}) = \text{Var}(S(\boldsymbol{\eta})) = \left(\frac{\partial \boldsymbol{\theta}}{\partial \boldsymbol{\eta}} \right)^T \mathcal{I}(\boldsymbol{\theta}) \left(\frac{\partial \boldsymbol{\theta}}{\partial \boldsymbol{\eta}} \right),$$

it is possible to obtain an explicit form of the approximate FIM as $\tilde{\mathcal{I}}(\boldsymbol{\eta}) = ((a_{ij}))$, where

$$a_{ij} = \begin{cases} m(1-\rho)^2(\beta_i + \beta_k) - (\pi_i^{-1} + \pi_k^{-1}), & i = j, \quad i, j \in \{1, \dots, k-1\} \\ m(1-\rho)^2\beta_k - \pi_k^{-1}, & i \neq j, \quad i, j \in \{1, \dots, k-1\} \\ m(1-\rho)(\gamma_i - \gamma_k), & j = k, \quad i \in \{1, \dots, k-1\} \\ \frac{m}{(1-\rho)} \sum_{i=1}^k \pi_i(1-\pi_i) [(1-\rho)\pi_i + \rho]^{-1}, & i = k, j = k \end{cases}$$

and

$$\beta_i = \frac{\pi_i}{(1-\rho)\pi_i + \rho} + \frac{1-\pi_i}{(1-\rho)\pi_i}, \quad \gamma_i = \frac{\pi_i(1-\pi_i)}{(1-\rho)\pi_i + \rho} + \frac{\pi_i}{(1-\rho)}, \quad i = 1, \dots, k.$$

It can be shown rigorously that $\tilde{\mathcal{I}}(\boldsymbol{\eta}) - \mathcal{I}(\boldsymbol{\eta}) \rightarrow \mathbf{0}$ as $m \rightarrow \infty$, as stated in (Morel and Nagaraj, 1993), and proved in detail in (Morel and Nagaraj, 1991). The proof is similar in spirit to the proof of Theorem 2.1. We then have AFSA iterations for RCM,

$$\boldsymbol{\eta}^{(g+1)} = \boldsymbol{\eta}^{(g)} + \tilde{\mathcal{I}}^{-1}(\boldsymbol{\eta}^{(g)}) S(\boldsymbol{\eta}^{(g)}), \quad g = 1, 2, \dots$$

The following example involves a mixture of multinomials where the response probabilities are functions of covariates. The idea is analogous to the usual multinomial with logit link, but with links corresponding to each component of the mixture.

Example 3.2. In practice there are often covariates to be linked into the model. As an example for how AFSA can be applied, consider the following fixed effect model for response $\mathbf{Y} \sim \text{MultMix}_k(\boldsymbol{\theta}(\mathbf{x}), m)$ with $d \times 1$ covariates \mathbf{x} and \mathbf{z} . To each \mathbf{p}_ℓ vector, a generalized logit link will be added

$$\log \frac{p_{\ell j}(\mathbf{x})}{p_{\ell k}(\mathbf{x})} = \eta_{\ell j}, \quad \eta_{\ell j} = \mathbf{x}^T \boldsymbol{\beta}_{\ell j}, \quad \text{for } \ell = 1, \dots, s \text{ and } j = 1, \dots, k-1.$$

A proportional odds model will be assumed for $\boldsymbol{\pi}$,

$$\log \frac{\pi_1(\mathbf{z}) + \dots + \pi_\ell(\mathbf{z})}{\pi_{\ell+1}(\mathbf{z}) + \dots + \pi_s(\mathbf{z})} = \eta_\ell^\pi, \quad \eta_\ell^\pi = \nu_\ell + \mathbf{z}^T \boldsymbol{\alpha}, \quad \text{for } \ell = 1, \dots, s-1,$$

taking $\eta_0^\pi := -\infty$ and $\eta_s^\pi := \infty$. The unknown parameters are the $d \times 1$ vectors $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}_{\ell j}$, and the scalars ν_ℓ . Denote these parameters collectively as

$$\mathbf{B} = \begin{pmatrix} \boldsymbol{\beta}_1 \\ \vdots \\ \boldsymbol{\beta}_s \\ \boldsymbol{\nu} \\ \boldsymbol{\alpha} \end{pmatrix}, \quad \text{where } \boldsymbol{\beta}_\ell = \begin{pmatrix} \beta_{\ell 1} \\ \vdots \\ \beta_{\ell, k-1} \end{pmatrix} \text{ and } \boldsymbol{\nu} = \begin{pmatrix} \nu_1 \\ \vdots \\ \nu_s \end{pmatrix}.$$

Expressions for the $\boldsymbol{\theta}$ parameters can be obtained as

$$p_{\ell j}(\mathbf{x}) = \frac{e^{\eta_{\ell j}}}{1 + \sum_{b=1}^{k-1} e^{\eta_{\ell b}}} \quad \text{for } \ell = 1, \dots, s \text{ and } j = 1, \dots, k-1,$$

$$\pi_{\ell}(\mathbf{z}) = \frac{e^{\eta_{\ell}^{\pi}}}{1 + e^{\eta_{\ell}^{\pi}}} - \frac{e^{\eta_{\ell-1}^{\pi}}}{1 + e^{\eta_{\ell-1}^{\pi}}} \quad \text{for } \ell = 1, \dots, s.$$

To implement AFSA, a score vector and approximate FIM are needed. For the score vector we have

$$S(\mathbf{B}) = \frac{\partial}{\partial \mathbf{B}} \log f(\mathbf{y}; \boldsymbol{\theta}) = \left(\frac{\partial \mathbf{N}}{\partial \mathbf{B}} \right)^T \left(\frac{\partial \boldsymbol{\theta}}{\partial \mathbf{N}} \right)^T \frac{\partial}{\partial \boldsymbol{\theta}} \log f(\mathbf{y}; \boldsymbol{\theta})$$

where $\mathbf{N} = (\boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_s, \boldsymbol{\eta}_{\pi})$, $\boldsymbol{\eta}_{\ell} = (\eta_{\ell 1}, \dots, \eta_{\ell, k-1})$, and $\boldsymbol{\eta}_{\pi} = (\eta_{\pi 1}^{\pi}, \dots, \eta_{\pi, s-1}^{\pi})$. For the FIM we have

$$\mathcal{I}(\mathbf{B}) = \text{Var}(S(\mathbf{B})) = \left(\frac{\partial \mathbf{N}}{\partial \mathbf{B}} \right)^T \left(\frac{\partial \boldsymbol{\theta}}{\partial \mathbf{N}} \right)^T \mathcal{I}(\boldsymbol{\theta}) \left(\frac{\partial \boldsymbol{\theta}}{\partial \mathbf{N}} \right) \left(\frac{\partial \mathbf{N}}{\partial \mathbf{B}} \right).$$

Finding expressions for the two Jacobians is tedious but straightforward.

Propositions 3.3 and 3.4 and Theorem 3.5 state consequences of the main approximation result, which have significant implications on the computation of MLEs. We have already seen that the approximate FIM is equivalent to a complete data FIM from EM. There is also an interesting connection between AFSA and EM, in that the iterations are algebraically related. To see this connection, explicit forms for AFSA and EM iterations are first presented, with proofs given in Appendix B.

Proposition 3.3 (AFSA Iterations). *The AFSA iterations*

$$\boldsymbol{\theta}^{(g+1)} = \boldsymbol{\theta}^{(g)} + \tilde{\mathcal{I}}^{-1}(\boldsymbol{\theta}^{(g)}) S(\boldsymbol{\theta}^{(g)}), \quad g = 1, 2, \dots \quad (3.3)$$

can be written explicitly as

$$\pi_{\ell}^{(g+1)} = \pi_{\ell}^{(g)} \frac{1}{n} \sum_{i=1}^n \frac{P_{\ell}(\mathbf{x}_i)}{P(\mathbf{x}_i)}, \quad \ell = 1, \dots, s$$

$$p_{\ell j}^{(g+1)} = \frac{1}{M} \sum_{i=1}^n \frac{P_{\ell}(\mathbf{x}_i)}{P(\mathbf{x}_i)} x_{ij} - p_{\ell j}^{(g)} \left[1 - \frac{1}{M} \sum_{i=1}^n m_i \frac{P_{\ell}(\mathbf{x}_i)}{P(\mathbf{x}_i)} \right], \quad \ell = 1, \dots, s, \quad j = 1, \dots, k.$$

where $M = m_1 + \dots + m_n$.

Proposition 3.4 (EM Iterations). *Consider the complete data*

$$Z_i = \begin{cases} 1 & \text{with probability } \pi_1 \\ \vdots & \\ s & \text{with probability } \pi_s, \end{cases} \quad \text{and} \quad (\mathbf{X}_i | Z_i = \ell) \sim \text{Mult}_k(\mathbf{p}_{\ell}, m_i),$$

where (\mathbf{X}_i, Z_i) are independent for $i = 1, \dots, n$. Denote $\gamma_{i\ell}^{(g)} := \mathbb{P}(Z_i = \ell \mid \mathbf{x}_i, \boldsymbol{\theta}^{(g)})$ as the posterior probability that the i th observation belongs to the ℓ th group. Iterations for an EM algorithm are given by

$$\begin{aligned}\pi_\ell^{(g+1)} &= \frac{1}{n} \sum_{i=1}^n \gamma_{i\ell}^{(g)} = \frac{1}{n} \pi_\ell^{(g)} \sum_{i=1}^n \frac{\mathbb{P}_\ell(\mathbf{x}_i)}{\mathbb{P}(\mathbf{x}_i)}, \quad \ell = 1, \dots, s, \\ p_{\ell j}^{(g+1)} &= \frac{\sum_{i=1}^n x_{ij} \gamma_{i\ell}^{(g)}}{\sum_{i=1}^n m_i \gamma_{i\ell}^{(g)}} = \frac{\sum_{i=1}^n x_{ij} \frac{\mathbb{P}_\ell(\mathbf{x}_i)}{\mathbb{P}(\mathbf{x}_i)}}{\sum_{i=1}^n m_i \frac{\mathbb{P}_\ell(\mathbf{x}_i)}{\mathbb{P}(\mathbf{x}_i)}}, \quad \ell = 1, \dots, s, \quad j = 1, \dots, k.\end{aligned}$$

The iterations for AFSA or EM are repeated for $g = 1, 2, \dots$, with a given initial guess $\boldsymbol{\theta}^{(0)}$, until

$$|\log L(\boldsymbol{\theta}^{(g+1)}) - \log L(\boldsymbol{\theta}^{(g)})| < \varepsilon,$$

where $\varepsilon > 0$ is a given tolerance, which is taken to be the stopping criteria for the remainder of this paper.

Theorem 3.5. *Denote the estimator from EM by $\hat{\boldsymbol{\theta}}$, and the estimator from AFSA by $\tilde{\boldsymbol{\theta}}$. Suppose cluster sizes are equal, so that $m_1 = \dots = m_n = m$. If the two algorithms start at the g th iteration with $\boldsymbol{\theta}^{(g)}$, then for the $(g+1)$ th iteration,*

$$\tilde{\pi}_\ell^{(g+1)} = \hat{\pi}_\ell^{(g+1)} \quad \text{and} \quad \tilde{p}_{\ell j}^{(g+1)} = \left(\frac{\hat{\pi}_\ell^{(g+1)}}{\pi_\ell^{(g)}} \right) \hat{p}_{\ell j}^{(g+1)} + \left(1 - \frac{\hat{\pi}_\ell^{(g+1)}}{\pi_\ell^{(g)}} \right) p_{\ell j}^{(g)}$$

for $\ell = 1, \dots, s$ and $j = 1, \dots, k$.

Proof of Theorem 3.5. It is immediate from Propositions 3.3 and 3.4 that $\tilde{\pi}_\ell^{(g+1)} = \hat{\pi}_\ell^{(g+1)}$, and that

$$\frac{\hat{\pi}_\ell^{(g+1)}}{\pi_\ell^{(g)}} = \frac{1}{n} \sum_{i=1}^n \frac{\mathbb{P}_\ell(\mathbf{x}_i)}{\mathbb{P}(\mathbf{x}_i)}.$$

Now,

$$\begin{aligned}& \left(\frac{\hat{\pi}_\ell^{(g+1)}}{\pi_\ell^{(g)}} \right) \hat{p}_{\ell j}^{(g+1)} + \left(1 - \frac{\hat{\pi}_\ell^{(g+1)}}{\pi_\ell^{(g)}} \right) p_{\ell j}^{(g)} \\ &= \frac{\sum_{i=1}^n x_{ij} \frac{\mathbb{P}_\ell(\mathbf{x}_i)}{\mathbb{P}(\mathbf{x}_i)}}{m \sum_{i=1}^n \frac{\mathbb{P}_\ell(\mathbf{x}_i)}{\mathbb{P}(\mathbf{x}_i)}} \frac{1}{n} \sum_{i=1}^n \frac{\mathbb{P}_\ell(\mathbf{x}_i)}{\mathbb{P}(\mathbf{x}_i)} + p_{\ell j}^{(g)} \left(1 - \frac{1}{n} \sum_{i=1}^n \frac{\mathbb{P}_\ell(\mathbf{x}_i)}{\mathbb{P}(\mathbf{x}_i)} \right) \\ &= \frac{1}{mn} \sum_{i=1}^n \frac{\mathbb{P}_\ell(\mathbf{x}_i)}{\mathbb{P}(\mathbf{x}_i)} x_{ij} + p_{\ell j}^{(g)} \left(1 - \frac{1}{n} \sum_{i=1}^n \frac{\mathbb{P}_\ell(\mathbf{x}_i)}{\mathbb{P}(\mathbf{x}_i)} \right) \\ &= \tilde{p}_{\ell j}^{(g+1)}.\end{aligned} \tag{3.4}$$

