

General Disclaimer

One or more of the Following Statements may affect this Document

- This document has been reproduced from the best copy furnished by the organizational source. It is being released in the interest of making available as much information as possible.
- This document may contain data, which exceeds the sheet parameters. It was furnished in this condition by the organizational source and is the best copy available.
- This document may contain tone-on-tone or color graphs, charts and/or pictures, which have been reproduced in black and white.
- This document is paginated as submitted by the original source.
- Portions of this document are not fully legible due to the historical nature of some of the material. However, it is the best reproduction available from the original submission.



DEPARTMENT OF MATHEMATICS
UNIVERSITY OF HOUSTON

HOUSTON, TEXAS

NASA CR-

CR-147857

(NASA-CR-147857) THE NUMERICAL EVALUATION
OF MAXIMUM-LIKELIHOOD ESTIMATES OF THE
PARAMETERS FOR A MIXTURE OF NORMAL
DISTRIBUTIONS FROM PARTIALLY IDENTIFIED
SAMPLES (Houston Univ.) 17 p HC \$3.50

N76-30891

Unclas

G3/65 01742

THE NUMERICAL EVALUATION OF MAX
LIKELIHOOD ESTIMATES OF THE
PARAMETERS FOR A MIXTURE OF
NORMAL DISTRIBUTIONS FROM
PARTIALLY IDENTIFIED SAMPLES

HOMER F. WALKER
REPORT #54 JUNE, 1976



PREPARED FOR
EARTH OBSERVATION DIVISION, JSC
UNDER
CONTRACT NAS-9-15000

HOUSTON, TEXAS 77004

*The Numerical Evaluation of Maximum-Likelihood
Estimates of the Parameters for a Mixture of Normal Distributions
from Partially Identified Samples*

by

Homer F. Walker

Department of Mathematics, University of Houston

Houston, Texas 77004

June, 1976

Report #54

The Numerical Evaluation of Maximum-Likelihood
Estimates of the Parameters for a Mixture of Normal Distributions
from Partially Identified Samples

by

Homer F. Walker

Department of Mathematics, University of Houston
Houston, Texas 77004

1. Introduction.

Let π_1, \dots, π_m be populations whose multivariate observations in \mathbb{R}^n are distributed with respective normal density functions

$$p_i(x) = \frac{1}{(2\pi)^{n/2} |\Sigma_i^0|^{1/2}} e^{-\frac{1}{2}(x-\mu_i^0)^T \Sigma_i^{0-1} (x-\mu_i^0)}, \quad i = 1, \dots, m.$$

If π_0 is a given mixture of members of these populations, then observations on π_0 are distributed in \mathbb{R}^n with density function

$$p(x) = \sum_{i=1}^m \alpha_i^0 p_i(x)$$

for an appropriate set of proportions $\{\alpha_i^0\}_{i=1, \dots, m}$. These proportions necessarily satisfy $\sum_{i=1}^m \alpha_i^0 = 1$ and $\alpha_i^0 \geq 0$, $i = 1, \dots, m$. In this note, we also assume that each α_i^0 is strictly positive.

We address here the problem of numerically approximating the maximum-likelihood estimates of the parameters $\{\alpha_i^0, \mu_i^0, \Sigma_i^0\}_{i=1, \dots, m}$ determined by samples of two types. Samples of both types consist of sets $\{x_{ik}\}_{k=1, \dots, N_i}$

of independent observations on π_i , $i = 0, \dots, m$. (The sets $\{x_{ik}\}_{k=1, \dots, N_i}$, $i = 1, \dots, m$, comprise the identified observations of such samples, and such samples are said to be partially identified.) We distinguish samples of the two types according to whether the numbers N_i of identified observations contain information about the proportions α_i^0 , $i = 1, \dots, m$. If the numbers of identified observations contain no information about the proportions, then the sample is of the first type; otherwise, the sample is of the second type. The following are examples of how samples of the first and second types, respectively, might be obtained:

- (1) For $i = 0, \dots, m$, numbers N_i are arbitrarily chosen and independent observations $\{x_{ik}\}_{k=1, \dots, N_i}$ are obtained from π_i .
- (2) A number K_0 of observations are obtained from π_0 . For some $N_0 < K_0$, N_0 of these observations are left unidentified, while the remaining $K_0 - N_0$ observations are identified. For $i = 1, \dots, m$, a subset $\{x_{ik}\}_{k=1, \dots, N_i}$ of the identified observations is determined whose member observations come from π_i .

