

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-

144518

(NASA-CR-144518) ON THE NUMERICAL
EVALUATION OF THE MAXIMUM-LIKELIHOOD
ESTIMATE OF MIXTURE MEANS (Houston Univ.)

16 p HC \$3.25

CSSL 12A

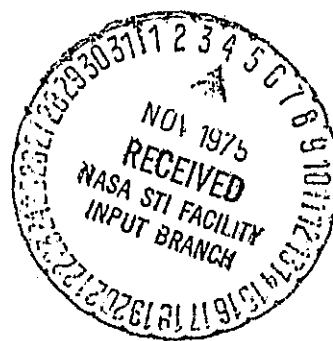
N76-10815

Unclas

G3/67 39654

ON THE NUMERICAL EVALUATION OF
THE MAXIMUM LIKELIHOOD ESTIMATE
OF MIXTURE MEANS
BY HOMER F. WALKER
REPORT #44 JULY, 1975

PREPARED FOR
EARTH OBSERVATION DIVISION, JSC
UNDER
CONTRACT NAS-9-12777



3801 CULLEN BLVD.
HOUSTON, TEXAS 77004

On the Numerical Evaluation of the Maximum-
Likelihood Estimate of Mixture Means

July, 1975

by

Homer F. Walker
Department of Mathematics
University of Houston

Report 44
NAS-9-12777

On the Numerical Evaluation of the Maximum-
Likelihood Estimate of Mixture Means

by

Homer F. Walker

Department of Mathematics, University of Houston
Houston, Texas

1. Introduction.

Let x be an n -dimensional random variable whose density function p is a convex combination of normal densities, i.e.,

$$p(x) = \sum_{i=1}^m \alpha_i^0 p_i(x) \quad \text{for } x \in \mathbb{R}^n,$$

where

$$\alpha_i^0 > 0, \quad \sum_{i=1}^m \alpha_i^0 = 1,$$

and

$$p_i(x) = \frac{1}{(2\pi)^{n/2} |\Sigma_i^0|^{1/2}} e^{-1/2 (x - \mu_i^0)^T \Sigma_i^0^{-1} (x - \mu_i^0)}$$

If $\{x_k\}_{k=1, \dots, N} \subseteq \mathbb{R}^n$ is an independent sample of observations on \mathbf{x} , then a maximum-likelihood estimate of the parameters $\{\alpha_i^0, \mu_i^0, \Sigma_i^0\}_{i=1, \dots, m}$ is a choice of parameters $\{\alpha_i^0, \mu_i^0, \Sigma_i^0\}_{i=1, \dots, m}$ which locally maximizes the log-likelihood function

$$L = \sum_{k=1}^N \log p(x_k),$$

in which p is evaluated with the true parameters $\{\alpha_i^0, \mu_i^0, \Sigma_i^0\}_{i=1, \dots, m}$ replaced by the estimate $\{\alpha_i, \mu_i, \Sigma_i\}_{i=1, \dots, m}$. (In the following, it is clear from the context which parameters are used in evaluating the density functions p_i and p . Therefore, these parameters are not explicitly pointed out.)

Clearly, L is a differentiable function of the parameters to be estimated. Equating to zero the partial derivatives of L with respect to these parameters, one obtains, after a straightforward calculations the following necessary conditions for a maximum-likelihood estimate:

$$(1.a) \quad \alpha_i = \frac{\alpha_i}{N} \sum_{k=1}^N \frac{p_i(x_k)}{p(x_k)}$$

$$(1.b) \quad \mu_i = \left\{ \frac{1}{N} \sum_{k=1}^N x_k \frac{p_i(x_k)}{p(x_k)} \right\} / \left\{ \frac{1}{N} \sum_{k=1}^N \frac{p_i(x_k)}{p(x_k)} \right\} \quad \left. \vphantom{\frac{1}{N} \sum_{k=1}^N} \right\} i=1, \dots, m.$$

$$(1.c) \quad \Sigma_i = \left\{ \frac{1}{N} \sum_{k=1}^N (x_k - \mu_i)(x_k - \mu_i)^T \frac{p_i(x_k)}{p(x_k)} \right\} / \left\{ \frac{1}{N} \sum_{k=1}^N \frac{p_i(x_k)}{p(x_k)} \right\}$$

These are known as the likelihood equations, and we shall assume that the parameters under consideration here are restricted to sets in which these equations are sufficient, as well as necessary, for a maximum-likelihood estimate.