□

The $(g + 1)$ th AFSA iterate can then be seen as a linear combination of the g th iterate and the $(g + 1)$ th step of EM. The coefficient $\hat{\pi}_\ell^{(g+1)}/\pi_\ell^{(g)}$ is non-negative but may be larger than 1. Therefore $\tilde{p}_{\ell j}^{(g+1)}$ need not lie strictly between $\hat{p}_{\ell j}^{(g+1)}$ and $p_{\ell j}^{(g)}$. Figure 1 shows a plot of $\tilde{p}_{\ell j}^{(g+1)}$ as the ratio $\hat{\pi}_\ell^{(g+1)}/\pi_\ell^{(g)}$ varies. However, suppose that at g th step the EM algorithm is close to convergence. Then

$$\hat{\pi}_\ell^{(g+1)} \approx \hat{\pi}_\ell^{(g)} \iff \frac{\hat{\pi}_\ell^{(g+1)}}{\hat{\pi}_\ell^{(g)}} \approx 1, \quad \text{for } \ell = 1, \dots, s.$$

From (3.4) we will also have

$$\tilde{p}_{\ell j}^{(g+1)} \approx \hat{p}_{\ell j}^{(g+1)}, \quad \text{for } \ell = 1, \dots, s, \text{ and } j = 1, \dots, k.$$

From this point on, AFSA and EM iterations are approximately the same. Hence, in the vicinity of a solution, AFSA and EM will produce the same estimate. Note that this result holds for any m , and does not require a large cluster size justification. For the case of varying cluster sizes m_1, \dots, m_n ,

$$\begin{aligned} & \frac{\hat{\pi}_\ell^{(g+1)}}{\pi_\ell^{(g)}} \hat{p}_{\ell j}^{(g+1)} + \left(1 - \frac{\hat{\pi}_\ell^{(g+1)}}{\pi_\ell^{(g)}}\right) p_{\ell j}^{(g)} \\ &= \frac{\sum_{i=1}^n x_{ij} \frac{P_\ell(\mathbf{x}_i)}{P(\mathbf{x}_i)}}{\sum_{i=1}^n m_i \frac{P_\ell(\mathbf{x}_i)}{P(\mathbf{x}_i)}} \frac{1}{n} \sum_{i=1}^n \frac{P_\ell(\mathbf{x}_i)}{P(\mathbf{x}_i)} + p_{\ell j}^{(g)} \left(1 - \frac{1}{n} \sum_{i=1}^n \frac{P_\ell(\mathbf{x}_i)}{P(\mathbf{x}_i)}\right), \end{aligned} \quad (3.5)$$

which does not simplify to $\tilde{p}_{\ell j}^{(g+1)}$ as in the proof of Theorem 3.5. However, this illustrates that EM and AFSA are still closely related. This also suggests an *ad-hoc* revision to AFSA, letting $\tilde{p}_{\ell j}^{(g+1)}$ equal (3.5) so that the algebraic relationship to EM would be maintained as in (3.4) for the balanced case.

A more general connection is known between EM and iterations of the form

$$\boldsymbol{\theta}^{(g+1)} = \boldsymbol{\theta}^{(g)} + \mathcal{I}_c^{-1}(\boldsymbol{\theta}^{(g)}) S(\boldsymbol{\theta}^{(g)}), \quad g = 1, 2, \dots, \quad (3.6)$$

where $\mathcal{I}_c(\boldsymbol{\theta})$ is a complete data FIM. Titterton (1984) shows that the two iterations are approximately equivalent under appropriate regularity conditions. The equivalence is exact when the complete data likelihood is a regular exponential family

$$L(\boldsymbol{\mu}) = \exp \left\{ b(\mathbf{x}) + \boldsymbol{\eta}^T \mathbf{t} + a(\boldsymbol{\eta}) \right\}, \quad \boldsymbol{\eta} = \boldsymbol{\eta}(\boldsymbol{\mu}), \quad \mathbf{t} = \mathbf{t}(\mathbf{x}),$$

and $\boldsymbol{\mu} := E(\mathbf{t}(\mathbf{X}))$ is the parameter of interest. The complete data likelihood for our multinomial mixture is indeed a regular exponential family, but the parameter of interest $\boldsymbol{\theta}$ is a transformation of $\boldsymbol{\mu}$ rather than $\boldsymbol{\mu}$ itself. Therefore the equivalence is approximate, as we have seen in Theorem 3.5. The justification for AFSA leading to this paper followed the historical approach of Blischke (1964), and not from the role of $\tilde{\mathcal{I}}(\boldsymbol{\theta})$ as a complete data FIM. But the relationship between EM and the iterations (3.6) suggests that AFSA is a reasonable approach for finite mixtures beyond the multinomial setting,

AFSA step compared to previous iterate and EM step

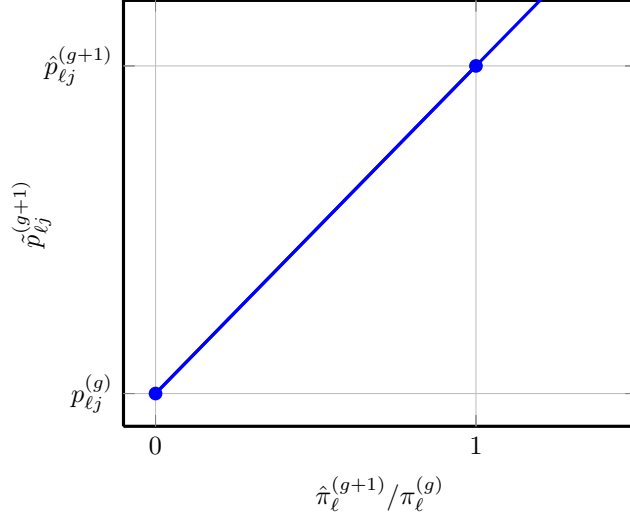


Figure 1: The next AFSA $\hat{p}_{\ell_j}^{(g+1)}$ iteration is a linear combination of $\hat{p}_{\ell_j}^{(g+1)}$ and $p_{\ell_j}^{(g)}$, which depends on the ratio $\hat{\pi}_\ell^{(g+1)} / \pi_\ell^{(g)}$.

4 Simulation Studies

The main result stated in Theorem 2.1 allows us to approximate the matrix $\mathcal{I}(\boldsymbol{\theta})$ by $\tilde{\mathcal{I}}(\boldsymbol{\theta})$, which is much more easily computed. Theorem 2.5 justifies $\tilde{\mathcal{I}}^{-1}(\boldsymbol{\theta})$ as an approximation for the inverse FIM. In the present section, simulation studies investigate the quality of the two approximations as a function of m . We also present studies to demonstrate the convergence speed and solution quality of AFSA.

4.1 Distance between true and approximate FIM

Consider two concepts of distance to compare the closeness of the exact and approximate matrices. Based on the Frobenius norm $\|\mathbf{A}\|_F = \sqrt{\sum_i \sum_j a_{ij}^2}$, a distance metric

$$d_F(\mathbf{A}, \mathbf{B}) = \|\mathbf{A} - \mathbf{B}\|_F$$

can be constructed using the sum of squared differences of corresponding elements. This distance will be larger in general when the magnitudes of the elements are larger, so we will also consider a scaled version

$$d_S(\mathbf{A}, \mathbf{B}) = \frac{d_F(\mathbf{A}, \mathbf{B})}{\|\mathbf{B}\|_F} = \sqrt{\frac{\sum_i \sum_j (a_{ij} - b_{ij})^2}{\sum_i \sum_j b_{ij}^2}},$$

noting that this is not a true distance metric since it is not symmetric. Using these two metrics, we compare the distance between true and approximate FIMs, and also the dis-

tance between their inverses. Consider a mixture $\text{MultMix}_2(\boldsymbol{\theta}, m)$ of three binomials, with parameters

$$\mathbf{p} = (1/7 \quad 1/3 \quad 2/3) \quad \text{and} \quad \boldsymbol{\pi} = (1/6 \quad 2/6 \quad 3/6).$$

Figure 2 plots the two distance types for both the FIM and inverse FIM as m varies. Note that distances are plotted on a log scale, so the vertical axis represents orders of magnitude. To see more concretely what is being compared, for the moderate cluster size $m = 20$ we have, respectively for the approximate and exact FIMs,

$$\begin{pmatrix} 27.222 & 0 & 0 & 0 & 0 \\ 0 & 30 & 0 & 0 & 0 \\ 0 & 0 & 45 & 0 & 0 \\ 0 & 0 & 0 & 8 & 2 \\ 0 & 0 & 0 & 2 & 5 \end{pmatrix} \text{ vs. } \begin{pmatrix} 14.346 & -2.453 & -0.184 & -3.341 & 1.625 \\ -2.453 & 12.605 & -6.749 & -4.440 & -0.944 \\ -0.184 & -6.749 & 34.175 & -1.205 & -2.914 \\ -3.341 & -4.440 & -1.205 & 6.022 & 2.536 \\ 1.625 & -0.944 & -2.914 & 2.536 & 3.621 \end{pmatrix}$$

and for the approximate and exact inverse FIMs,

$$\begin{pmatrix} 0.037 & 0 & 0 & 0 & 0 \\ 0 & 0.033 & 0 & 0 & 0 \\ 0 & 0 & 0.022 & 0 & 0 \\ 0 & 0 & 0 & 0.139 & -0.056 \\ 0 & 0 & 0 & -0.056 & 0.222 \end{pmatrix} \text{ vs. } \begin{pmatrix} 0.216 & 0.160 & 0.020 & 0.366 & -0.295 \\ 0.160 & 0.251 & 0.043 & 0.383 & -0.240 \\ 0.020 & 0.043 & 0.040 & 0.053 & -0.003 \\ 0.366 & 0.383 & 0.053 & 0.953 & -0.690 \\ -0.295 & -0.240 & -0.003 & -0.690 & 0.827 \end{pmatrix}.$$

Since the approximations are block diagonal matrices they have no way of capturing the off-diagonal blocks, which are present in the exact matrices but are eventually dominated by the block-diagonal elements as $m \rightarrow \infty$. This emphasizes one obvious disadvantage of the approximate FIM, which is that it cannot be used to estimate all the asymptotic covariances for the MLEs for a fixed cluster size. For this $m = 20$ case, the block-diagonal elements for both pairs of matrices are not very close, although they are at least the same order of magnitude with the same signs. The magnitudes of elements in the inverse FIMs are in general much smaller than those in the FIMs, so the unscaled distance will naturally be smaller between the inverses.

Now in Figure 2 consider the distance $d_F(\tilde{\mathcal{I}}(\boldsymbol{\theta}), \mathcal{I}(\boldsymbol{\theta}))$ as m is varied. For the FIM, the distance appears to be moderate at first, then increasing with m , and finally beginning to vanish as m becomes large. What is not reflected here is that the magnitudes of the elements themselves are increasing; this is inflating the distance until the convergence of Theorem 2.1 begins to kick in. Considering the scaled distance $d_S(\tilde{\mathcal{I}}(\boldsymbol{\theta}), \mathcal{I}(\boldsymbol{\theta}))$ helps to suppress the effect of the element magnitudes and gives a clearer picture of the convergence.

Focusing next on the inverse FIM, consider the distance $d_F(\tilde{\mathcal{I}}^{-1}(\boldsymbol{\theta}), \mathcal{I}^{-1}(\boldsymbol{\theta}))$. For $m < 5$ the exact FIM is computationally singular, so its inverse cannot be computed. Note that in this case the conditions for identifiability are not satisfied (see Appendix A). This is not just a coincidence; there is a known relationship between model non-identifiability and singularity of the FIM (Rothenberg, 1971). For m between 5 and about 23, the distance is very large at first because of near-singularity of the FIM, but quickly returns to a reasonable magnitude. As m increases further, the distance quickly vanishes toward zero. We also consider the

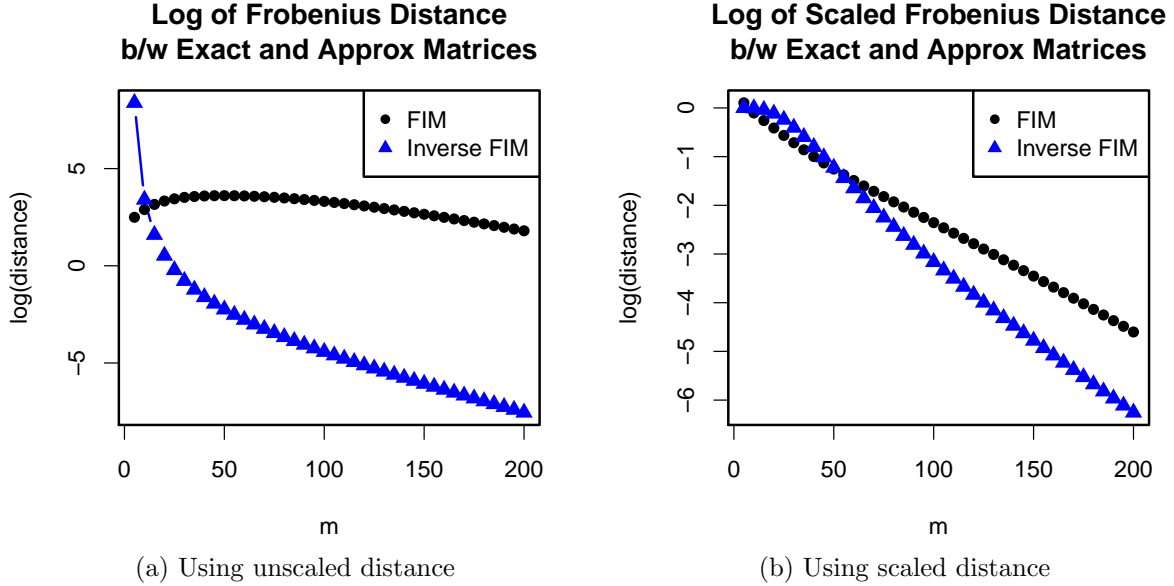


Figure 2: Distance between exact and approximate FIM and its inverse, as m is varied.

scaled distance $d_S(\tilde{\mathcal{I}}^{-1}(\boldsymbol{\theta}), \mathcal{I}^{-1}(\boldsymbol{\theta}))$. Again, this helps to remove the effects of the element magnitudes, which are becoming very small as m increases. Even after taking into account the scale of the elements, the distance between the inverse matrices appears to be converging more quickly than the distance between the FIM and approximate FIM.