In the following, we consider likelihood equations determined by the two types of samples which are necessary conditions for a maximum-likelihood estimate. These equations, which were derived by Coberly [1], suggest certain successive-approximations iterative procedures for obtaining maximum-likelihood estimates. These procedures, which are generalized steepest ascent (deflected gradient) procedures, contain those of Hosmer [2] as a special case. Using arguments that parallel those of [3], we show that, with probability 1 as

N_0 approaches infinity (regardless of the relative sizes of N_0 and N_i , $i = 1, \dots, m$), these procedures converge locally to the strongly consistent maximum-likelihood estimates* whenever the step-size is between 0 and 2. Furthermore, the value of the step-size which yields optimal local convergence rates is bounded from below by a number which always lies between 1 and 2.

2. Samples of the first type.

We first assume that numbers $\{N_i\}_{i=0, \dots, m}$ are given and that, for $i = 0, \dots, m$, N_i independent observations $\{x_{ik}\}_{k=1, \dots, N_i}$ are drawn on π_i . The log-likelihood function for a sample of this type is

$$L_1(\theta) = \sum_{i=1}^m \sum_{k=1}^{N_i} \log p_i(x_{ik}) + \sum_{k=1}^{N_0} \log p(x_{0k}) .$$

In this expression, the parameter vector θ (with components $\alpha_i, \mu_i, \Sigma_i$, $i = 1, \dots, m$) belongs to the vector space $\mathcal{A} \oplus \mathcal{M} \oplus \mathcal{S}$ defined in [3], and the density functions on the right-hand side are evaluated with the true parameter vector θ^0 (with components $\alpha_i^0, \mu_i^0, \Sigma_i^0$, $i = 1, \dots, m$) replaced by θ .

*As in [3], one can show that, given any sufficiently small neighborhood of the true parameters, there is, with probability 1 as N_0 approaches infinity (regardless of the relative sizes of N_0 and N_i , $i = 1, \dots, m$), a unique solution of the likelihood equations for either type of sample in that neighborhood, and this solution is a maximum-likelihood estimate.

Differentiating $L_1(\theta)$ and setting its partial derivatives to zero gives the likelihood equations

$$(1.a) \quad \alpha_i = A_i(\theta) \equiv \frac{\alpha_i}{N_0} \sum_{k=1}^{N_0} \frac{p_i(x_{ok})}{p(x_{ok})}$$

$$(1.b) \quad \mu_i = M_i(\theta) \equiv \left\{ \sum_{k=1}^{N_i} x_{ik} + \sum_{k=1}^{N_0} x_{ok} \frac{\alpha_i p_i(x_{ok})}{p(x_{ok})} \right\} / \left\{ N_i + \sum_{k=1}^{N_0} \frac{\alpha_i p_i(x_{ok})}{p(x_{ok})} \right\}$$

$$(1.c) \quad \Sigma_i = S_i(\theta) \equiv \left\{ \sum_{k=1}^{N_i} (x_{ik} - \mu_i)(x_{ik} - \mu_i)^T + \sum_{k=1}^{N_0} (x_{ok} - \mu_i)(x_{ok} - \mu_i)^T \frac{\alpha_i p_i(x_{ok})}{p(x_{ok})} \right\} / \left\{ N_i + \sum_{k=1}^{N_0} \frac{\alpha_i p_i(x_{ok})}{p(x_{ok})} \right\}$$

for $i = 1, \dots, m$.