The likelihood equations suggest the following iterative procedure for obtaining a solution: Beginning with some set of starting values, obtain successive approximations to a solution by inserting the preceding approximations in the expressions on the right-hand sides of (1.a), (1.b), and (1.c). This scheme is attractive for its relative ease of implementation, and it has been investigated by a number of authors. Empirical studies of Day [1], Duda and Hart [2], and Hasselblad [3] suggest that this scheme is convergent and that convergence is particularly fast when the component normal densities in p are "widely separated" in a certain sense. No proof of convergence is given in these papers, although Peters and Walker [8] have shown that, with probability approaching 1 as N approaches infinity, a related procedure (which includes this one as a special case) converges locally to the consistent maximum-likelihood estimate whenever a certain "step-size" is sufficiently small. (An iterative procedure is said to converge locally to a limit if the iterates converge to that limit whenever the starting values are sufficiently near that limit.)

Peters and Coberly [7] have proved that, if all of the parameters μ_i and Σ_i are held fixed, then the iterative procedure suggested by the equations (1.a) alone converges locally to a maximum-likelihood estimate of the parameters α_i , $i=1, \dots, m$. They also report on numerical studies in which the computational feasibility of this procedure is demonstrated. In this note, we provide sufficient conditions for the iterative procedure suggested by the equations (1.b) alone, for fixed parameters α_i and Σ_i , to converge locally to a maximum-likelihood estimate of the means μ_i , $i=1, \dots, m$. These conditions are, roughly, that either $m = 2$ or the component normal densities in p be "widely separated" in a certain sense.

2. Preliminary discussion.

We denote by \mathcal{M} the m -fold direct sum of \mathbb{R}^n with itself, and we represent its elements as columns

$$\bar{\mu} = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_m \end{pmatrix} \in \mathcal{M}.$$

(Of course, \mathcal{M} is isomorphic to \mathbb{R}^{mn} .) We also find it convenient to represent parameter sets $\{\alpha_i\}_{i=1,\dots,m}$ and $\{\Sigma_i\}_{i=1,\dots,m}$ as columns

$$\bar{\alpha} = \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_m \end{pmatrix} \quad \text{and} \quad \bar{\Sigma} = \begin{pmatrix} \Sigma_1 \\ \vdots \\ \Sigma_m \end{pmatrix},$$

and, in the following, we use the fact that $\bar{\alpha}$ and $\bar{\Sigma}$ belong to normed vector spaces without explicitly introducing these spaces or their norms.

Setting

$$M_i(\bar{\alpha}, \bar{\mu}, \bar{\Sigma}) = \left\{ \frac{1}{N} \sum_{k=1}^N x_k \frac{p_1(x_k)}{p(x_k)} \right\} / \left\{ \frac{1}{N} \sum_{k=1}^N \frac{p_1(x_k)}{p(x_k)} \right\}, \quad i=1,\dots,m,$$

we define

$$M(\bar{\alpha}, \bar{\mu}, \bar{\Sigma}) = \begin{pmatrix} M_1(\bar{\alpha}, \bar{\mu}, \bar{\Sigma}) \\ \vdots \\ M_m(\bar{\alpha}, \bar{\mu}, \bar{\Sigma}) \end{pmatrix},$$

which we regard as a function from \mathcal{M} to itself depending on parameters $\bar{\alpha}$ and $\bar{\Sigma}$. The equations (1.b) can now be written as

$$(2) \quad \bar{\mu} = M(\bar{\alpha}, \bar{\mu}, \bar{\Sigma}),$$

and the iterative procedure under consideration is the following: Beginning with some starting value $\bar{\mu}^{(1)}$, define successive iterates inductively by