4.2 Effectiveness of AFSA method

4.2.1 Convergence Speed

We first observe the convergence speed of AFSA and several of its competitors. Consider the mixture of two trinomials

$$\begin{aligned} \mathbf{Y}_i &\stackrel{\text{iid}}{\sim} \text{MultMix}_3(\boldsymbol{\theta}, m = 20), & i = 1, \dots, n = 500 \\ \mathbf{p}_1 &= (1/3 \quad 1/3 \quad 1/3), & \mathbf{p}_2 = (0.1 \quad 0.3 \quad 0.6), & \pi = 0.75. \end{aligned}$$

We fit the MLE using AFSA, FSA, and EM. After the g th iteration, the quantity

$$\delta^{(g)} = \log L(\boldsymbol{\theta}^{(g)}) - \log L(\boldsymbol{\theta}^{(g-1)})$$

is measured. The sequence $\log |\delta^{(g)}|$ is plotted for each algorithm in Figure 3. Note that $\delta^{(g)}$ may be negative, except for example in EM which guarantees an improvement to the log-likelihood in every step. A negative $\delta^{(g)}$ can be interpreted as negative progress, at least from a local maximum. The absolute value is taken to make plotting possible on the log scale, but some steps with negative progress have been obscured. The resulting estimates and

standard errors for all algorithms are shown in Table 1, and additional summary information is shown in Table 2.

We see that AFSA and EM have almost exactly the same rate of convergence toward the same solution, as suggested by Theorem 3.5. FSA had severe problems, and was not able to converge within 100 iterations; i.e. $\delta^{(g)} < 10^{-8}$ was not attained. The situation for FSA is worse than it appears in the plot. Although $\log |\delta^{(g)}|$ is becoming small, FSA’s steps result in both positive and negative $\delta^{(g)}$ ’s until the iteration limit is reached. This indicates a failure to approach any maximum of the log-likelihood.

We also considered an FSA hybrid with a “warmup period”, where for a given $\varepsilon_0 > 0$ the approximate FIM is used until the first time $\delta^{(g)} < \varepsilon_0$ is crossed. Notice that $\varepsilon_0 = \infty$ corresponds to “no warmup period”. A similar idea has been considered by Neerchal and Morel (2005), who proposed a two-stage procedure for AFSA in the RCM setting of Example 3.1. The first stage consisted of running AFSA iterations until convergence, and in the second stage one additional iteration of exact Fisher Scoring was performed. The purpose of the FSA iteration was to improve standard error estimates, which were previously found to be inaccurate when computed directly from the approximate FIM (Neerchal and Morel, 1998). Here we note that FSA also offers a faster convergence rate than AFSA, given an initial path to a solution. Therefore, AFSA can be used in early iterations to move to the vicinity of a solution, then a switch to FSA will give an accelerated converge to the solution. This approach depends on the exact FIM being feasible to compute, so the sample space cannot be too large. For the present simulations, we make use of the naive summation (2.3). Hence, there is a trade-off in the choice of ε_0 between energy spent on computing the exact FIM and a larger number of iterations required for AFSA. Figure 3 shows that the hybrid strategy is effective, addressing the erratic behavior of FSA from an arbitrary starting value and the slower convergence rates of EM and AFSA. Table 2 shows that even a very limited warmup period such as $\varepsilon_0 = 10$ can be sufficient.

The Newton-Raphson algorithm, which has not been shown here, performed similarly to Fisher Scoring but has issues with singularity in some samples. Standard errors for AFSA were obtained as $\sqrt{a^{11}}, \dots, \sqrt{a^{qq}}$, denoting $\tilde{\mathcal{I}}^{-1}(\hat{\boldsymbol{\theta}}) = ((a^{ij}))$. For FSA and FSA-Hybrid, the inverse of the exact FIM was used instead. The basic EM algorithm does not yield standard error estimates. Several extensions have been proposed to address this, such as by Louis (1982) and Meng and Rubin (1991). In light of Theorem 3.5, standard errors from $\tilde{\mathcal{I}}^{-1}(\boldsymbol{\theta})$ evaluated at EM estimates could also be used to obtain similar results to AFSA.

4.2.2 Monte Carlo Study

We next consider a Monte Carlo study of the difference between AFSA and EM estimators. Observations were generated from

$$\mathbf{Y}_i \stackrel{\text{iid}}{\sim} \text{MultMix}_k(\boldsymbol{\theta}, m_i), \quad i = 1, \dots, n = 500,$$

given varying cluster sizes m_1, \dots, m_n which themselves were generated as

$$Z_1, \dots, Z_n \stackrel{\text{iid}}{\sim} \text{Gamma}(\alpha, \beta), \quad m_i = \lceil Z_i \rceil. \quad (4.1)$$

Convergence of competing algorithms

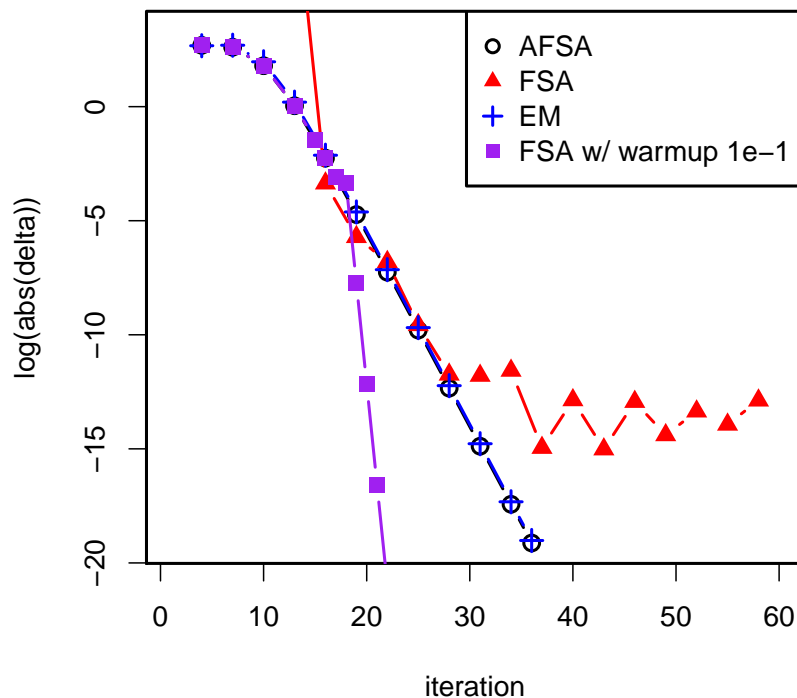


Figure 3: Convergence of several competing algorithms for a small test problem

Table 1: Estimates and standard errors for the competing algorithms. FSA Hybrid produced the same results with ε_0 set to 0.001, 0.01, 0.1, 1, and 10.

	FSA	AFSA	EM	FSA Hybrid
\hat{p}_{11}	0.2690	0.3315	0.3315	0.3315
SE	0.0044	0.0055	—	0.0065
\hat{p}_{12}	0.3341	0.3403	0.3403	0.3403
SE	0.0047	0.0055	—	0.0057
\hat{p}_{21}	0.3341	0.0928	0.0928	0.0928
SE	0.0908	0.0057	—	0.0081
\hat{p}_{22}	0.0002	0.3167	0.3167	0.3167
SE	0.0055	0.0091	—	0.0102
$\hat{\pi}$	0.9990	0.7381	0.7381	0.7381
SE	0.0014	0.0197	—	0.0253

Table 2: Convergence of several competing algorithms. Hybrid FSA is shown with several choices of the warmup tolerance ε_0 . Exact FSA uses $\varepsilon_0 = \infty$.

method	ε_0	logLik	tol	iter
AFSA	—	-2234.655	4.94×10^{-09}	36
EM	—	-2234.655	5.50×10^{-09}	36
FSA	∞	-2423.864	-1.26×10^{-07}	100
FSA	10	-2234.655	4.46×10^{-10}	16
FSA	1	-2234.655	1.34×10^{-10}	20
FSA	0.1	-2234.655	7.84×10^{-10}	22
FSA	0.01	-2234.655	1.42×10^{-10}	24
FSA	0.001	-2234.655	9.05×10^{-10}	26

Several different settings of θ are considered, with $s = 2$ mixing components and proportion $\pi = 0.75$ for the first component. The parameters α and β were chosen such that $E(Z_i) = \alpha\beta = 20$. This gives $\beta = 20/\alpha$ so only α is free, and $\text{Var}(Z_i) = \alpha\beta^2 = 400/\alpha$ can be chosen as desired. The expectation and variance of m_i are intuitively similar to Z_i , and their exact values may be computed numerically.

Once the n observations are generated, an AFSA estimator $\tilde{\theta}$ and an EM estimator $\hat{\theta}$ are fit. This process is repeated 1000 times yielding $\tilde{\theta}^{(r)}$ and $\hat{\theta}^{(r)}$ for $r = 1, \dots, 1000$. A default initial value was selected for each setting of θ , and used for both algorithms in every repetition. To measure the closeness of the two estimators, a maximum relative difference is taken over all components of θ , then averaged over all repetitions:

$$\bar{D} = \frac{1}{1000} \sum_{r=1}^{1000} D_r, \quad \text{where } D_r = \bigvee_{j=1}^q \left| \frac{\tilde{\theta}_j^{(r)} - \hat{\theta}_j^{(r)}}{\tilde{\theta}_j^{(r)}} \right|.$$

Here \bigvee represents the “maximum” operator. Notice that obtaining a good result for \bar{D} depends on the vectors $\hat{\theta}$ and $\tilde{\theta}$ being ordered in the same way. To help ensure this, we add the constraint

$$\pi_1 > \dots > \pi_s,$$

which is enforced in both algorithms by reordering the estimates for π_1, \dots, π_s and $\mathbf{p}_1, \dots, \mathbf{p}_s$ accordingly after every iteration. Table 3 shows the results of the simulation. Nine different scenarios for θ are considered. The cluster sizes m_1, \dots, m_n are selected in three different ways: a balanced case where $m_i = 20$ for $i = 1, \dots, n$, cluster sizes selected at random with small variability (using $\alpha = 100$), and cluster sizes selected at random with moderate variability (using $\alpha = 25$).

Both algorithms are susceptible to finding local maxima of the likelihood, but in this experiment AFSA encountered the problem much more frequently. These cases stood out because the local maxima occurred with one of the mixing proportions or category probabilities close to zero, i.e. a convergence to the boundary of the parameter space. This is an especially bad situation for our Monte Carlo statistic \bar{D} , which can become very large if

Table 3: Closeness between AFSA and EM estimates, over 1000 trials

		Cluster sizes equal	$\alpha = 100$	$\alpha = 25$	
\mathbf{p}_1	\mathbf{p}_2	$m_i = 20$	$\text{Var}(m_i) \approx 4.083$	$\text{Var}(m_i) \approx 16.083$	
A.	(0.1)	(0.5)	2.178×10^{-6}	2.019×10^{-6}	2.080×10^{-6}
B.	(0.3)	(0.5)	4.073×10^{-5}	3.501×10^{-5}	3.890×10^{-5}
C.	(0.35)	(0.5)	8.683×10^{-4}	2.625×10^{-4}	2.738×10^{-4}
D.	(0.4)	(0.5)	9.954×10^{-3}	6.206×10^{-2}	6.563×10^{-2}
E.	(0.1, 0.3)	(1/3, 1/3)	1.342×10^{-3}	1.009×10^{-3}	1.878×10^{-3}
F.	(0.1, 0.5)	(1/3, 1/3)	1.408×10^{-6}	1.338×10^{-6}	1.334×10^{-6}
G.	(0.3, 0.5)	(1/3, 1/3)	3.884×10^{-6}	3.943×10^{-6}	3.885×10^{-6}
H.	(0.1, 0.1, 0.3)	(0.25, 0.25, 0.25)	8.389×10^{-7}	8.251×10^{-7}	8.440×10^{-7}
I.	(0.1, 0.2, 0.3)	(0.25, 0.25, 0.25)	1.523×10^{-6}	1.472×10^{-6}	1.408×10^{-6}

this occurs even once for a given scenario. The problem occurred most frequently for the case $\mathbf{p}_1 = (0.1, 0.3)$ and $\mathbf{p}_2 = (1/3, 1/3)$. To counter this, we restarted AFSA with a random starting value whenever a solution with any estimate less than 0.01 was obtained. For this experiment, no more than 15 out of 1000 trials required a restart, and no more than two restarts were needed for the same trial. In practice, we recommend starting AFSA with several initial values to ensure that any solutions on the boundary are not missteps taken by the algorithm.

The entries in Table 3 show that small to moderate variation of the cluster sizes does not have a significant impact on the equivalence of AFSA and EM. On the other hand, as \mathbf{p}_1 and \mathbf{p}_2 are moved closer together, the quantity \bar{D} tends to become larger. Theorem 2.1 depends on the distinctness of the category probability vectors, so the quality of the FIM approximation at moderate cluster size may begin to suffer in this case. The estimation problem itself also intuitively becomes more difficult as \mathbf{p}_1 and \mathbf{p}_2 become closer. Recall that the dimension of \mathbf{p}_i is $k - 1$; it can be seen from Table 3 that increasing k from 2 to 4 does not necessarily have a negative effect on the results.

In Scenario E, \mathbf{p}_1 and \mathbf{p}_2 are not too close together, yet \bar{D} has a similar magnitude to Scenario D where the two vectors are closer. Figure 4 shows a plot of the individual D_r for Scenarios D and E. Notice that in Scenario E, one particular simulation in each case is responsible for the large magnitude of \bar{D} . Upon removal of these simulations, the order of \bar{D} is reduced from 10^{-3} to 10^{-6} . However, many large D_r were present in the Scenario D results.