We set

$$\Lambda(\theta) = \begin{pmatrix} A_1(\theta) \\ \vdots \\ A_m(\theta) \end{pmatrix}, \quad M(\theta) = \begin{pmatrix} M_1(\theta) \\ \vdots \\ M_m(\theta) \end{pmatrix}, \quad S(\theta) = \begin{pmatrix} S_1(\theta) \\ \vdots \\ S_m(\theta) \end{pmatrix}$$

and define an operator Φ_ϵ on $\mathcal{A} \oplus \mathcal{M} \oplus \mathcal{S}$ by

$$\Phi_\epsilon(\theta) = (1 - \epsilon)\theta + \epsilon \begin{pmatrix} \Lambda(\theta) \\ M(\theta) \\ S(\theta) \end{pmatrix}.$$

Clearly, for any non-zero ϵ , the likelihood equations are satisfied by a vector $\theta \in \mathcal{A} \oplus \mathcal{M} \oplus \mathcal{S}$ if and only if $\theta = \Phi_\epsilon(\theta)$.

We consider the following iterative procedure: Beginning with some starting value $\theta^{(1)}$, define successive iterates inductively by

$$(2) \quad \theta^{(j+1)} = \Phi_\epsilon(\theta^{(j)})$$

for $j = 1, 2, 3, \dots$. Our local convergence result for this iterative procedure, as stated in the introduction, follows immediately from the theorem below.

Theorem 1: With probability 1 as N_0 approaches infinity, ϕ_ϵ is a locally contractive operator (in some norm on $\mathcal{A}(\Theta)$) near the strongly consistent maximum-likelihood estimate whenever $0 < \epsilon < 2$.

In saying that ϕ_ϵ is a locally contractive operator near a point $\theta \in \mathcal{A}(\Theta)$, we mean that there is a vector norm $\|\cdot\|$ on $\mathcal{A}(\Theta)$ and a number λ , $0 \leq \lambda < 1$, such that

$$\|\phi_\epsilon(\theta') - \theta\| \leq \lambda \|\theta' - \theta\|$$

whenever θ' lies sufficiently near θ .

Proof of Theorem 1: Let

$$\theta = \begin{pmatrix} \bar{\alpha} \\ \bar{\mu} \\ \bar{\Sigma} \end{pmatrix} = \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_m \\ \mu_1 \\ \vdots \\ \mu_m \\ \Sigma_1 \\ \vdots \\ \Sigma_m \end{pmatrix}$$

be the strongly consistent maximum-likelihood estimate. We assume that

$\alpha_i \neq 0$, $i = 1, \dots, m$. (As N_0 approaches infinity, the probability is 1 that this is the case.) As in [3], it suffices to show that, with probability 1, $\forall \phi_\epsilon(\theta)$ converges to an operator which has operator norm less than 1 with respect to a suitable vector norm on $\alpha \otimes \mathcal{M} \otimes \mathcal{S}$.

Now

$$\forall \phi_\epsilon(\theta) = (1 - \epsilon)I + \epsilon \forall \begin{pmatrix} A(\theta) \\ M(\theta) \\ S(\theta) \end{pmatrix},$$

and we write

$$\forall \begin{pmatrix} A \\ M \\ S \end{pmatrix} = \begin{pmatrix} \forall_\alpha A & \forall_\mu A & \forall_\Sigma A \\ \forall_\alpha M & \forall_\mu M & \forall_\Sigma M \\ \forall_\alpha S & \forall_\mu S & \forall_\Sigma S \end{pmatrix}.$$

Define inner products $\langle \cdot, \cdot \rangle'_i$ on \mathcal{M} , $\langle \cdot, \cdot \rangle''_i$ on \mathcal{S} , and $\langle \cdot, \cdot \rangle$ on $\alpha \otimes \mathcal{M} \otimes \mathcal{S}$ as in [3]. Setting

$$\beta_i(x) = \frac{p_i(x)}{p(x)}, \gamma_i(x) = (x - \mu_i), \delta_i(x) = [\Sigma_i^{-1}(x - \mu_i)(x - \mu_i)^T - I], K_i = N_i + \alpha_i N_0$$