$$(3) \quad \bar{\mu}^{(k+1)} = M(\bar{\alpha}, \bar{\mu}^{(k)}, \bar{\Sigma})$$

for $k = 1, 2, \dots$

In our results concerning the convergence of the procedure (3), the Fréchet derivative of M with respect to $\bar{\mu}$, which we denote by $\nabla_{\bar{\mu}} M$, is of central importance. (For questions concerning the definition and properties of Fréchet derivatives, see Luenberger [6].) Indeed, if $\bar{\alpha}, \bar{\mu}$, and $\bar{\Sigma}$ satisfy (2) and if $|| \cdot ||$ is any norm on \mathcal{M} , then one can write

$$M(\bar{\alpha}, \bar{\mu}', \bar{\Sigma}) - \bar{\mu} = \nabla_{\bar{\mu}} M(\bar{\alpha}, \bar{\mu}, \bar{\Sigma}) (\bar{\mu}' - \bar{\mu}) + o(||\bar{\mu}' - \bar{\mu}||^2)$$

for $\bar{\mu}'$ near $\bar{\mu}$. Consequently, if there exists a norm $|| \cdot ||$ on \mathcal{M} with respect to which $\nabla_{\bar{\mu}} M(\bar{\alpha}, \bar{\mu}, \bar{\Sigma})$ has operator norm less than 1, then M is locally contractive in that norm near $\bar{\mu}$, i.e., there is a number λ , $0 \leq \lambda < 1$, such that

$$(4) \quad ||M(\bar{\alpha}, \bar{\mu}', \bar{\Sigma}) - \bar{\mu}|| \leq \lambda ||\bar{\mu}' - \bar{\mu}||$$

whenever $\bar{\mu}'$ is sufficiently near $\bar{\mu}$. Since an inequality of the form (4) implies the local convergence of the iterative procedure (3) to $\bar{\mu}$, our objectives will be met by giving sufficient conditions for $\nabla_{\bar{\mu}} M(\bar{\alpha}, \bar{\mu}, \bar{\Sigma})$ to have operator norm less than 1 (with respect to some norm on \mathcal{M}) at parameter vectors $\bar{\alpha}, \bar{\mu}$, and $\bar{\Sigma}$ which satisfy (2).

We now calculate $\nabla_{\bar{\mu}} M$ at a set of parameter vectors $\hat{\alpha}, \hat{\mu}$, and $\hat{\Sigma}$ (with components $\hat{\alpha}_i, \hat{\mu}_i$, and $\hat{\Sigma}_i$, $i=1, \dots, m$) which satisfy the likelihood equations. We first define inner products $\langle \cdot, \cdot \rangle_i$ on \mathbb{R}^n by

$$\langle x, y \rangle_i = x^T (\hat{\alpha}_i \hat{\Sigma}_i^{-1}) y \quad \text{for } x, y \in \mathbb{R}^n, i=1, \dots, m.$$

Then, denoting the Fréchet derivative of M_i with respect to μ_j by $\nabla_{\mu_j} M_i$, one verifies with the aid of the likelihood equations that

$$\nabla_{\mu_j} M_i(\hat{\alpha}, \hat{\mu}, \hat{\Sigma}) = \begin{cases} -\frac{1}{N} \sum_{k=1}^N \frac{p_i(x_k)}{p(x_k)} (x_k - \hat{\mu}_i) \langle \frac{p_j(x_k)}{p(x_k)} (x_k - \hat{\mu}_j), \cdot \rangle_j & \text{if } i \neq j \\ I - \frac{1}{N} \sum_{k=1}^N \frac{p_i(x_k)}{p(x_k)} (x_k - \hat{\mu}_i) \langle \frac{p_i(x_k)}{p(x_k)} (x_k - \hat{\mu}_i), \cdot \rangle_i & \text{if } i = j. \end{cases}$$