5 Conclusions

A large cluster approximation was presented for the FIM of the finite mixture of multinomials model (Theorem 2.1). This matrix has a convenient block diagonal form, where each non-zero block is the FIM of a standard multinomial observation. Furthermore, the approximation is equivalent to the “complete data” FIM, had population labels been recorded for each

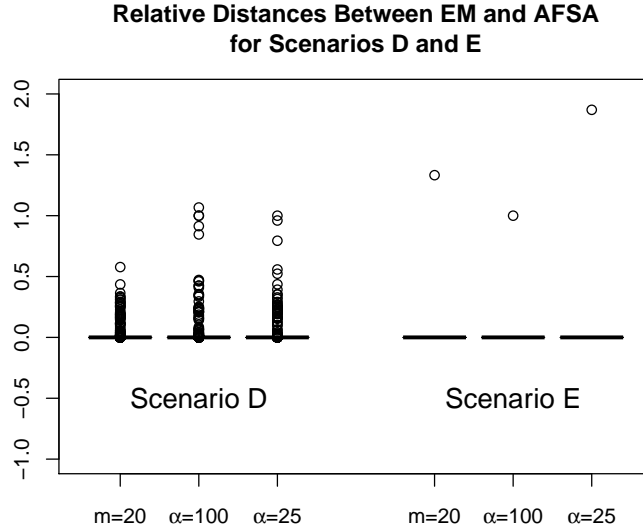


Figure 4: Boxplots for Scenarios D and E of Monte Carlo study. At this scale, the boxes appear as thin horizontal lines.

observation (Proposition 2.2). Using this approximation to the FIM, we formulated the Approximate Fisher Scoring Algorithm (AFSA), and showed that its iterations are closely related to the well known Expectation-Maximization (EM) algorithm for finite mixtures (Theorem 3.5). Simulations show that a rather large cluster size is needed before the exact and approximate FIM are close; this is not surprising given that a block diagonal matrix is being used to approximate a dense matrix. A large cluster size is also needed for a close approximation of the inverse, although the inverses are seen to converge together more quickly. Therefore, the approximate FIM and its inverse are not well-suited to replace the exact matrices for general usage. This means, for example, that one should be cautious about computing standard errors for the MLE from the approximate inverse FIM.

As another example of a general use for the approximate FIM, consider approximate $(1 - \alpha)$ level Wald-type and Score-type confidence regions,

$$\left\{ \boldsymbol{\theta}_0 : (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)^T \tilde{\mathcal{I}}(\hat{\boldsymbol{\theta}}) (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \leq \chi_{q,\alpha}^2 \right\} \quad \text{and} \quad \left\{ \boldsymbol{\theta}_0 : S(\boldsymbol{\theta}_0)^T \tilde{\mathcal{I}}^{-1}(\boldsymbol{\theta}_0) S(\boldsymbol{\theta}_0) \leq \chi_{q,\alpha}^2 \right\}, \quad (5.1)$$

respectively, using the approximate FIM in place of the exact FIM. Such regions are very practical to compute, but will likely not have the desired coverage for $\boldsymbol{\theta}$. However, we might expect the Score region to perform better for moderate cluster sizes because it involves the inverse matrix. On the other hand, the approximate FIM works well as a tool for estimation in the AFSA algorithm. This is interesting because the more standard Fisher Scoring and Newton-Raphson algorithms do not work well on their own. For Newton-Raphson, the invertibility of the Hessian depends on the sample as well as the current iterate $\boldsymbol{\theta}^{(g)}$ and the model. Fisher Scoring can be computed when the cluster size is not too small (so that

the FIM is non-singular), but it is often unable to make progress at all from an arbitrarily chosen starting point. In this case, AFSA (or EM) is useful for giving FSA some initial help. If FSA has a sufficiently good starting point, it can converge very quickly. Therefore we recommend a hybrid approach: use AFSA iterations for an initial warmup period, then switch to FSA once a path toward a solution has been established. This approach may also help to reduce the number of exact FIM computations needed, which may be expensive. Although AFSA and EM are closely related and often tend toward the same solution, AFSA is not restricted to the parameter space. Additional precautions may therefore be needed to prevent AFSA iterations from drifting outside of the space. AFSA also tended to converge to the boundary of the space more often than EM; hence, we reiterate the usual advice of trying several initial values as a good practice. AFSA may be preferable to EM in situations where it is more natural to formulate. Derivation of the E-step conditional log-likelihood may involve evaluating a complicated expectation, but is not required for AFSA. A trade-off for AFSA is that the score vector for the observed data must be computed; this may involve a messy differentiation, but is arguably easier to address numerically than the E-step. AFSA iterations were obtained for the Random-Clumped Multinomial in Example 3.1, starting from a general multinomial mixture and using an appropriate transformation of the parameters.

It is interesting to note the relationship between FSA, AFSA, and EM as Newton-type algorithms. Fisher Scoring is a classic algorithm where the Hessian is replaced by its expectation. In AFSA the Hessian is replaced instead by a complete data FIM. EM can be considered a Newton-type algorithm also, where the entire likelihood is replaced by a complete data likelihood with missing data integrated out. In this light, EM and AFSA iterations are seen to be approximately equivalent.

Several interesting questions can be raised at this point. There is a relationship between AFSA and EM which extends beyond the multinomial mixture; we wonder if the relationship between the exact and complete data information matrix generalizes as well. Also, for the present multinomial mixture, perhaps there is a small cluster bias correction that could be applied to improve the approximation. This might allow standard errors and confidence regions such as (5.1) to safely be derived from the approximate FIM.

6 Acknowledgements

The hardware used in the computational studies is part of the UMBC High Performance Computing Facility (HPCF). The facility is supported by the U.S. National Science Foundation through the MRI program (grant no. CNS-0821258) and the SCREMS program (grant no. DMS-0821311), with additional substantial support from the University of Maryland, Baltimore County (UMBC). See www.umbc.edu/hpcf for more information on HPCF and the projects using its resources. The first author additionally acknowledges financial support as HPCF RA.

References

- W. R. Blischke. Moment estimators for the parameters of a mixture of two binomial distributions. *The Annals of Mathematical Statistics*, 33(2):444–454, 1962.
- W. R. Blischke. Estimating the parameters of mixtures of binomial distributions. *Journal of the American Statistical Association*, 59(306):510–528, 1964.
- H. Bozdogan. Model selection and Akaike’s Information Criterion (AIC): The general theory and its analytical extensions. *Psychometrika*, 52(3):345–370, 1987.
- S. Chandra. On the mixtures of probability distributions. *Scandinavian Journal of Statistics*, 4:105–112, 1977.
- 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*, 39(1):1–38, 1977.
- A. B. Kabir. Estimation of parameters of a finite mixture of distributions. *Journal of the Royal Statistical Society, Series B*, 30:472–482, 1968.
- K. Krolikowska. Estimation of the parameters of any finite mixture of geometric distributions. *Demonstratio Mathematica*, 9:573–582, 1976.
- K. Lange. *Numerical Analysis for Statisticians*. Springer, 2nd edition, 2010.
- E. L. Lehmann and G. Casella. *Theory of Point Estimation*. Springer, 2nd edition, 1998.
- T. A. Louis. Finding the observed information matrix when using the EM algorithm. *Journal of the Royal Statistical Society, Series B*, 44:226–233, 1982.
- G. McLachlan and D. Peel. *Finite Mixture Models*. Wiley-Interscience, 2000.
- X. L. Meng and D. B. Rubin. Using EM to obtain asymptotic variance-covariance matrices: the SEM algorithm. *Journal of the American Statistical Association*, 86(416):899–909, 1991.
- C. D. Meyer. *Matrix Analysis and Applied Linear Algebra*. SIAM, 2001.
- J. G. Morel and N. K. Nagaraj. A finite mixture distribution for modeling multinomial extra variation. Technical Report Research report 91–03, Department of Mathematics and Statistics, University of Maryland, Baltimore County, 1991.
- J. G. Morel and N. K. Nagaraj. A finite mixture distribution for modelling multinomial extra variation. *Biometrika*, 80(2):363–371, 1993.
- J. G. Morel and N. K. Nagaraj. *Overdispersion Models in SAS*. SAS Institute, 2012.

- N. K. Neerchal and J. G. Morel. Large cluster results for two parametric multinomial extra variation models. *Journal of the American Statistical Association*, 93(443):1078–1087, 1998.
- N. K. Neerchal and J. G. Morel. An improved method for the computation of maximum likelihood estimates for multinomial overdispersion models. *Computational Statistics & Data Analysis*, 49(1):33–43, 2005.
- M. Okamoto. Some inequalities relating to the partial sum of binomial probabilities. *Annals of the Institute of Statistical Mathematics*, 10:29–35, 1959.
- A. M. Raim, M. K. Gobbert, N. K. Neerchal, and J. G. Morel. Maximum likelihood estimation of the random-clumped multinomial model as prototype problem for large-scale statistical computing. Accepted, 2012.
- C. R. Rao. *Linear statistical inference and its applications*. John Wiley and Sons Inc, 1965.
- H. Robbins. Mixture of distributions. *The Annals of Mathematical Statistics*, 19(3):360–369, 1948.
- T. J. Rothenberg. Identification in parametric models. *Econometrica*, 39:577–591, 1971.
- H. Teicher. On the mixture of distributions. *The Annals of Mathematical Statistics*, 31(1):55–73, 1960.
- D. M. Titterton. Recursive parameter estimation using incomplete data. *Journal of the Royal Statistical Society. Series B*, 46:257–267, 1984.
- H. Zhou and K. Lange. MM algorithms for some discrete multivariate distributions. *Journal of Computational and Graphical Statistics*, 19(3):645–665, 2010.

A Appendix: Preliminaries and Notation

Given an independent sample $\mathbf{X}_1, \dots, \mathbf{X}_n$ with joint likelihood $L(\boldsymbol{\theta})$ and $\boldsymbol{\theta}$ having dimension $q \times 1$, the score vector is

$$S(\boldsymbol{\theta}) = \frac{\partial}{\partial \boldsymbol{\theta}} \log L(\boldsymbol{\theta}) = \sum_{i=1}^n \frac{\partial}{\partial \boldsymbol{\theta}} \log f(\mathbf{x}; \boldsymbol{\theta}).$$

For $\mathbf{X}_i \sim \text{Mult}_k(\mathbf{p}, m)$ the score vector for a single observation can be obtained from

$$\begin{aligned} \frac{\partial}{\partial p_a} \log f(\mathbf{x}; \mathbf{p}, m) &= \frac{\partial}{\partial p_a} \left[x_1 \log p_1 + \dots + x_{k-1} \log p_{k-1} + x_k \log \left(1 - \sum_{j=1}^{k-1} p_j \right) \right] \\ &= x_a/p_a - x_k/p_k, \end{aligned} \tag{A.1}$$

so that

$$\frac{\partial}{\partial \mathbf{p}} \log f(\mathbf{x}; \mathbf{p}, m) = \begin{pmatrix} x_1/p_1 \\ \vdots \\ x_{k-1}/p_{k-1} \end{pmatrix} - \begin{pmatrix} x_k/p_k \\ \vdots \\ x_k/p_k \end{pmatrix} = \mathbf{D}^{-1} \mathbf{x}_{-k} - \frac{x_k}{p_k} \mathbf{1},$$

denoting $\mathbf{D} := \text{diag}(p_1, \dots, p_{k-1})$ and $\mathbf{x}_{-k} := (x_1, \dots, x_{k-1})$.

The score vector for a single observation $\mathbf{X} \sim \text{MultMix}_k(\boldsymbol{\theta}, m)$ can also be obtained,

$$\begin{aligned} \frac{\partial \log P(\mathbf{x})}{\partial \mathbf{p}_a} &= \frac{\partial \log \{ \sum_{\ell=1}^s \pi_\ell P_\ell(\mathbf{x}) \}}{\partial \mathbf{p}_a} \\ &= \frac{1}{P(\mathbf{x})} \pi_a \frac{\partial P_a(\mathbf{x})}{\partial \mathbf{p}_a} \\ &= \frac{\pi_a P_a(\mathbf{x})}{P(\mathbf{x})} \frac{\partial \log P_a(\mathbf{x})}{\partial \mathbf{p}_a} \\ &= \frac{\pi_a P_a(\mathbf{x})}{P(\mathbf{x})} \left[\mathbf{D}_a^{-1} \mathbf{x}_{-k} - \frac{x_k}{p_{ak}} \mathbf{1} \right], \quad a = 1, \dots, s, \end{aligned}$$

where $\mathbf{D}_a := \text{diag}(p_{a1}, \dots, p_{a,k-1})$, and

$$\begin{aligned} \frac{\partial \log P(\mathbf{x})}{\partial \pi_a} &= \frac{\partial \log \{ \sum_{\ell=1}^s \pi_\ell P_\ell(\mathbf{x}) \}}{\partial \pi_a} \\ &= \frac{P_a(\mathbf{x}) - P_s(\mathbf{x})}{P(\mathbf{x})}, \quad a = 1, \dots, s-1. \end{aligned}$$

Next, consider the $q \times q$ FIM for the independent sample $\mathbf{X}_1, \dots, \mathbf{X}_n$

$$\begin{aligned} \mathcal{I}(\boldsymbol{\theta}) &= \text{Var}(S(\boldsymbol{\theta})) = \text{E} \left[\left\{ \frac{\partial}{\partial \boldsymbol{\theta}} \log L(\boldsymbol{\theta}) \right\} \left\{ \frac{\partial}{\partial \boldsymbol{\theta}} \log L(\boldsymbol{\theta}) \right\}^T \right] \\ &= \text{E} \left[-\frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \log L(\boldsymbol{\theta}) \right]. \end{aligned}$$

The last equality holds under appropriate regularity conditions. For the multinomial FIM, we may use (A.1) to obtain

$$\frac{\partial}{\partial p_a} \frac{\partial}{\partial p_b} \log f(\mathbf{x}; \mathbf{p}, m) = \begin{cases} x_k/p_k^2 & \text{if } a \neq b \\ -x_a/p_a^2 - x_k/p_k^2 & \text{otherwise} \end{cases}$$

and so

$$\frac{\partial}{\partial \mathbf{p} \partial \mathbf{p}^T} \log f(\mathbf{x}; \mathbf{p}, m) = \text{diag} \left(-\frac{x_1}{p_1^2}, \dots, -\frac{x_{k-1}}{p_{k-1}^2} \right) - \frac{x_k}{p_k^2} \mathbf{1} \mathbf{1}^T.$$