for $i = 1, \dots, m$, one calculates

$$\forall_\alpha A(\theta) = I - (\text{diag } \alpha_i) \frac{1}{N_0} \frac{N_0}{\Sigma_0} \frac{1}{1} \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_m \end{pmatrix} \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_m \end{pmatrix}^T$$

$$\forall_\mu A(\theta) = - (\text{diag } \alpha_i) \frac{1}{N_0} \frac{N_0}{\Sigma_0} \frac{1}{1} \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_m \end{pmatrix} \begin{pmatrix} \langle \beta_1 \gamma_1, \cdot \rangle'_1 \\ \vdots \\ \langle \beta_m \gamma_m, \cdot \rangle'_m \end{pmatrix}^T$$

$$\forall_\Sigma A(\theta) = - (\text{diag } \alpha_i) \frac{1}{N_0} \frac{N_0}{\Sigma_0} \frac{1}{1} \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_m \end{pmatrix} \begin{pmatrix} \langle \beta_1 \delta_1, \cdot \rangle''_1 \\ \vdots \\ \langle \beta_m \delta_m, \cdot \rangle''_m \end{pmatrix}^T$$

$$V_{\alpha}^{-M}(0) = \left(\text{diag} \frac{1}{K_i} \sum_1^{N_0} \beta_i \gamma_i \right) - \left(\text{diag} \frac{\alpha_i}{K_i} \right) \left\{ \sum_1^{N_0} \begin{pmatrix} \beta_1 \gamma_1 \\ \vdots \\ \beta_m \gamma_m \end{pmatrix} \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_m \end{pmatrix}^T \right\}$$

$$V_{\mu}^{-M}(0) = \left(\text{diag} \frac{\alpha_i}{K_i} \sum_1^{N_0} \gamma_i \gamma_i^T \Sigma_i^{-1} \beta_i \right) - \left(\text{diag} \frac{\alpha_i}{K_i} \right) \left\{ \sum_1^{N_0} \begin{pmatrix} \beta_1 \gamma_1 \\ \vdots \\ \beta_m \gamma_m \end{pmatrix} \begin{pmatrix} \langle \beta_1 \gamma_1, \cdot \rangle'_1 \\ \vdots \\ \langle \beta_m \gamma_m, \cdot \rangle'_m \end{pmatrix}^T \right\}$$

$$V_{\Sigma}^{-M}(0) = \left(\text{diag} \frac{1}{K_i} \sum_1^{N_0} \beta_i \gamma_i \langle \delta_i, \cdot \rangle''_i \right) - \left(\text{diag} \frac{\alpha_i}{K_i} \right) \left\{ \sum_1^{N_0} \begin{pmatrix} \beta_1 \gamma_1 \\ \vdots \\ \beta_m \gamma_m \end{pmatrix} \begin{pmatrix} \langle \beta_1 \delta_1, \cdot \rangle''_1 \\ \vdots \\ \langle \beta_m \delta_m, \cdot \rangle''_m \end{pmatrix}^T \right\}$$

$$V_{\alpha}^{-S}(0) = \left(\text{diag} \frac{\Sigma_i}{K_i} \sum_1^{N_0} \beta_i \delta_i \right) - \left(\text{diag} \frac{\alpha_i \Sigma_i}{K_i} \right) \left\{ \sum_1^{N_0} \begin{pmatrix} \beta_1 \delta_1 \\ \vdots \\ \beta_m \delta_m \end{pmatrix} \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_m \end{pmatrix}^T \right\}$$

$$V_{\mu}^{-S}(0) = \left(\text{diag} \frac{1}{K_i} \left\{ -\sum_1^{N_i} [(\cdot) \gamma_i^T + \gamma_i (\cdot)^T] - \alpha_i \sum_1^{N_0} [(\cdot) \gamma_i^T + \gamma_i (\cdot)^T] \beta_i + \Sigma_i \sum_1^{N_0} \delta_i \langle \beta_i \gamma_i, \cdot \rangle'_i \right\} \right) -$$