This yields the following expression, in the form of a matrix of Fréchet derivatives, for $\nabla_{\bar{\mu}} M$ at a solution of the likelihood equations:

$$(5) \quad \nabla_{\bar{\mu}} M(\hat{\alpha}, \hat{\mu}, \hat{\Sigma}) = \begin{pmatrix} \nabla_{\mu_1} M_1(\hat{\alpha}, \hat{\mu}, \hat{\Sigma}) & \dots & \nabla_{\mu_m} M_1(\hat{\alpha}, \hat{\mu}, \hat{\Sigma}) \\ \vdots & \ddots & \vdots \\ \nabla_{\mu_1} M_m(\hat{\alpha}, \hat{\mu}, \hat{\Sigma}) & \dots & \nabla_{\mu_m} M_m(\hat{\alpha}, \hat{\mu}, \hat{\Sigma}) \end{pmatrix}$$

$$= I - \left\{ \frac{1}{N} \sum_{k=1}^N \begin{pmatrix} \frac{p_1(x_k)}{p(x_k)} (x_k - \hat{\mu}_1) \\ \vdots \\ \frac{p_m(x_k)}{p(x_k)} (x_k - \hat{\mu}_m) \end{pmatrix} \begin{pmatrix} \frac{p_1(x_k)}{p(x_k)} (x_k - \hat{\mu}_1), \cdot - 1 \\ \vdots \\ \frac{p_m(x_k)}{p(x_k)} (x_k - \hat{\mu}_m), \cdot - m \end{pmatrix}^T \right\}.$$

The inner products $\langle \cdot, \cdot \rangle_{\mu}$ induce an inner product $\langle \cdot, \cdot \rangle$ on \mathcal{M} . In the following, $|| \cdot ||$ will denote both the vector norm and the operator norm defined by this inner product. It is apparent from (5) that, at a solution of the likelihood equations, $\nabla_{\mu} M$ is of the form $I - Q$, where Q is positive semi-definite and symmetric with respect to the inner product $\langle \cdot, \cdot \rangle$. In fact, we prove in an appendix that Q is positive-definite with probability 1 whenever $N \geq mn$. It follows that, with probability 1 for $N \geq mn$, $||\nabla_{\mu} M|| < 1$ at a solution of the likelihood equations if and only if $||Q|| < 2$. We conclude these preliminary remarks with the following

Lemma: $||Q|| < m$.

Proof: Since Q is symmetric with respect to $\langle \cdot, \cdot \rangle$, one has

$$||Q|| = \sup_{\bar{v} \in \mathcal{M}} \langle \bar{v}, Q\bar{v} \rangle.$$

$$||\bar{v}|| \leq 1$$

If $\{v_i\}_{i=1, \dots, m} \subseteq \mathbb{R}^n$ is such that

$$\bar{v} = \begin{pmatrix} v_1 \\ \vdots \\ v_m \end{pmatrix} \in \mathcal{M} \text{ satisfies } ||\bar{v}|| \leq 1, \text{ then}$$

$$\begin{aligned}
\langle \bar{v}, Q\bar{v} \rangle &= \sum_{i=1}^m \sum_{j=1}^m \frac{1}{N} \sum_{k=1}^N \langle v_i, \frac{p_i(x_k)}{p(x_k)} (x_k - \hat{\mu}_i) \rangle_i \langle v_j, \frac{p_j(x_k)}{p(x_k)} (x_k - \hat{\mu}_j) \rangle_j \\
&\leq \sum_{i=1}^m \sum_{j=1}^m \left\{ \frac{1}{N} \sum_{k=1}^N \langle v_i, \frac{p_i(x_k)}{p(x_k)} (x_k - \hat{\mu}_i) \rangle_i^2 \right\}^{1/2} \left\{ \frac{1}{N} \sum_{k=1}^N \langle v_j, \frac{p_j(x_k)}{p(x_k)} (x_k - \hat{\mu}_j) \rangle_j^2 \right\}^{1/2} \\
&< \sum_{i=1}^m \sum_{j=1}^m \langle v_i, [\frac{1}{N} \sum_{k=1}^N (x_k - \hat{\mu}_i)(x_k - \hat{\mu}_i)^T \frac{p_i(x_k)}{p(x_k)}] \hat{\Sigma}_i^{-1} v_i \rangle_i^{1/2} \\
&\quad \cdot \langle v_j, [\frac{1}{N} \sum_{k=1}^N (x_k - \hat{\mu}_j)(x_k - \hat{\mu}_j)^T \frac{p_j(x_k)}{p(x_k)}] \hat{\Sigma}_j^{-1} v_j \rangle_j^{1/2}
\end{aligned}$$