Therefore, we have

$$\begin{aligned}\mathcal{I}(\mathbf{p}) &= \mathbb{E} \left(-\frac{\partial}{\partial \mathbf{p} \partial \mathbf{p}^T} \log f(\mathbf{x}; \mathbf{p}, m) \right) \\ &= \text{diag} \left(\frac{mp_1}{p_1^2}, \dots, \frac{mp_{k-1}}{p_{k-1}^2} \right) + \frac{mp_k}{p_k^2} \mathbf{1}\mathbf{1}^T \\ &= m (\mathbf{D}^{-1} + p_k^{-1} \mathbf{1}\mathbf{1}^T).\end{aligned}$$

The score vector and Hessian of the log-likelihood can be used to implement the Newton-Raphson algorithm, where the $(g+1)$ th iteration is given by

$$\boldsymbol{\theta}^{(g+1)} = \boldsymbol{\theta}^{(g)} - \left\{ \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \log L(\boldsymbol{\theta}^{(g)}) \right\}^{-1} S(\boldsymbol{\theta}^{(g)}).$$

The Hessian may be replaced with the FIM to implement Fisher Scoring

$$\boldsymbol{\theta}^{(g+1)} = \boldsymbol{\theta}^{(g)} + \mathcal{I}^{-1}(\boldsymbol{\theta}^{(g)}) S(\boldsymbol{\theta}^{(g)}).$$

In order for the estimation problem to be well-defined in the first place, the model must be identifiable. For finite mixtures, this is taken to mean that the equality

$$\sum_{\ell=1}^s \pi_\ell f(\mathbf{x}; \boldsymbol{\theta}_\ell) \stackrel{as}{=} \sum_{\ell=1}^v \lambda_\ell f(\mathbf{x}; \boldsymbol{\xi}_\ell)$$

implies $s = v$ and terms within the sums are equal, except the indices may be permuted (McLachlan and Peel, 2000, section 1.14). Chandra (1977) provides some insight into the identifiability issue, and shows that a family of multivariate mixtures is identifiable if any of the corresponding marginal mixtures are identifiable. In the present case, the multivariate mixtures consist of multinomial densities, and the marginal densities are binomials. It is well known that a finite mixture of s components from

$$\{ \text{Binomial}(m, \theta) : \theta \in (0, 1) \}$$

is identifiable if and only if $m \geq 2s - 1$; see, for example, Blischke (1964). Then a sufficient condition for model (2.2) to be identifiable is that $m_i \geq 2s - 1$ for at least one observation. This can be seen by the following lemma.

Lemma A.1. *Suppose $\mathbf{X}_i \stackrel{ind}{\sim} f_i(\mathbf{x}; \boldsymbol{\theta})$, $i = 1, \dots, n$, where f_i 's are densities, and for at least one $r \in \{1, \dots, n\}$ the family $\{f_r(\cdot; \boldsymbol{\theta}) : \boldsymbol{\theta} \in \Theta\}$ is identifiable. Then the joint model is identifiable.*

Proof of Lemma A.1. WLOG assume that $r = 1$, and suppose we have

$$\prod_{i=1}^n f_i(\mathbf{x}_i; \boldsymbol{\theta}) \stackrel{as}{=} \prod_{i=1}^n f_i(\mathbf{x}_i; \boldsymbol{\xi}).$$

Integrating both sides with respect to $\mathbf{x}_2, \dots, \mathbf{x}_n$, using the appropriate dominating measure,

$$f_1(\mathbf{x}_1; \boldsymbol{\theta}) \stackrel{as}{=} f_1(\mathbf{x}_1; \boldsymbol{\xi}).$$

Since the family $\{f_1(\cdot; \boldsymbol{\theta}) : \boldsymbol{\theta} \in \Theta\}$ is identifiable, this implies $\boldsymbol{\theta} = \boldsymbol{\xi}$. Hence the joint family $\{\prod_{i=1}^n f_i(\cdot; \boldsymbol{\theta}) : \boldsymbol{\theta} \in \Theta\}$ is identifiable. \square

B Appendix: Additional Proofs

To prove Theorem 2.1, we will first establish a key inequality. A similar strategy was used by Morel and Nagaraj (1991), but they considered the special case $k = s$, so that the number of mixture components is equal to the number of categories within each component. Here we generalize their argument to the general case where $k = s$ need not hold. The original proof was inspired by the following inequality from Okamoto (1959) for the tail probability of the binomial distribution, which was also considered by Blischke (1962).

Lemma B.1. *Suppose $X \sim \text{Binomial}(m, p)$ and let $f(x; m, p)$ be its density. Then for $c \geq 0$,*

- i. $P(X/m - p \geq c) \leq e^{-2mc^2}$,
- ii. $P(X/m - p \leq -c) \leq e^{-2mc^2}$.

Theorem B.2. *For a given index $b \in \{1, \dots, s\}$ we have*

$$\sum_{\mathbf{x} \in \Omega} \sum_{a \neq b}^s \frac{\pi_a P_a(\mathbf{x}) P_b(\mathbf{x})}{P(\mathbf{x})} \leq \frac{2}{\pi_b} \sum_{a \neq b}^s e^{-\frac{m}{2} \delta_{ab}^2},$$

where $\delta_{ab} = \sqrt{\sum_{j=1}^{k-1} (p_{aj} - p_{bj})^2}$.

Proof of Theorem B.2. For $a, b \in \{1, \dots, s\}$, assume WLOG that

$$\delta_{ab} := \sqrt{\sum_{j=1}^{k-1} (p_{aj} - p_{bj})^2} = (p_{aL} - p_{bL}), \quad \text{for some } L \in \{1, \dots, k-1\}$$

is positive. Denote as $\Omega(x_j)$ the multinomial sample space when the j th element of \mathbf{x} is fixed at a number x_j . Then we have

$$\begin{aligned} \sum_{\mathbf{x} \in \Omega} \frac{\pi_a P_a(\mathbf{x}) P_b(\mathbf{x})}{P(\mathbf{x})} &= \sum_{x_L=0}^m \sum_{\mathbf{x} \in \Omega(x_L)} \frac{\pi_a P_a(\mathbf{x}) P_b(\mathbf{x})}{P(\mathbf{x})} \\ &= \sum_{x_L \leq \frac{m}{2}(p_{aL} + p_{bL})} \sum_{\mathbf{x} \in \Omega(x_L)} \pi_a P_a(\mathbf{x}) \frac{P_b(\mathbf{x})}{P(\mathbf{x})} + \sum_{x_L > \frac{m}{2}(p_{aL} + p_{bL})} \sum_{\mathbf{x} \in \Omega(x_L)} \frac{\pi_a P_a(\mathbf{x})}{P(\mathbf{x})} P_b(\mathbf{x}) \\ &\leq \sum_{x_L \leq \frac{m}{2}(p_{aL} + p_{bL})} \sum_{\mathbf{x} \in \Omega(x_L)} \frac{\pi_a}{\pi_b} P_a(\mathbf{x}) + \sum_{x_L > \frac{m}{2}(p_{aL} + p_{bL})} \sum_{\mathbf{x} \in \Omega(x_L)} P_b(\mathbf{x}) \\ &= \frac{\pi_a}{\pi_b} \sum_{x_L \leq \frac{m}{2}(p_{aL} + p_{bL})} \sum_{\mathbf{x} \in \Omega(x_L)} P_a(\mathbf{x}) + \sum_{x_L > \frac{m}{2}(p_{aL} + p_{bL})} \sum_{\mathbf{x} \in \Omega(x_L)} P_b(\mathbf{x}). \end{aligned} \tag{B.1}$$

Notice that the last statement above consists of marginal probabilities for the L th coordinate of k -dimensional multinomials, which are binomial probabilities. Following Blischke (1962), suppose $A \sim \text{Binomial}(m, p_{aL})$ and $B \sim \text{Binomial}(m, p_{bL})$, then (B.1) is equal to

$$\frac{\pi_a}{\pi_b} P \left\{ A \leq \frac{m}{2}(p_{aL} + p_{bL}) \right\} + P \left\{ B > \frac{m}{2}(p_{aL} + p_{bL}) \right\}. \tag{B.2}$$

Taking $c = \frac{1}{2}(p_{aL} - p_{bL})$ yields

$$\begin{aligned} m(p_{aL} - c) &= \frac{m}{2}(p_{aL} + p_{bL}), \\ m(p_{bL} + c) &= \frac{m}{2}(p_{aL} + p_{bL}), \end{aligned}$$

and (B.2) is equivalent to

$$\begin{aligned} &\frac{\pi_a}{\pi_b} \mathbb{P}\{A \leq m(p_{aL} - c)\} + \mathbb{P}\{B > m(p_{bL} + c)\} \\ &= \frac{\pi_a}{\pi_b} \mathbb{P}\{A/m - p_{aL} \leq -c\} + \mathbb{P}\{B/m - p_{bL} > c\} \\ &\leq \frac{\pi_a}{\pi_b} e^{-2mc^2} + e^{-2mc^2}, \quad \text{by Lemma B.1} \\ &= \left(\frac{\pi_a + \pi_b}{\pi_b}\right) e^{-\frac{1}{2}m\delta_{ab}^2}. \end{aligned}$$

Now we have

$$\sum_{\mathbf{x} \in \Omega} \sum_{a \neq b}^s \frac{\pi_a P_a(\mathbf{x}) P_b(\mathbf{x})}{P(\mathbf{x})} = \sum_{a \neq b}^s \sum_{\mathbf{x} \in \Omega} \frac{\pi_a P_a(\mathbf{x}) P_b(\mathbf{x})}{P(\mathbf{x})} \leq \sum_{a \neq b}^s \frac{\pi_a + \pi_b}{\pi_b} e^{-\frac{m}{2}\delta_{ab}^2} \leq \frac{2}{\pi_b} \sum_{a \neq b}^s e^{-\frac{m}{2}\delta_{ab}^2}.$$

□

Corollary B.3. *The following intermediate result was obtained in the proof of Theorem B.2*

$$\sum_{\mathbf{x} \in \Omega} \frac{\pi_a P_a(\mathbf{x}) P_b(\mathbf{x})}{P(\mathbf{x})} \leq \left(\frac{\pi_a + \pi_b}{\pi_b}\right) e^{-\frac{1}{2}m\delta_{ab}^2} \leq \frac{2}{\pi_b} e^{-\frac{1}{2}m\delta_{ab}^2}.$$

We are now prepared to prove Theorem 2.1. Following the strategy of Morel and Nagaraj (1991), we consider the difference between the $\mathcal{I}(\boldsymbol{\theta})$ and the limiting matrix $\tilde{\mathcal{I}}(\boldsymbol{\theta})$ element by element for finite cluster sizes and obtain bounds which converge to zero as $m \rightarrow \infty$. The bound used by Morel and Nagaraj (1991) is slightly different than ours, since we do not require that $k = s$.

Proof of Theorem 2.1. Partition the exact FIM as

$$\mathcal{I}(\boldsymbol{\theta}) = \begin{pmatrix} \mathbf{C}_{11} & \mathbf{C}_{12} \\ \mathbf{C}_{21} & \mathbf{C}_{22} \end{pmatrix}$$

where

$$\mathbf{C}_{11} = \begin{pmatrix} \mathbf{A}_{11} & \cdots & \mathbf{A}_{1s} \\ \vdots & \ddots & \vdots \\ \mathbf{A}_{s1} & \cdots & \mathbf{A}_{ss} \end{pmatrix}, \quad \mathbf{C}_{12} = \begin{pmatrix} \mathbf{A}_{1\pi} \\ \vdots \\ \mathbf{A}_{s\pi} \end{pmatrix} = \mathbf{C}_{21}^T, \quad \mathbf{C}_{22} = \mathbf{A}_{\pi\pi},$$

and

$$\begin{aligned}\mathbf{A}_{ab} &= \mathbb{E} \left(\left\{ \frac{\partial \log f(\mathbf{x}; \boldsymbol{\theta})}{\partial \mathbf{p}_a} \right\} \left\{ \frac{\partial \log f(\mathbf{x}; \boldsymbol{\theta})}{\partial \mathbf{p}_b} \right\}^T \right), \quad \text{for } a = 1, \dots, s \text{ and } b = 1, \dots, s, \\ \mathbf{A}_{\pi b} &= \mathbb{E} \left(\left\{ \frac{\partial \log f(\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\pi}} \right\} \left\{ \frac{\partial \log f(\mathbf{x}; \boldsymbol{\theta})}{\partial \mathbf{p}_b} \right\}^T \right), \quad \text{for } b = 1, \dots, s \\ &= \mathbf{A}_{b\pi}^T, \\ \mathbf{A}_{\pi\pi} &= \mathbb{E} \left(\left\{ \frac{\partial \log f(\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\pi}} \right\} \left\{ \frac{\partial \log f(\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\pi}} \right\}^T \right).\end{aligned}$$

We must show that as $m \rightarrow \infty$,

$$\mathbf{C}_{11} - \text{Blockdiag}(\pi_1 \mathbf{F}_1, \dots, \pi_s \mathbf{F}_s) \rightarrow \mathbf{0}, \quad (\text{B.3})$$

$$\mathbf{C}_{21}^T = \mathbf{C}_{12} \rightarrow \mathbf{0}, \quad (\text{B.4})$$

$$\mathbf{C}_{22} - \mathbf{F}_\pi \rightarrow \mathbf{0}. \quad (\text{B.5})$$

We will show (B.3) and (B.5) only; (B.4) is similar. The reader may refer to Morel and Nagaraj (1991) for more details, shown for the $k = s$ case.