$$- \left(\text{diag} \frac{\alpha_i \Sigma_i}{K_i} \right) \left\{ \sum_1^{N_0} \begin{pmatrix} \beta_1 \delta_1 \\ \vdots \\ \beta_m \delta_m \end{pmatrix} \begin{pmatrix} \langle \beta_1 \gamma_1, \cdot \rangle'_1 \\ \vdots \\ \langle \beta_m \gamma_m, \cdot \rangle'_m \end{pmatrix}^T \right\}$$

$$V_{\Sigma}^{-S}(0) = \left(\text{diag} \frac{\Sigma_i}{K_i} \sum_1^{N_0} \beta_i \delta_i \langle \delta_i, \cdot \rangle''_i \right) - \left(\text{diag} \frac{\alpha_i \Sigma_i}{K_i} \right) \left\{ \sum_1^{N_0} \begin{pmatrix} \beta_1 \delta_1 \\ \vdots \\ \beta_m \delta_m \end{pmatrix} \begin{pmatrix} \langle \beta_1 \delta_1, \cdot \rangle''_1 \\ \vdots \\ \langle \beta_m \delta_m, \cdot \rangle''_m \end{pmatrix}^T \right\}$$

Here, the arguments of β_i, γ_i and δ_i can be determined from the indices of summation, e.g.,

$$\sum_1^{N_0} \beta_i \gamma_i = \sum_{k=1}^{N_0} \beta_i(x_{ok}) \gamma_i(x_{ok}) .$$

Setting

$$V = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_m \\ \beta_1 \gamma_1 \\ \vdots \\ \beta_m \gamma_m \\ \beta_1 \delta_1 \\ \vdots \\ \beta_m \delta_m \end{pmatrix}$$

one obtains at 0

$$V \begin{pmatrix} A \\ M \\ S \end{pmatrix} = \begin{pmatrix} I & 0 & 0 \\ B_{21} & B_{22} & B_{23} \\ B_{31} & B_{32} & B_{33} \end{pmatrix} - \begin{pmatrix} (\text{diag } \frac{\alpha_i}{N_0}) & 0 & 0 \\ 0 & (\text{diag } \frac{\alpha_i}{K_i}) & 0 \\ 0 & 0 & (\text{diag } \frac{\alpha_i \gamma_i}{K_i}) \end{pmatrix} \begin{matrix} N_0 \\ \{ \sum_{k=1}^{N_0} V(x_{ok}) \cdot V(x_{ok}), \dots \} \\ 1 \end{matrix}$$

where

$$B_{21} = (\text{diag } \frac{1}{K_i} \sum_{l=1}^{N_0} \beta_l \gamma_l)$$

$$B_{22} = (\text{diag } \frac{\alpha_i}{K_i} \sum_{l=1}^{N_0} \gamma_l \gamma_l^T \Sigma_l^{-1} \beta_l)$$

$$B_{23} = (\text{diag } \frac{1}{K_i} \sum_{l=1}^{N_0} \beta_l \gamma_l \langle \delta_l, \cdot \rangle_l)$$

$$B_{31} = (\text{diag } \frac{\sum_l \beta_l}{K_i} \sum_{l=1}^{N_0} \beta_l \delta_l)$$

$$B_{32} = (\text{diag } \frac{1}{K_i} \{ - \sum_{l=1}^{N_0} [(\cdot) \gamma_l^T + \gamma_l (\cdot)^T] - \alpha_l \sum_{l=1}^{N_0} [(\cdot) \gamma_l^T + \gamma_l (\cdot)^T] \beta_l + \sum_{l=1}^{N_0} \delta_l \langle \beta_l \gamma_l, \cdot \rangle_l \})$$

$$B_{33} = (\text{diag } \frac{\sum_l \beta_l}{K_i} \sum_{l=1}^{N_0} \beta_l \delta_l \langle \delta_l, \cdot \rangle_l)$$