since $\frac{\hat{\alpha}_i p_i(x)}{p(x)} < 1$ for $i=1, \dots, m$. From the likelihood equations, one concludes that

$$\langle \bar{v}, Q\bar{v} \rangle < \sum_{i=1}^m \sum_{j=1}^m \langle v_i, v_i \rangle_i^{1/2} \langle v_j, v_j \rangle_j^{1/2} = (\sum_{i=1}^m \langle v_i, v_i \rangle_i^{1/2})^2 \leq m,$$

and the lemma is proved.

3. Sufficient conditions for local convergence.

Sufficient conditions will now be given for local convergence of the procedure (3) to a solution of (2). Our first condition is given by the theorem below.

Theorem 1: Suppose that $m = 2$ and $N \geq 2n$, and let $\frac{\hat{\alpha}}{\alpha, \hat{\mu}, \hat{\Sigma}}$ be vectors of parameters which satisfy the likelihood equations. If $\bar{\alpha}, \bar{\mu}$, and $\bar{\Sigma}$, satisfy (2) and lie sufficiently near $\frac{\hat{\alpha}}{\alpha, \hat{\mu}}$, and $\frac{\hat{\Sigma}}{\Sigma}$, then the iterative procedure (3) converges locally to $\bar{\mu}$ with probability 1.

Proof: From the preliminary discussion, we know that the procedure (3) converges locally to $\bar{\mu}$ if $V_{\bar{\mu}}^M(\bar{\alpha}, \bar{\mu}, \bar{\Sigma})$ has operator norm less than 1 with respect to some vector norm on \mathcal{M} . Then, since $V_{\bar{\mu}}^M$ depends continuously on $\bar{\alpha}, \bar{\mu}$, and $\bar{\Sigma}$, it suffices to find a norm on \mathcal{M} with respect to which $V_{\bar{\mu}}^M(\hat{\alpha}, \hat{\mu}, \hat{\Sigma})$ has operator norm less than 1 in order to prove the theorem.

Now $V_{\bar{\mu}}^M(\hat{\alpha}, \hat{\mu}, \hat{\Sigma}) = I - Q$, where Q is the operator introduced in the preliminary discussion. With probability 1, Q is positive-definite as well as symmetric with respect to $\langle \cdot, \cdot \rangle$, and, from the Lemma, $\|Q\| < m = 2$. Consequently, $\|V_{\bar{\mu}}^M(\hat{\alpha}, \hat{\mu}, \hat{\Sigma})\| < 1$ with probability 1, and the proof is complete.

We now define an operator Q^0 on \mathcal{M} by

$$Q^0 = \int_{\mathbb{R}^n} \begin{pmatrix} \frac{p_1(x)}{p(x)}(x - \mu_1^0) \\ \vdots \\ \frac{p_m(x)}{p(x)}(x - \mu_m^0) \end{pmatrix} \begin{pmatrix} \langle \frac{p_1(x)}{p(x)}(x - \mu_1^0), \cdot \rangle_1 \\ \vdots \\ \langle \frac{p_m(x)}{p(x)}(x - \mu_m^0), \cdot \rangle_m \end{pmatrix}^T p(x) dx,$$

in which the true parameters (whose vectors we denote by $\bar{\alpha}^0, \bar{\mu}^0$, and $\bar{\Sigma}^0$) are used in evaluating the functions p_i and p and the inner products $\langle \cdot, \cdot \rangle_i$. The operator Q^0 can be thought of as an $m \times m$ array of operators on \mathbb{R}^n , the ij th operator of which is

$$\int_{\mathbb{R}^n} \frac{p_i(x)}{p(x)}(x - \mu_i^0) \langle \frac{p_j(x)}{p(x)}(x - \mu_j^0), \cdot \rangle_j p(x) dx.$$