Case (i). First consider the (i, i) th block of $\mathbf{C}_{11} - \text{Blockdiag}(\pi_1 \mathbf{F}_1, \dots, \pi_s \mathbf{F}_s)$

$$\begin{aligned}\mathbf{D}_i (\mathbf{A}_{ii} - \pi_i \mathbf{F}_i) \mathbf{D}_i &= \mathbf{D}_i \left\{ \mathbb{E} \left[\left\{ \frac{\partial}{\partial \mathbf{p}_i} \log P(\mathbf{x}) \right\} \left\{ \frac{\partial}{\partial \mathbf{p}_i} \log P(\mathbf{x}) \right\}^T \right] - \pi_i \mathbf{F}_i \right\} \mathbf{D}_i \\ &= \pi_i^2 \mathbf{D}_i \mathbb{E} \left[\frac{P_i^2(\mathbf{x})}{P^2(\mathbf{x})} \frac{\partial \log P_i(\mathbf{x})}{\partial \mathbf{p}_i} \frac{\partial \log P_i(\mathbf{x})}{\partial \mathbf{p}_i^T} \right] \mathbf{D}_i - \pi_i \mathbf{D}_i \mathbf{F}_i \mathbf{D}_i \\ &= \pi_i^2 \sum_{\mathbf{x} \in \Omega} \frac{P_i(\mathbf{x})}{P(\mathbf{x})} \left(\mathbf{x}_{-k} - \frac{x_k}{p_{ik}} \mathbf{p}_i \right) \left(\mathbf{x}_{-k} - \frac{x_k}{p_{ik}} \mathbf{p}_i \right)^T P_i(\mathbf{x}) \\ &\quad - \pi_i^2 \sum_{\mathbf{x} \in \Omega} \frac{1}{\pi_i} \left(\mathbf{x}_{-k} - \frac{x_k}{p_{ik}} \mathbf{p}_i \right) \left(\mathbf{x}_{-k} - \frac{x_k}{p_{ik}} \mathbf{p}_i \right)^T P_i(\mathbf{x})\end{aligned} \quad (\text{B.6})$$

$$\begin{aligned}&= \pi_i^2 \sum_{\mathbf{x} \in \Omega} \left(\mathbf{x}_{-k} - \frac{x_k}{p_{ik}} \mathbf{p}_i \right) \left(\mathbf{x}_{-k} - \frac{x_k}{p_{ik}} \mathbf{p}_i \right)^T \left(\frac{P_i(\mathbf{x})}{P(\mathbf{x})} - \frac{1}{\pi_i} \right) P_i(\mathbf{x}) \\ &= \frac{\pi_i}{p_{ik}^2} \sum_{\mathbf{x} \in \Omega} (p_{ik} \mathbf{x}_{-k} - x_k \mathbf{p}_i) (p_{ik} \mathbf{x}_{-k} - x_k \mathbf{p}_i)^T \left(\frac{\pi_i P_i(\mathbf{x}) - P(\mathbf{x})}{P(\mathbf{x})} \right) P_i(\mathbf{x}).\end{aligned} \quad (\text{B.7})$$

where x_k is the k th element of \mathbf{x} and $\mathbf{x}_{-k} = (x_1, \dots, x_{k-1})$. We have pre and post-multiplied by \mathbf{D}_i so that Theorem B.2 can be applied. But note that since \mathbf{D}_i does not vary over m ,

$$\mathbf{D}_i \{ \mathbf{A}_{ii} - \pi_i \mathbf{F}_i \} \mathbf{D}_i \rightarrow \mathbf{0} \quad \implies \quad \mathbf{A}_{ii} - \pi_i \mathbf{F}_i \rightarrow \mathbf{0}, \quad \text{as } m \rightarrow \infty.$$

We have also used the fact in step (B.6) that

$$\mathbf{D}_i \frac{\partial \log P_i(\mathbf{x})}{\partial \mathbf{p}_i} = \mathbf{D}_i \left\{ \mathbf{D}_i^{-1} \mathbf{x}_{-k} - \frac{x_k}{p_{ik}} \mathbf{1} \right\} = \mathbf{x}_{-k} - \frac{x_k}{p_{ik}} \mathbf{p}_i.$$

We next have for $r, s \in \{1, \dots, k-1\}$

$$[p_{ik}x_r - x_k p_{ir}]^2 \leq [x_r + m p_{ir}]^2 \leq 4m^2.$$

Also,

$$\begin{aligned} 0 &\leq \left[[p_{ik}x_r - x_k p_{ir}] + [p_{ik}x_s - x_k p_{is}] \right]^2 \\ &= [p_{ik}x_r - x_k p_{ir}]^2 + [p_{ik}x_s - x_k p_{is}]^2 + 2[p_{ik}x_r - x_k p_{ir}][p_{ik}x_s - x_k p_{is}] \end{aligned}$$

and similarly

$$\begin{aligned} 0 &\leq \left[[p_{ik}x_r - x_k p_{ir}] - [p_{ik}x_s - x_k p_{is}] \right]^2 \\ &= [p_{ik}x_r - x_k p_{ir}]^2 + [p_{ik}x_s - x_k p_{is}]^2 - 2[p_{ik}x_r - x_k p_{ir}][p_{ik}x_s - x_k p_{is}], \end{aligned}$$

which implies that

$$\begin{aligned} \left| [p_{ik}x_r - x_k p_{ir}][p_{ik}x_s - x_k p_{is}] \right| &\leq \frac{1}{2} \left\{ [x_r + m p_{ir}]^2 + [x_s + m p_{is}]^2 \right\} \\ &\leq 4m^2. \end{aligned}$$

Notice that this bound is free of r and s , so it holds uniformly over all $r, s \in \{1, \dots, k-1\}$. If we denote the (r, s) th element of the matrix given in (B.7) by ε_{rs} , we have

$$\begin{aligned} |\varepsilon_{rs}| &\leq \frac{4\pi_i m^2}{p_{ik}^2} \sum_{\mathbf{x} \in \Omega} \frac{P(\mathbf{x}) - \pi_i P_i(\mathbf{x})}{P(\mathbf{x})} P_i(\mathbf{x}) = \frac{4\pi_i m^2}{p_{ik}^2} \sum_{\mathbf{x} \in \Omega} \sum_{j \neq i} \frac{\pi_j P_i(\mathbf{x}) P_j(\mathbf{x})}{P(\mathbf{x})} \\ &\leq \frac{8m^2}{p_{ik}^2} \sum_{j \neq i} e^{-\frac{m}{2} \delta_{ij}^2}, \end{aligned}$$

by Theorem B.2. By assumption, $\delta_{ij}^2 > 0$ for $i \neq j$, and therefore $\varepsilon_{rs} \rightarrow 0$ as $m \rightarrow \infty$.

Case (ii). Next, consider the (i, j) th block of $\mathbf{C}_{11} - \text{Blockdiag}(\pi_1 \mathbf{F}_1, \dots, \pi_s \mathbf{F}_s)$ where $i \neq j$.

$$\begin{aligned}
& \mathbf{D}_i \mathbf{A}_{ij} \mathbf{D}_j \\
&= \mathbf{D}_i \left\{ \mathbb{E} \left[\left\{ \frac{\partial}{\partial \mathbf{p}_i} \log P(\mathbf{x}) \right\} \left\{ \frac{\partial}{\partial \mathbf{p}_j} \log P(\mathbf{x}) \right\}^T \right] \right\} \mathbf{D}_j \\
&= \mathbf{D}_i \left[\mathbb{E} \left(\frac{\pi_i \pi_j}{P^2(\mathbf{x})} \frac{\partial P_i(\mathbf{x})}{\partial \mathbf{p}_i} \frac{\partial P_j(\mathbf{x})}{\partial \mathbf{p}_j^T} \right) \right] \mathbf{D}_j \\
&= \pi_i \pi_j \mathbf{D}_i \left[\mathbb{E} \left(\frac{P_i(\mathbf{x}) P_j(\mathbf{x})}{P^2(\mathbf{x})} \frac{\partial \log P_i(\mathbf{x})}{\partial \mathbf{p}_i} \frac{\partial \log P_j(\mathbf{x})}{\partial \mathbf{p}_j^T} \right) \right] \mathbf{D}_j \\
&= \pi_i \pi_j \sum_{\mathbf{x} \in \Omega} \frac{P_i(\mathbf{x}) P_j(\mathbf{x})}{P^2(\mathbf{x})} \left(\mathbf{x}_{-k} - \frac{x_k}{p_{ik}} \mathbf{p}_i \right) \left(\mathbf{x}_{-k} - \frac{x_k}{p_{jk}} \mathbf{p}_j \right)^T P(\mathbf{x}) \\
&= \frac{\pi_i \pi_j}{p_{ik} p_{jk}} \sum_{\mathbf{x} \in \Omega} \frac{P_i(\mathbf{x}) P_j(\mathbf{x})}{P(\mathbf{x})} (p_{ik} \mathbf{x}_{-k} - x_k \mathbf{p}_i) (p_{jk} \mathbf{x}_{-k} - x_k \mathbf{p}_j)^T. \tag{B.8}
\end{aligned}$$

If we now denote the (r, s) th element of the matrix given in (B.8) by ε_{rs} , we have

$$|\varepsilon_{rs}| \leq \frac{4\pi_i \pi_j m^2}{p_{ik} p_{jk}} \sum_{\mathbf{x} \in \Omega} \frac{P_i(\mathbf{x}) P_j(\mathbf{x})}{P(\mathbf{x})} \leq \frac{8m^2}{p_{ik} p_{jk}} e^{-\frac{m}{2} \delta_{ij}^2}$$

for all (r, s) , applying Theorem B.3 and a similar argument to Case (i). Since $\delta_{ij}^2 > 0$ for $i \neq j$, $\varepsilon_{rs} \rightarrow 0$ as $m \rightarrow \infty$.

Case (iii). Now consider the matrix

$$\begin{aligned}
& \mathbf{A}_{\pi\pi} - \mathbf{F}_\pi \tag{B.9} \\
&= \mathbb{E} \left[\left\{ \frac{\partial}{\partial \boldsymbol{\pi}} \log P(\mathbf{x}) \right\} \left\{ \frac{\partial}{\partial \boldsymbol{\pi}} \log P(\mathbf{x}) \right\}^T \right] - \mathbf{F}_\pi \\
&= \mathbb{E} \left[\frac{1}{P^2(\mathbf{x})} \left\{ \begin{pmatrix} P_1(\mathbf{x}) \\ \vdots \\ P_{s-1}(\mathbf{x}) \end{pmatrix} - P_s(\mathbf{x}) \cdot \mathbf{1} \right\} \left\{ \begin{pmatrix} P_1(\mathbf{x}) \\ \vdots \\ P_{s-1}(\mathbf{x}) \end{pmatrix} - P_s(\mathbf{x}) \cdot \mathbf{1} \right\}^T \right] - (\mathbf{D}_\pi^{-1} + \pi_s^{-1} \mathbf{1}\mathbf{1}^T).
\end{aligned}$$

Pick out the (a, a) th entry which we will denote as ε_{aa} . We have

$$\begin{aligned}
\varepsilon_{aa} &= \mathbb{E} \left[\frac{[\mathbf{P}_a(\mathbf{x}) - \mathbf{P}_s(\mathbf{x})]^2}{\mathbf{P}^2(\mathbf{x})} \right] - (\pi_a^{-1} + \pi_s^{-1}) \\
&= \sum_{\mathbf{x} \in \Omega} \frac{\mathbf{P}_a^2(\mathbf{x}) - 2\mathbf{P}_a(\mathbf{x})\mathbf{P}_s(\mathbf{x}) + \mathbf{P}_s^2(\mathbf{x})}{\mathbf{P}(\mathbf{x})} - (\pi_a^{-1} + \pi_s^{-1}) \\
&= \sum_{\mathbf{x} \in \Omega} \left(\frac{\mathbf{P}_a^2(\mathbf{x})}{\mathbf{P}(\mathbf{x})} - \frac{\mathbf{P}_a(\mathbf{x})}{\pi_a} \right) + \sum_{\mathbf{x} \in \Omega} \left(\frac{\mathbf{P}_s^2(\mathbf{x})}{\mathbf{P}(\mathbf{x})} - \frac{\mathbf{P}_s(\mathbf{x})}{\pi_s} \right) - 2 \sum_{\mathbf{x} \in \Omega} \frac{\mathbf{P}_a(\mathbf{x})\mathbf{P}_s(\mathbf{x})}{\mathbf{P}(\mathbf{x})} \\
&= \frac{1}{\pi_a} \sum_{\mathbf{x} \in \Omega} \frac{\pi_a \mathbf{P}_a(\mathbf{x}) - \mathbf{P}(\mathbf{x})}{\mathbf{P}(\mathbf{x})} \mathbf{P}_a(\mathbf{x}) + \frac{1}{\pi_s} \sum_{\mathbf{x} \in \Omega} \frac{\pi_s \mathbf{P}_s(\mathbf{x}) - \mathbf{P}(\mathbf{x})}{\mathbf{P}(\mathbf{x})} \mathbf{P}_s(\mathbf{x}) - 2 \sum_{\mathbf{x} \in \Omega} \frac{\mathbf{P}_a(\mathbf{x})\mathbf{P}_s(\mathbf{x})}{\mathbf{P}(\mathbf{x})} \\
&= -\frac{1}{\pi_a} \sum_{\mathbf{x} \in \Omega} \sum_{\ell \neq a}^s \frac{\pi_\ell \mathbf{P}_\ell(\mathbf{x}) \mathbf{P}_a(\mathbf{x})}{\mathbf{P}(\mathbf{x})} - \frac{1}{\pi_s} \sum_{\mathbf{x} \in \Omega} \sum_{\ell \neq s}^s \frac{\pi_\ell \mathbf{P}_\ell(\mathbf{x}) \mathbf{P}_s(\mathbf{x})}{\mathbf{P}(\mathbf{x})} - \frac{2}{\pi_a} \sum_{\mathbf{x} \in \Omega} \frac{\pi_a \mathbf{P}_a(\mathbf{x}) \mathbf{P}_s(\mathbf{x})}{\mathbf{P}(\mathbf{x})}
\end{aligned}$$

Then by the triangle inequality,

$$|\varepsilon_{aa}| \leq \frac{2}{\pi_a^2} \sum_{\ell \neq a}^s e^{-\frac{m}{2}\delta_{\ell a}^2} + \frac{2}{\pi_s^2} \sum_{\ell \neq s}^s e^{-\frac{m}{2}\delta_{\ell s}^2} + \frac{4}{\pi_a \pi_s} e^{-\frac{m}{2}\delta_{as}^2},$$

applying Theorem B.2 to the first two terms, and Corollary B.3 to the last term. Since $\delta_{ij}^2 > 0$ for $i \neq j$, we have $\varepsilon_{aa} \rightarrow 0$ for $a \in \{1, \dots, s-1\}$ as $m \rightarrow \infty$.