We have assumed that $\hat{\theta}$ is the strongly consistent maximum-likelihood estimate. Then, regardless of the relative sizes of N_i and N_0 , one can show as in [3] that, with probability 1, $\{\nabla\phi_c(\hat{\theta}) - E(\nabla\phi_c(\hat{\theta}^0))\}$ converges to zero as N_0 approaches infinity. Now

$$E\left(\nabla \begin{pmatrix} A(\hat{\theta}^0) \\ M(\hat{\theta}^0) \\ S(\hat{\theta}^0) \end{pmatrix}\right) = \begin{pmatrix} I & 0 & 0 \\ 0 & (\text{diag } \frac{\alpha_i^{0N_0}}{K_i} I) & 0 \\ 0 & 0 & (\text{diag } \frac{\alpha_i^{0N_0}}{K_i} I) \end{pmatrix} -$$

$$- \begin{pmatrix} (\text{diag } \alpha_i^0) & 0 & 0 \\ 0 & (\text{diag } \frac{\alpha_i^{0N_0}}{K_i} I) & 0 \\ 0 & 0 & (\text{diag } \frac{\alpha_i^{0N_0}}{K_i} \Sigma_i^e) \end{pmatrix} \left\{ \int_{\mathbb{R}^n} V(x) \langle V(x), \cdot \rangle p(x) dx \right\}$$

$$= B(1 - QR),$$

where

$$B = \begin{pmatrix} I & 0 & 0 \\ 0 & (\text{diag } \frac{\alpha_i^{0N_0}}{K_i} I) & 0 \\ 0 & 0 & (\text{diag } \frac{\alpha_i^{0N_0}}{K_i} I) \end{pmatrix}$$

$$Q = \begin{pmatrix} (\text{diag } \alpha_i^0) & 0 & 0 \\ 0 & I & 0 \\ 0 & 0 & (\text{diag } \Sigma_i^0) \end{pmatrix}$$

$$R = \int_{\mathbb{R}^n} V(x) \langle V(x), \cdot \rangle p(x) dx .$$

It was shown in [3] that QR is positive-definite and symmetric with operator norm less than 1 with respect to the inner product $\langle \cdot, Q^{-1} \cdot \rangle$ on $\mathcal{X} \oplus \mathcal{Y} \oplus \mathcal{Z}$. It follows that $I-QR$ is positive-definite and symmetric with norm less than 1 with respect to $\langle \cdot, Q^{-1} \cdot \rangle$. Since B and Q commute, $\langle \cdot, Q^{-1}B^{-1} \cdot \rangle$ is an inner product on $\mathcal{X} \oplus \mathcal{Y} \oplus \mathcal{Z}$, and one sees that $\langle W, Q^{-1}W \rangle \leq \langle W, Q^{-1}B^{-1}W \rangle$ for $W \in \mathcal{X} \oplus \mathcal{Y} \oplus \mathcal{Z}$. Consequently, $B(I-QR)$ is positive-definite and symmetric with norm less than 1 with respect to the inner product $\langle \cdot, Q^{-1}B^{-1} \cdot \rangle$. One concludes that

$$E(\nabla \phi_{\epsilon}(\theta^0)) = (1 - \epsilon)I + \epsilon E\left(\nabla \begin{pmatrix} A(\theta^0) \\ M(\theta^0) \\ S(\theta^0) \end{pmatrix}\right)$$

has norm less than 1 with respect to $\langle \cdot, Q^{-1}B^{-1} \cdot \rangle$ whenever $0 < \epsilon < 2$.

This completes the proof of the theorem.

We remark that, reasoning as in [3], one may determine a particular value of ϵ (the "optimal ϵ ") which yields, with probability 1 as N_0 approaches infinity, the fastest asymptotic uniform rates of local convergence of the iterative procedure (2) near θ . This optimal ϵ is given by

$$\epsilon = \frac{2}{2 - (\tau + \rho)}$$

where ρ and τ are, respectively the largest and smallest eigenvalues of $B(I-QR)$ regarded as an operator on $\mathcal{E} \oplus \mathcal{W} \oplus \mathcal{Z}$ (\mathcal{E} is the subspace of \mathcal{X} whose components sum to zero.) Since ρ and τ lie between zero and 1, one sees that the optimal ϵ is always greater than 1. If the component populations are "widely separated," then ρ and τ are near zero and,

hence, the optimal ϵ is near 1. If two or more of the component populations are nearly indistinguishable and if N_0 is large relative to the N_i 's, then τ is near zero, and the optimal ϵ cannot be much smaller than 2.