If the component normal densities in p are "widely separated" in the sense that each pair of parameters μ_i^0 and μ_j^0 differs greatly from every other

pair, then the off-diagonal operators in this array are near zero. On the other hand, regardless of the "separation" of the component densities, the diagonal operators define an operator on \mathcal{M} which lies strictly between the zero operator and the identity operator in the ordering on symmetric operators defined by the inner product $\langle \cdot, \cdot \rangle$. Consequently, if the component normal densities in p are sufficiently "widely separated" in this sense, then the operator $I - Q^0$ has spectral radius less than 1, and, hence, there exists a norm on \mathcal{M} with respect to which $I - Q^0$ has operator norm less than 1. (See Householder [4].) This motivates our second condition.

Theorem 2: Suppose that the component normal densities in p are sufficiently "widely separated" that the spectral radius of $I - Q^0$ is less than 1. Then the probability is 1 that, for sufficiently large N , there exist neighborhoods of $\bar{\alpha}^0, \bar{\mu}^0$, and $\bar{\Sigma}^0$ such that, if $\bar{\alpha}, \bar{\mu}$, and $\bar{\Sigma}$ lie in these neighborhoods and satisfy (2), then the iterative procedure (3) converges locally to $\bar{\mu}$.

Proof: A straightforward calculation and an application of the Strong Law of Large Numbers (see Loève [5]) yields that $\nabla_{\bar{\mu}} M(\bar{\alpha}^0, \bar{\mu}^0, \bar{\Sigma}^0)$ converges with probability 1 to $I - Q^0$ as N approaches infinity. Since $I - Q^0$ is assumed to have spectral radius less than 1, it follows that, with probability 1, if N is sufficiently large, then $\nabla_{\bar{\mu}} M(\bar{\alpha}, \bar{\mu}, \bar{\Sigma})$ has operator norm less than 1 with respect to some norm on \mathcal{M} whenever $\bar{\alpha}, \bar{\mu}$, and $\bar{\Sigma}$ lie near $\bar{\alpha}^0, \bar{\mu}^0$, and $\bar{\Sigma}^0$. If $\nabla_{\bar{\mu}} M(\bar{\alpha}, \bar{\mu}, \bar{\Sigma})$ has operator norm less than 1 and $\bar{\alpha}, \bar{\mu}$, and $\bar{\Sigma}$ also satisfy (2), then the iterative procedure (3) converges locally to $\bar{\mu}$. This completes the proof.

Appendix

We now prove that the operator

$$Q = \frac{1}{N} \sum_{k=1}^N \begin{pmatrix} \frac{p_1(x_k)}{p(x_k)}(x_k^{-\mu_1}) \\ \vdots \\ \frac{p_m(x_k)}{p(x_k)}(x_k^{-\mu_m}) \end{pmatrix} \begin{pmatrix} \langle \frac{p_1(x_k)}{p(x_k)}(x_k^{-\mu_1}), \cdot \rangle_1 \\ \vdots \\ \langle \frac{p_m(x_k)}{p(x_k)}(x_k^{-\mu_m}), \cdot \rangle_m \end{pmatrix}^T$$

is positive-definite on \mathcal{M} with probability 1 whenever $N \geq mn$. Clearly, it suffices to show that the vectors

$$v(x_k) \equiv \begin{pmatrix} \frac{p_1(x_k)}{p(x_k)}(x_k^{-\mu_1}) \\ \vdots \\ \frac{p_m(x_k)}{p(x_k)}(x_k^{-\mu_m}) \end{pmatrix}, \quad k = 1, \dots, N,$$

span \mathcal{M} with probability 1 whenever $N \geq mn$. This follows from the more general result below.