Case (iv). Consider again the matrix $\mathbf{A}_{\pi\pi} - \mathbf{F}_\pi$ from (B.9), but now the case where $a \neq b$. We have

$$\begin{aligned}
\varepsilon_{ab} &= \mathbb{E} \left[\frac{[\mathbf{P}_a(\mathbf{x}) - \mathbf{P}_s(\mathbf{x})][\mathbf{P}_b(\mathbf{x}) - \mathbf{P}_s(\mathbf{x})]}{\mathbf{P}^2(\mathbf{x})} - \pi_s^{-1} \right] \\
&= \sum_{\mathbf{x} \in \Omega} \frac{\mathbf{P}_a(\mathbf{x})\mathbf{P}_b(\mathbf{x})}{\mathbf{P}(\mathbf{x})} - \sum_{\mathbf{x} \in \Omega} \frac{\mathbf{P}_a(\mathbf{x})\mathbf{P}_s(\mathbf{x})}{\mathbf{P}(\mathbf{x})} - \sum_{\mathbf{x} \in \Omega} \frac{\mathbf{P}_b(\mathbf{x})\mathbf{P}_s(\mathbf{x})}{\mathbf{P}(\mathbf{x})} + \sum_{\mathbf{x} \in \Omega} \frac{\mathbf{P}_s^2(\mathbf{x})}{\mathbf{P}(\mathbf{x})} - \pi_s^{-1}. \quad (\text{B.10})
\end{aligned}$$

We can use Corollary B.3 to handle the first three terms. For the last term, notice that

$$\sum_{\mathbf{x} \in \Omega} \frac{\mathbf{P}_s^2(\mathbf{x})}{\mathbf{P}(\mathbf{x})} - \frac{1}{\pi_s} = \sum_{\mathbf{x} \in \Omega} \left(\frac{\mathbf{P}_s(\mathbf{x})}{\mathbf{P}(\mathbf{x})} - \frac{1}{\pi_s} \right) \mathbf{P}_s(\mathbf{x}) = -\frac{1}{\pi_s} \sum_{\mathbf{x} \in \Omega} \sum_{\ell \neq s} \frac{\pi_\ell \mathbf{P}_\ell(\mathbf{x}) \mathbf{P}_s(\mathbf{x})}{\mathbf{P}(\mathbf{x})}.$$

Now, applying the triangle inequality to (B.10),

$$|\varepsilon_{ab}| \leq \frac{2}{\pi_a \pi_b} e^{-\frac{m}{2}\delta_{ab}^2} + \frac{2}{\pi_a \pi_s} e^{-\frac{m}{2}\delta_{as}^2} + \frac{2}{\pi_b \pi_s} e^{-\frac{m}{2}\delta_{bs}^2} + \frac{2}{\pi_s^2} \sum_{\ell \neq s} e^{-\frac{m}{2}\delta_{\ell s}^2}$$

Since $\delta_{ij}^2 > 0$ for $i \neq j$, we have $\varepsilon_{ab} \rightarrow 0$ for $a \neq b$ in $\{1, \dots, s-1\}$ as $m \rightarrow \infty$. \square

Proof of Theorem 2.5. This proof uses properties of matrix norms; refer to Lange (2010, Chapter 6) or Meyer (2001, Chapter 5) for background. Notice that for non-singular $q \times q$ matrices \mathbf{A} and \mathbf{B} ,

$$\mathbf{B}^{-1} - \mathbf{A}^{-1} = \mathbf{A}^{-1}(\mathbf{A} - \mathbf{B})\mathbf{B}^{-1}.$$

Then for any matrix norm satisfying the sub-multiplicative property,

$$\|\mathbf{A}^{-1} - \mathbf{B}^{-1}\| \leq \|\mathbf{A}^{-1}\| \cdot \|\mathbf{A} - \mathbf{B}\| \cdot \|\mathbf{B}^{-1}\|. \quad (\text{B.11})$$

Fix $\boldsymbol{\theta} \in \Theta$, take $\mathbf{A} = \tilde{\mathcal{I}}(\boldsymbol{\theta})$ and $\mathbf{B} = \mathcal{I}(\boldsymbol{\theta})$, and take $\|\cdot\|$ to the Frobenius matrix norm for convenience. Then (B.11) becomes

$$\|\mathcal{I}^{-1}(\boldsymbol{\theta}) - \tilde{\mathcal{I}}^{-1}(\boldsymbol{\theta})\|_F \leq \|\mathcal{I}^{-1}(\boldsymbol{\theta})\|_F \cdot \|\tilde{\mathcal{I}}^{-1}(\boldsymbol{\theta})\|_F \cdot \|\mathcal{I}(\boldsymbol{\theta}) - \tilde{\mathcal{I}}(\boldsymbol{\theta})\|_F,$$

where $\|\mathbf{A}\|_F = \sqrt{\sum_{i=1}^q \sum_{j=1}^q a_{ij}^2}$, and a_{ij} denote the elements of \mathbf{A} . To show that the RHS converges to 0 as $m \rightarrow \infty$, we will handle the three terms separately. Since $\mathcal{I}(\boldsymbol{\theta}) - \tilde{\mathcal{I}}(\boldsymbol{\theta}) \rightarrow \mathbf{0}$ as $m \rightarrow \infty$ by Theorem 2.1, $\|\mathcal{I}(\boldsymbol{\theta}) - \tilde{\mathcal{I}}(\boldsymbol{\theta})\|_F \rightarrow 0$. Next, we address the $\|\tilde{\mathcal{I}}^{-1}(\boldsymbol{\theta})\|_F$ term. Using the explicit form in Corollary 2.4, we have

$$\begin{aligned} 0 \leq \|\tilde{\mathcal{I}}^{-1}(\boldsymbol{\theta})\|_F^2 &= \sum_{\ell=1}^s \|\pi_\ell^{-1} \mathbf{F}_\ell^{-1}\|_F^2 + \|\mathbf{F}_\pi^{-1}\|_F^2 \\ &= \sum_{\ell=1}^s m^{-2} \pi_\ell^{-2} \|\mathbf{D}_\ell - \mathbf{p}_\ell \mathbf{p}_\ell^T\|_F^2 + \|\mathbf{D}_\pi - \boldsymbol{\pi} \boldsymbol{\pi}^T\|_F^2. \end{aligned}$$

Then $\|\tilde{\mathcal{I}}^{-1}(\boldsymbol{\theta})\|_F$ is decreasing in m , and hence is bounded in m .

We will now consider the term $\|\mathcal{I}^{-1}(\boldsymbol{\theta})\|$, with the 2-norm in place of the Frobenius norm. Let $\lambda_1(m) \geq \dots \geq \lambda_q(m)$ be the eigenvalues of $\mathcal{I}(\boldsymbol{\theta})$ for a fixed m , all assumed to be positive. Since the 2-norm of a symmetric positive definite matrix is its largest eigenvalue, we have

$$\begin{aligned} 0 \leq \|\mathcal{I}^{-1}(\boldsymbol{\theta})\|_2 &= \frac{1}{\lambda_q(m)} = \frac{1}{\min_{\|\mathbf{x}\|=1} \mathbf{x}^T \mathcal{I}(\boldsymbol{\theta}) \mathbf{x}} \\ &= \frac{1}{\min_{\|\mathbf{x}\|=1} \left\{ \mathbf{x}^T \left[\mathcal{I}(\boldsymbol{\theta}) - \tilde{\mathcal{I}}(\boldsymbol{\theta}) \right] \mathbf{x} + \mathbf{x}^T \tilde{\mathcal{I}}(\boldsymbol{\theta}) \mathbf{x} \right\}}. \end{aligned}$$

Notice that

$$\min_{\|\mathbf{x}\|=1} \mathbf{x}^T \left[\mathcal{I}(\boldsymbol{\theta}) - \tilde{\mathcal{I}}(\boldsymbol{\theta}) \right] \mathbf{x} + \min_{\|\mathbf{x}\|=1} \mathbf{x}^T \tilde{\mathcal{I}}(\boldsymbol{\theta}) \mathbf{x} \leq \min_{\|\mathbf{x}\|=1} \left\{ \mathbf{x}^T \left[\mathcal{I}(\boldsymbol{\theta}) - \tilde{\mathcal{I}}(\boldsymbol{\theta}) \right] \mathbf{x} + \mathbf{x}^T \tilde{\mathcal{I}}(\boldsymbol{\theta}) \mathbf{x} \right\}$$

since both LHS and RHS are lower bounds for $\mathbf{x}^T \left[\mathcal{I}(\boldsymbol{\theta}) - \tilde{\mathcal{I}}(\boldsymbol{\theta}) \right] \mathbf{x} + \mathbf{x}^T \tilde{\mathcal{I}}(\boldsymbol{\theta}) \mathbf{x}$, and the RHS is the greatest such bound. Therefore

$$\begin{aligned} 1/\lambda_q(m) &\leq \frac{1}{\min_{\|\mathbf{x}\|=1} \mathbf{x}^T \left[\mathcal{I}(\boldsymbol{\theta}) - \tilde{\mathcal{I}}(\boldsymbol{\theta}) \right] \mathbf{x} + \min_{\|\mathbf{x}\|=1} \mathbf{x}^T \tilde{\mathcal{I}}(\boldsymbol{\theta}) \mathbf{x}} \\ &= \frac{1}{\beta_q(m) + \tilde{\lambda}_q(m)} \end{aligned}$$

denoting the eigenvalues of $\tilde{\mathcal{I}}(\boldsymbol{\theta})$ as $\tilde{\lambda}_1(m) \geq \dots \geq \tilde{\lambda}_q(m)$ (all positive), and the eigenvalues of $\mathcal{I}(\boldsymbol{\theta}) - \tilde{\mathcal{I}}(\boldsymbol{\theta})$ as $\beta_1(m) \geq \dots \geq \beta_q(m)$. It is well known that the mapping from a matrix to its eigenvalues is a continuous function of its elements (Meyer, 2001, Chapter 7). Therefore

$$\mathcal{I}(\boldsymbol{\theta}) - \tilde{\mathcal{I}}(\boldsymbol{\theta}) \rightarrow \mathbf{0} \text{ as } m \rightarrow \infty \quad \implies \quad \beta_q(m) \rightarrow 0 \text{ as } m \rightarrow \infty.$$

Now for any $\varepsilon > 0$, there exists a positive integer m_0 such that $|\beta_q(m)| < \varepsilon$ for all $m \geq m_0$, and so we have

$$0 \leq \|\mathcal{I}^{-1}(\boldsymbol{\theta})\|_2 \leq \frac{1}{\beta_q(m) + \tilde{\lambda}_q(m)} \leq \frac{1}{\tilde{\lambda}_q(m) - \varepsilon} \quad (\text{B.12})$$

for all $m \geq m_0$. Because $\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_F$, the result in Part (ii) gives that for all m there exists a $K > 0$ such that,

$$1/\tilde{\lambda}_q(m) = \|\tilde{\mathcal{I}}^{-1}(\boldsymbol{\theta})\|_2 \leq \|\tilde{\mathcal{I}}^{-1}(\boldsymbol{\theta})\|_F \leq K \quad \iff \quad \tilde{\lambda}_q(m) \geq 1/K.$$

WLOG assume that ε has been chosen so that $\tilde{\lambda}_q(m) \geq 1/K > \varepsilon$, to avoid division by zero. The RHS of (B.12) is bounded above by $(1/K - \varepsilon)^{-1}$ for all $m \geq m_0$, which implies $\|\mathcal{I}^{-1}(\boldsymbol{\theta})\|_2$ is bounded when $m \geq m_0$.