3. Samples of the second type.

We now assume that K_0 observations are obtained from the mixture population π_0 , and that, for some $N_0 < K_0$, N_0 of these observations are left unidentified, while the remaining $K_0 - N_0$ observations are identified. For $i = 1, \dots, m$, let $\{x_{ik}\}_{k=1, \dots, N_i}$ denote the subset of the identified observations which come from π_i , and let $\{x_{ok}\}_{k=1, \dots, N_0}$ be the set of unidentified observations from π_0 . The log-likelihood function for this sample is

$$L_2(\theta) = \log \left\{ \frac{(\sum_{i=1}^m N_i)!}{N_1! \dots N_m!} \alpha_1^{N_1} \dots \alpha_m^{N_m} \right\} + \sum_{i=1}^m \sum_{k=1}^{N_i} \log p_i(x_{ik}) + \sum_{k=1}^{N_0} \log p(x_{ok})$$

$$= \log \left\{ \frac{(\sum_{i=1}^m N_i)!}{N_1! \dots N_m!} \right\} + \sum_{i=1}^m \sum_{k=1}^{N_i} \log[\alpha_i p_i(x_{ik})] + \sum_{k=1}^{N_0} \log p(x_{ok}) .$$

Differentiating L_2 and setting its partial derivatives to zero gives the likelihood equations

$$(3.a) \quad \alpha_i = \tilde{\Lambda}_i(\theta) \equiv \frac{N_i}{K_0} + \frac{\alpha_i}{K_0} \sum_{k=1}^{N_0} \frac{p_i(x_{ok})}{p(x_{ok})}$$

$$(3.b) \quad \mu_i = M_i(\theta)$$

$$(3.c) \quad \Sigma_i = S_i(\theta)$$

for $i = 1, \dots, m$.

We set

$$\tilde{A}(\theta) = \begin{pmatrix} \tilde{A}_1(\theta) \\ \vdots \\ \tilde{A}_m(\theta) \end{pmatrix}$$

and define an operator $\tilde{\Phi}_\epsilon$ on $\mathcal{A}(\theta) \times \mathcal{S}$ by

$$\tilde{\Phi}_\epsilon(\theta) = (1 - \epsilon)\theta + \epsilon \begin{pmatrix} A(\theta) \\ M(\theta) \\ S(\theta) \end{pmatrix} .$$

Our iterative procedure is the following: Beginning with some starting value $\theta^{(1)}$, define successive iterates inductively by

$$(4) \quad \theta^{(j+1)} = \tilde{\Phi}_\epsilon(\theta^{(j)})$$

for $j = 1, 2, 3, \dots$. As before, the desired local convergence result for this iterative procedure follows from the theorem below.

Theorem 2: With probability 1 as N_0 approaches infinity, $\tilde{\Phi}_\epsilon$ is a locally contractive operator (in some norm on $\mathcal{A}(\theta) \times \mathcal{S}$) near the strongly consistent maximum-likelihood estimate whenever $0 < \epsilon < 2$.

Proof of Theorem 2: If θ is the strongly consistent maximum-likelihood estimate, then, as before, it suffices to show that, with probability 1, $\tilde{\Phi}_\epsilon(\theta)$ converges as N_0 approaches infinity to an operator which has operator norm less than 1 with respect to some vector norm on $\mathcal{A}(\theta) \times \mathcal{S}$. Proceeding as before, one sees that

$$V_{\bar{\alpha}} \tilde{\Lambda}(\Theta) = \left(\text{diag} \left(1 - \frac{N_i}{\alpha_i K_0} \right) \right) - \left(\text{diag} \frac{\alpha_i}{K_0} \right) \left\{ \sum_1^{N_0} \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_m \end{pmatrix} \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_m \end{pmatrix}^T \right\}$$