Lemma. Let $\{x_k\}_{k=1, \dots, N}$ be an independent sample of observations on a random variable x in \mathbb{R}^s which is distributed with a probability density function p . If V is a real-analytic function from \mathbb{R}^s to \mathbb{R}^t whose component functions are linearly independent, then the vectors $v(x_k)$, $k=1, \dots, N$, span \mathbb{R}^t with probability 1 whenever $N \geq t$.

Proof: Denoting the j^{th} component function of V by v_j , we define a real-analytic function V_j from \mathbb{R}^s to \mathbb{R}^j by

$$V_j(x) = \begin{pmatrix} v_1(x) \\ \vdots \\ v_j(x) \end{pmatrix}$$

for $j = 1, \dots, t$. Our proof of the lemma consists of showing inductively that, for $j = 1, \dots, t$, the set $\{V_j(x_k)\}_{k=1, \dots, j}$ spans \mathbb{R}^j with probability 1. We make the preliminary observation that, since the real-analytic functions v_j are assumed to be linearly independent, any non-zero linear combination of them vanishes only on a set of Lebesgue measure zero in \mathbb{R}^s .

From the observation above, $V_1(x_1)$ is non-zero with probability 1; hence $V_1(x_1)$ spans \mathbb{R}^1 with probability 1. Suppose now that, for some j , $1 \leq j < t$, the set $\{V_j(x_k)\}_{k=1, \dots, j}$ spans \mathbb{R}^j with probability 1. Then, with probability 1, the set $\{V_{j+1}(x_k)\}_{k=1, \dots, j+1}$ fails to span \mathbb{R}^{j+1} if and only if

$$(*) \quad V_{j+1}(x_{j+1}) = \sum_{k=1}^j c_k V_{j+1}(x_k)$$

for some set of constants $\{c_k\}_{k=1, \dots, j}$. If $(*)$ holds, the constants c_k are determined by

$$\begin{pmatrix} c_1 \\ \vdots \\ c_j \end{pmatrix} = V_j^{-1} V_j(x_{j+1})$$

with probability 1, where \tilde{V}_j is the $j \times j$ matrix whose k^{th} column is $V_j(x_k)$. Thus, with probability 1, (*) holds if and only if

$$[\tilde{V}_j^{-1} V_j(x_{j+1})]^T \begin{pmatrix} v_{j+1}(x_1) \\ \vdots \\ v_{j+1}(x_j) \end{pmatrix} - v_{j+1}(x_{j+1}) = 0.$$

Now

$$[\tilde{V}_j^{-1} V_j(x)]^T \begin{pmatrix} v_{j+1}(x_1) \\ \vdots \\ v_{j+1}(x_j) \end{pmatrix} - v_{j+1}(x)$$

is a non-zero linear combination of the functions v_1, \dots, v_{j+1} and, hence, vanishes only on a set of Lebesgue measure zero in \mathbb{R}^S . One concludes that $\{V_{j+1}(x_k)\}_{k=1, \dots, j+1}$ fails to span \mathbb{R}^{j+1} with probability zero. This completes the induction, and the lemma is proved.

REFERENCES

1. N.E. Day, "Estimating the components of a mixture of normal distributions," Biometrika 56 (1969), pp. 463-474.
2. R.O. Duda and P.E. Hart, Pattern Classification and Scene Analysis, John Wiley and Sons, Inc., New York, 1973.
3. V. Hasselblad, "Estimation of parameters for a mixture of normal distributions," Technometrics 8 (1966), pp. 431-446.
4. A. Householder, Theory of Matrices and Numerical Analysis, Blaisdell Publishing Co., New York, 1964.
5. M. Loève, Probability Theory, D. Van Nostrand Co., New York, 1963.
6. D.G. Luenberger, Optimization by Vector Space Methods, John Wiley and Sons, Inc., New York, 1969.
7. B.C. Peters and W.H. Coberly, "The numerical evaluation of the maximum-likelihood estimate of mixture proportions," to appear.
8. B.C. Peters and H.F. Walker, "An iterative procedure for obtaining maximum-likelihood estimates of the parameters for a mixture of normal distributions," to appear.