



## CONTRIBUTED ARTICLE

# Convergence Results for the EM Approach to Mixtures of Experts Architectures

MICHAEL I. JORDAN AND LEI XU

Massachusetts Institute of Technology

(Received 27 September 1993; revised and accepted 10 January 1995)

**Abstract**—The Expectation-Maximization (EM) algorithm is an iterative approach to maximum likelihood parameter estimation. Jordan and Jacobs recently proposed an EM algorithm for the mixture of experts architecture of Jacobs, Jordan, Nowlan and Hinton (1991) and the hierarchical mixture of experts architecture of Jordan and Jacobs (1992). They showed empirically that the EM algorithm for these architectures yields significantly faster convergence than gradient ascent. In the current paper we provide a theoretical analysis of this algorithm. We show that the algorithm can be regarded as a variable metric algorithm with its searching direction having a positive projection on the gradient of the log likelihood. We also analyze the convergence of the algorithm and provide an explicit expression for the convergence rate. In addition, we describe an acceleration technique that yields a significant speedup in simulation experiments.

**Keywords**—Supervised learning, Statistical models, Maximum likelihood, EM algorithm, Convergence rate, Mixture models, Hierarchical models, Optimization.

## 1. INTRODUCTION

Although neural networks are capable in principle of representing complex nonlinear functions, the time required to train a complex network does not always scale well with problem size and the solution obtained does not always reveal the structure in the problem. Moreover, it is often difficult to express prior knowledge in the language of fully connected neural networks. Achieving better scaling behavior, better interpretability of solutions, and better ways of incorporating prior knowledge may require a more modular approach in which the learning problem is decomposed into subproblems. Such an approach has been used with success in the statistics literature and the machine learning literature, where decision tree

algorithms such as CART and ID3 and multivariate spline algorithms such as MARS have running times that can be orders of magnitude faster than neural network algorithms and often yield simple, interpretable solutions (Breiman et al., 1984; Friedman, 1991).

A general strategy for designing modular learning systems is to treat the problem as one of combining multiple models, each of which is defined over a local region of the input space. Jacobs et al. (1991) introduced such a strategy with their “mixture of experts” (ME) architecture for supervised learning. The architecture involves a set of function approximators (“expert networks”) that are combined by a classifier (“gating network”). These networks are trained simultaneously so as to split the input space into regions where particular experts can specialize. Jordan and Jacobs (1992) extended this approach to a recursively defined architecture in which a tree of gating networks combines the expert networks into successively larger groupings that are defined over nested regions of the input space. This “hierarchical mixture of experts” (HME) architecture is closely related to the decision tree and multivariate spline algorithms.

The problem of training a mixture of experts architecture can be treated as a maximum likelihood estimation problem. Both Jacobs et al. (1991) and

---

Acknowledgements: The project was supported in part by a grant from the McDonnell-Pew Foundation, by a grant from ATR Human Information Processing Research Laboratories, by a grant from Siemens Corporation, by grant IRI-9013991 from the National Science Foundation, and by grant N00014-90-J-1942 from the Office of Naval Research. The project was also supported by NSF grant ASC-9217041 in support of the Center for Biological and Computational Learning at MIT, including funds provided by DARPA under the HPCC program, and NSF grant ECS-9216531 to support an Initiative in Intelligent Control at MIT. Michael I. Jordan is a NSF Presidential Young Investigator.

Requests for reprints should be sent to Michael I. Jordan, E10-034D, MIT, Cambridge, MA 02139, USA.

Jordan and Jacobs (1992) derived learning algorithms by computing the gradient of the log likelihood for their respective architectures. Empirical tests revealed that although the gradient approach succeeded in finding reasonable parameter values in particular problems, the convergence rate was not significantly better than that obtained by using gradient methods in multilayered neural network architectures. The gradient approach did not appear to take advantage of the modularity of the architecture. An alternative to the gradient approach was proposed by Jordan and Jacobs (1994), who introduced an Expectation-Maximization (EM) algorithm for mixture of experts architectures. EM is a general technique for maximum likelihood estimation that can often yield simple and elegant algorithms (Baum et al., 1970; Dempster, Laird & Rubin, 1977). For mixture of experts architectures, the EM algorithm decouples the estimation process in a manner that fits well with the modular structure of the architecture. Moreover, Jordan and Jacobs (1994) observed a significant speedup over gradient techniques.