$$V_{\bar{\mu}} \tilde{\Lambda}(\Theta) = - \left(\text{diag} \frac{\alpha_i}{K_0} \right) \left\{ \sum_1^{N_0} \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_m \end{pmatrix} \begin{pmatrix} \langle \beta_1 \gamma_1, \cdot \rangle_1 \\ \vdots \\ \langle \beta_m \gamma_m, \cdot \rangle_m \end{pmatrix}^T \right\}$$

$$V_{\bar{\Sigma}} \tilde{\Lambda}(\Theta) = - \left(\text{diag} \frac{\alpha_i}{K_0} \right) \left\{ \sum_1^{N_0} \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_m \end{pmatrix} \begin{pmatrix} \langle \beta_1 \delta_1, \cdot \rangle_1 \\ \vdots \\ \langle \beta_m \delta_m, \cdot \rangle_m \end{pmatrix}^T \right\}$$

The remaining Fréchet derivatives, i.e., the derivatives at Θ of M and S with respect to $\bar{\alpha}$, $\bar{\mu}$, and $\bar{\Sigma}$, are unchanged, except that K_i must be replaced by $\alpha_i K_0$ wherever it appears.

One obtains at Θ

$$(4) \quad V \begin{pmatrix} \tilde{\Lambda} \\ M \\ S \end{pmatrix} = \begin{pmatrix} \left(\text{diag} \left(1 - \frac{N_i}{\alpha_i K_0} \right) \right) & 0 & 0 \\ \tilde{B}_{21} & \tilde{B}_{22} & \tilde{B}_{23} \\ \tilde{B}_{31} & \tilde{B}_{32} & \tilde{B}_{33} \end{pmatrix} - \begin{pmatrix} \left(\text{diag} \frac{\alpha_i}{K_0} \right) & 0 & 0 \\ 0 & \frac{1}{K_0} I & 0 \\ 0 & 0 & \left(\text{diag} \frac{\Sigma_i}{K_0} \right) \end{pmatrix} \left\{ \sum_{k=1}^{N_0} V(x_{0k}) \langle V(x_{0,k}), \cdot \rangle \right\}$$

In this expression, each \tilde{B}_{jk} is the same as the corresponding B_{jk} defined

previously, except that each K_1 in the latter is replaced by $\alpha_1 K_0$ in the former. One verifies that, with probability 1 as N_0 approaches infinity, (4) has the same limit as $\tilde{B}(I-QR)$, where Q and R are as before and $\tilde{B} = \frac{N_0}{K_0} I$. Repeating our earlier reasoning, one verifies that $\tilde{B}(I-QR)$ is positive-definite and symmetric with norm less than 1 with respect to the inner product $\langle \cdot, Q^{-1} \tilde{B}^{-1} \cdot \rangle$. Hence

$$V\tilde{\Phi}_\epsilon(\theta) = (1 - \epsilon) + \epsilon V \begin{pmatrix} \tilde{A}(\theta) \\ M(\theta) \\ S(\theta) \end{pmatrix}$$

converges to an operator which has norm less than 1 with respect to $\langle \cdot, Q^{-1} \tilde{B}^{-1} \cdot \rangle$ whenever $0 < \epsilon < 2$. This completes the proof of the theorem.

The remarks concerning the "optimal ϵ " at the conclusion of the preceding section are valid here verbatim.

BIBLIOGRAPHY

1. W. H. Coberly, private communication.
2. D. W. Hosmer, Jr., "A comparison of iterative maximum-likelihood estimates of the parameters of a mixture of two normal distributions under three different types of samples," Biometrics 29 (1973), pp. 761-770.
3. B. C. Peters, Jr., and H. F. Walker, "An iterative procedure for obtaining maximum-likelihood estimates of the parameters for a mixture of normal distributions," Report #51, NASA contract NAS-9-12777, University of Houston, Department of Mathematics.