To conclude the proof, note that in general $q^{-1/2}\|\mathbf{A}\|_F \leq \|\mathbf{A}\|_2$, so that

$$\begin{aligned} 0 &\leq \|\mathcal{I}^{-1}(\boldsymbol{\theta}) - \tilde{\mathcal{I}}^{-1}(\boldsymbol{\theta})\|_F \\ &\leq \|\mathcal{I}^{-1}(\boldsymbol{\theta})\|_F \cdot \|\tilde{\mathcal{I}}^{-1}(\boldsymbol{\theta})\|_F \cdot \|\mathcal{I}(\boldsymbol{\theta}) - \tilde{\mathcal{I}}(\boldsymbol{\theta})\|_F \\ &\leq \sqrt{q}\|\mathcal{I}^{-1}(\boldsymbol{\theta})\|_2 \cdot \|\tilde{\mathcal{I}}^{-1}(\boldsymbol{\theta})\|_F \cdot \|\mathcal{I}(\boldsymbol{\theta}) - \tilde{\mathcal{I}}(\boldsymbol{\theta})\|_F. \end{aligned}$$

It follows from the earlier steps that the RHS converges to zero as $m \rightarrow \infty$, and therefore $\|\mathcal{I}^{-1}(\boldsymbol{\theta}) - \tilde{\mathcal{I}}^{-1}(\boldsymbol{\theta})\|_F \rightarrow 0$, which implies $\mathcal{I}^{-1}(\boldsymbol{\theta}) - \tilde{\mathcal{I}}^{-1}(\boldsymbol{\theta}) \rightarrow \mathbf{0}$. □

Proof of Proposition 3.3. The general form for AFSA is given by

$$\boldsymbol{\theta}^{(g+1)} = \boldsymbol{\theta}^{(g)} + \text{Blockdiag}(\pi_1 \mathbf{F}_1^{-1}, \dots, \pi_s \mathbf{F}_s^{-1}, \mathbf{F}_\pi^{-1}) S(\boldsymbol{\theta}^{(g)})$$

so that the individual updates are

$$\begin{aligned} \mathbf{p}_\ell^{(g+1)} &= \mathbf{p}_\ell^{(g)} + \pi_\ell^{-1} \mathbf{F}_\ell^{-1} \frac{\partial}{\partial \mathbf{p}_\ell} \log L(\boldsymbol{\theta}^{(g)}), \quad \ell = 1, \dots, s \\ \boldsymbol{\pi}^{(g+1)} &= \boldsymbol{\pi}^{(g)} + \mathbf{F}_\pi^{-1} \frac{\partial}{\partial \boldsymbol{\pi}} \log L(\boldsymbol{\theta}^{(g)}). \end{aligned}$$

From Corollary 2.4 we have

$$\begin{aligned} \boldsymbol{\pi}^{(g+1)} &= \boldsymbol{\pi}^{(g)} + (n \mathbf{F}_\pi)^{-1} \sum_{i=1}^n \frac{\partial \log L(\boldsymbol{\theta}^{(g)} | \mathbf{x}_i)}{\partial \boldsymbol{\pi}} \\ &= \boldsymbol{\pi}^{(g)} + n^{-1} \left[\text{diag}\{\boldsymbol{\pi}^{(g)}\} - \boldsymbol{\pi}^{(g)} \boldsymbol{\pi}^{(g)T} \right] \sum_{i=1}^n \frac{\partial \log(\boldsymbol{\theta}^{(g)} | \mathbf{x}_i)}{\partial \boldsymbol{\pi}}. \end{aligned}$$

Then for $\ell = 1, \dots, s-1$,

$$\begin{aligned}
\pi_\ell^{(g+1)} &= \pi_\ell^{(g)} + n^{-1} \pi_\ell^{(g)} \sum_{i=1}^n \frac{P_\ell(\mathbf{x}_i) - P_s(\mathbf{x}_i)}{P(\mathbf{x}_i)} - n^{-1} \sum_{i=1}^n \sum_{t=1}^{s-1} \pi_\ell^{(g)} \pi_t^{(g)} \frac{P_t(\mathbf{x}_i) - P_s(\mathbf{x}_i)}{P(\mathbf{x}_i)} \\
&= \pi_\ell^{(g)} + n^{-1} \pi_\ell^{(g)} \sum_{i=1}^n \frac{P_\ell(\mathbf{x}_i) - P_s(\mathbf{x}_i)}{P(\mathbf{x}_i)} - n^{-1} \pi_\ell^{(g)} \sum_{i=1}^n \left\{ \frac{P(\mathbf{x}_i) - \pi_s^{(g)} P_s(\mathbf{x}_i) - (1 - \pi_s^{(g)}) P_s(\mathbf{x}_i)}{P(\mathbf{x}_i)} \right\} \\
&= \pi_\ell^{(g)} + n^{-1} \pi_\ell^{(g)} \sum_{i=1}^n \frac{P_\ell(\mathbf{x}_i) - P_s(\mathbf{x}_i)}{P(\mathbf{x}_i)} - n^{-1} \pi_\ell^{(g)} \sum_{i=1}^n \left\{ 1 - \frac{P_s(\mathbf{x}_i)}{P(\mathbf{x}_i)} \right\} \\
&= \pi_\ell^{(g)} \frac{1}{n} \sum_{i=1}^n \frac{P_\ell(\mathbf{x}_i)}{P(\mathbf{x}_i)}.
\end{aligned}$$

Next, to obtain explicit iterations for $p_{\ell j}$'s, the blocks for $\ell = 1, \dots, s$ are given by

$$\begin{aligned}
\mathbf{p}_\ell^{(g+1)} &= \mathbf{p}_\ell^{(g)} + \left(\pi_\ell^{(g)} \mathbf{F}_\ell \right)^{-1} \sum_{i=1}^n \frac{\partial}{\partial \mathbf{p}_\ell} \log L(\boldsymbol{\theta}^{(g)} \mid \mathbf{x}_i) \\
&= \mathbf{p}_\ell^{(g)} + \frac{1}{M \pi_\ell^{(g)}} \left[\text{diag}\{\mathbf{p}_\ell^{(g)}\} - \mathbf{p}_\ell^{(g)} \mathbf{p}_\ell^{(g)T} \right] \sum_{i=1}^n \frac{\partial}{\partial \mathbf{p}_\ell} \log L(\boldsymbol{\theta}^{(g)} \mid \mathbf{x}_i).
\end{aligned}$$

For $j = 1, \dots, k-1$,

$$\begin{aligned}
p_{\ell j}^{(g+1)} &= p_{\ell j}^{(g)} + \frac{1}{M} \sum_{i=1}^n p_{\ell j}^{(g)} \frac{P_\ell(\mathbf{x}_i)}{P(\mathbf{x}_i)} \left(\frac{x_{ij}}{p_{\ell j}^{(g)}} - \frac{x_{ik}}{p_{\ell k}^{(g)}} \right) - \frac{1}{M} \sum_{i=1}^n \sum_{t=1}^{k-1} p_{\ell j}^{(g)} p_{\ell t}^{(g)} \frac{P_\ell(\mathbf{x}_i)}{P(\mathbf{x}_i)} \left(\frac{x_{it}}{p_{\ell t}^{(g)}} - \frac{x_{ik}}{p_{\ell k}^{(g)}} \right) \\
&= p_{\ell j}^{(g)} + \frac{1}{M} \sum_{i=1}^n \frac{P_\ell(\mathbf{x}_i)}{P(\mathbf{x}_i)} \left(x_{ij} - \frac{p_{\ell j}^{(g)}}{p_{\ell k}^{(g)}} x_{ik} \right) - \frac{1}{M} \sum_{i=1}^n p_{\ell j}^{(g)} \frac{P_\ell(\mathbf{x}_i)}{P(\mathbf{x}_i)} \sum_{t=1}^{k-1} \left(x_{it} - \frac{p_{\ell t}^{(g)}}{p_{\ell k}^{(g)}} x_{ik} \right). \\
&= p_{\ell j}^{(g)} + \frac{1}{M} \sum_{i=1}^n \frac{P_\ell(\mathbf{x}_i)}{P(\mathbf{x}_i)} \left\{ \left(x_{ij} - \frac{p_{\ell j}^{(g)}}{p_{\ell k}^{(g)}} x_{ik} \right) - p_{\ell j}^{(g)} \sum_{t=1}^{k-1} \left(x_{it} - \frac{p_{\ell t}^{(g)}}{p_{\ell k}^{(g)}} x_{ik} \right) \right\} \quad (\text{B.13})
\end{aligned}$$

Since $\sum_{t=1}^k x_{it} = m_i$ and $\sum_{t=1}^k p_{\ell t}^{(g)} = 1$,

$$\sum_{t=1}^{k-1} \left(x_{it} - \frac{p_{\ell t}^{(g)}}{p_{\ell k}^{(g)}} x_{ik} \right) = (m_i - x_{ik}) - x_{ik} \frac{1 - p_{\ell k}^{(g)}}{p_{\ell k}^{(g)}} = m_i - x_{ik} / p_{\ell k}^{(g)}.$$

Applying this result to (B.13) and simplifying we get

$$\begin{aligned}
p_{\ell j}^{(g+1)} &= p_{\ell j}^{(g)} + \frac{1}{M} \sum_{i=1}^n \frac{P_\ell(\mathbf{x}_i)}{P(\mathbf{x}_i)} \left(x_{ij} - m_i p_{\ell j}^{(g)} \right) \\
&= p_{\ell j}^{(g)} + \frac{1}{M} \sum_{i=1}^n \frac{P_\ell(\mathbf{x}_i)}{P(\mathbf{x}_i)} x_{ij} - \frac{p_{\ell j}^{(g)}}{M} \sum_{i=1}^n m_i \frac{P_\ell(\mathbf{x}_i)}{P(\mathbf{x}_i)}.
\end{aligned}$$

□

Proof of Proposition 3.4. The complete data likelihood is

$$L(\boldsymbol{\theta} \mid \mathbf{x}, \mathbf{z}) = \prod_{i=1}^n \prod_{\ell=1}^s \left[\pi_{\ell} f(\mathbf{x}_i \mid \mathbf{p}_{\ell}, m_i) \right]^{\Delta_{i\ell}}.$$

where $\Delta_{i\ell} = I(z_i = \ell) \sim \text{Bernoulli}(\pi_{\ell})$. Then the corresponding log-likelihood is

$$\log L(\boldsymbol{\theta} \mid \mathbf{x}, \mathbf{z}) = \sum_{i=1}^n \sum_{\ell=1}^s \Delta_{i\ell} \log \left[\pi_{\ell} f(\mathbf{x}_i \mid \mathbf{p}_{\ell}, m_i) \right]. \quad (\text{B.14})$$

Since z_1, \dots, z_n are not observed, we instead use the expected log-likelihood, conditional on $\boldsymbol{\theta} = \boldsymbol{\theta}^{(g)}$ and \mathbf{x} . First note that

$$\begin{aligned} \gamma_{i\ell}^{(g)} &:= \mathbb{E}(\Delta_{i\ell} \mid \mathbf{x}_1, \dots, \mathbf{x}_n, \boldsymbol{\theta}^{(g)}) = \mathbb{P}(Z_i = \ell \mid \mathbf{x}_i, \boldsymbol{\theta}^{(g)}) \\ &= \frac{\mathbb{P}(Z_i = \ell \mid \boldsymbol{\theta}^{(g)}) \mathbb{P}(\mathbf{x}_i \mid Z_i = \ell, \boldsymbol{\theta}^{(g)})}{f(\mathbf{x}_i \mid \boldsymbol{\theta}^{(g)}, m_i)} = \frac{\pi_{\ell}^{(g)} f(\mathbf{x}_i \mid \mathbf{p}_{\ell}^{(g)}, m_i)}{\sum_{a=1}^s \pi_a^{(g)} f(\mathbf{x}_i \mid \mathbf{p}_a^{(g)}, m_i)} \end{aligned}$$

is the posterior probability of population ℓ , given \mathbf{x}_i and the previous iteration. Conditional on this information, the expectation of (B.14) becomes

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(g)}) := \sum_{i=1}^n \sum_{\ell=1}^s \gamma_{i\ell}^{(g)} \log \pi_{\ell} + \sum_{i=1}^n \sum_{\ell=1}^s \gamma_{i\ell}^{(g)} \log \left[f(\mathbf{x}_i \mid \mathbf{p}_{\ell}, m_i) \right].$$

Now to maximize this expression with respect to each parameter, equate partial derivatives to zero and solve for the parameter. For π_1, \dots, π_{s-1} we have

$$\begin{aligned} 0 &= \frac{\partial}{\partial \pi_a} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(g)}) = \sum_{i=1}^n \frac{\gamma_{ia}^{(g)}}{\pi_a} - \sum_{i=1}^n \frac{\gamma_{is}^{(g)}}{\pi_s} \\ &\iff \pi_s \sum_{i=1}^n \gamma_{ia}^{(g)} = \pi_a \sum_{i=1}^n \gamma_{is}^{(g)}. \end{aligned} \quad (\text{B.15})$$

Summing both sides of (B.15) over $a = 1, \dots, s$ we obtain

$$\begin{aligned} \pi_s \sum_{a=1}^s \sum_{i=1}^n \gamma_{ia}^{(g)} &= \sum_{i=1}^n \gamma_{is}^{(g)} \iff \pi_s n = \sum_{i=1}^n \gamma_{is}^{(g)} \\ &\iff \hat{\pi}_s^{(g+1)} = \frac{1}{n} \sum_{i=1}^n \gamma_{is}^{(g)} \end{aligned}$$

since the posterior probabilities $\gamma_{i1}^{(g)}, \dots, \gamma_{is}^{(g)}$ sum to 1. Replacing this back into (B.15) yields

$$\hat{\pi}_a^{(g+1)} = \frac{\hat{\pi}_s^{(g+1)} \sum_{i=1}^n \gamma_{ia}^{(g)}}{\sum_{i=1}^n \gamma_{is}^{(g)}} = \frac{1}{n} \sum_{i=1}^n \gamma_{ia}^{(g)}.$$

Similar steps yield the EM iterations for the p_{ab} 's. For p_{ab} where $a = 1, \dots, s$ and $b = 1, \dots, k - 1$,

$$\begin{aligned}
0 &= \frac{\partial}{\partial p_{ab}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(g)}) = \sum_{i=1}^n \gamma_{ia}^{(g)} \left(\frac{x_{ib}}{p_{ab}} - \frac{x_{ik}}{p_{ik}} \right) \\
&\iff p_{ak} \sum_{i=1}^n \gamma_{ia}^{(g)} x_{ib} = p_{ab} \sum_{i=1}^n \gamma_{ia}^{(g)} x_{ik}.
\end{aligned} \tag{B.16}$$

Summing both sides of (B.16) over $b = 1, \dots, k$ we obtain

$$p_{ak} \sum_{i=1}^n \gamma_{ia}^{(g)} m_i = \sum_{i=1}^n \gamma_{ia}^{(g)} x_{ik} \iff \hat{p}_{ak}^{(g+1)} = \frac{\sum_{i=1}^n x_{ik} \gamma_{ia}^{(g)}}{\sum_{i=1}^n m_i \gamma_{ia}^{(g)}}$$

since $x_{i1} + \dots + x_{ik} = m_i$. Replacing this back into (B.16) yields

$$\hat{p}_{ab}^{(g+1)} = \hat{p}_{ak}^{(g+1)} \frac{\sum_{i=1}^n x_{ib} \gamma_{ia}^{(g)}}{\sum_{i=1}^n x_{ik} \gamma_{ia}^{(g)}} = \frac{\sum_{i=1}^n x_{ib} \gamma_{ia}^{(g)}}{\sum_{i=1}^n m_i \gamma_{ia}^{(g)}}.$$

□