In this paper, we provide further insight into the EM approach to mixtures of experts architectures via a set of convergence theorems.<sup>1</sup> We study a particular variant of the EM algorithm proposed by Jordan and Jacobs (1994) and demonstrate a relationship between this algorithm and gradient ascent. We also provide theorems on the convergence rate of the algorithm and provide explicit formulas for the constants.

The remainder of the paper is organized as follows. Section 2 introduces the ME model. The EM algorithm for this architecture is derived and two convergence theorems are presented. Section 3 presents an analogous derivation and a set of convergence results for the HME model. Section 4 introduces two acceleration techniques for improving convergence and presents the results of numerical experiments. Section 5 presents our conclusions.

## 2. THEORETICAL ANALYSIS OF AN EM ALGORITHM FOR THE MIXTURE OF EXPERTS ARCHITECTURE

### 2.1. Network Learning Based on Maximum Likelihood Estimation

We begin by studying the nonhierarchical case. As shown in Figure 1, the mixture of experts (ME)

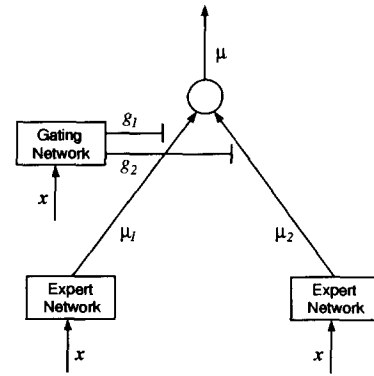


FIGURE 1. The mixture of experts architecture. The total output  $\mu$  is the weighted sum of the expert network outputs:  $\mu = g_1\mu_1 + g_2\mu_2$ , where the weights are the gating network outputs  $g_1$  and  $g_2$ .

architecture is comprised of  $K$  expert networks, each of which solves a function approximation problem over a local region of the input space. To each expert network we associate a probabilistic model that relates input vectors  $\mathbf{x} \in R^n$  to output vectors  $\mathbf{y} \in R^m$ . We denote these probabilistic models as follows

$$P(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}_j), \quad j = 1, 2, \dots, K,$$

where the  $\boldsymbol{\theta}_j$  are parameter vectors. Each of these probability densities can be generally assumed to belong to the exponential family of densities (Jordan & Jacobs, 1994); in the current paper, however, we restrict our analysis to Gaussian models for simplicity. The  $j$ th expert network produces as output a parameter vector  $\boldsymbol{\mu}_j$

$$\boldsymbol{\mu}_j = \mathbf{f}_j(\mathbf{x}, \boldsymbol{\theta}_j), \quad j = 1, 2, \dots, K,$$

which is the location parameter for the  $j$ th probability density. In the current paper, as in Jordan and Jacobs (1994), we treat the case in which the functions  $\mathbf{f}_j$  are linear in the parameters. We extend our results to the case of experts that are nonlinear in the parameters in the Appendix.

For Gaussian probability densities  $P(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}_j)$  the location parameter  $\boldsymbol{\mu}_j$  is simply the mean. We also associate a nonsingular covariance matrix  $\Sigma_j$  with each expert network, yielding the following probabilistic model for expert  $j$

$$P(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}_j) = \frac{1}{(2\pi)^{m/2} |\Sigma_j|^{1/2}} \times \exp\left\{-\frac{1}{2} [\mathbf{y} - \mathbf{f}_j(\mathbf{x}, \boldsymbol{\theta}_j)]^T \Sigma_j^{-1} [\mathbf{y} - \mathbf{f}_j(\mathbf{x}, \boldsymbol{\theta}_j)]\right\}. \quad (1)$$

The ME architecture also utilizes an auxiliary network known as a *gating network*, whose job it is to

<sup>1</sup> Our focus in this paper is the convergence of the EM algorithm to (local) maxima of the likelihood function for a fixed data set. We do not address the asymptotic convergence of these local maxima to "true" parameter values as the number of data points goes to infinity. For results on asymptotic convergence that are relevant to ME architectures, see Redner and Walker (1984).

partition the input space into regions corresponding to the various expert networks. This is done by assigning a probability vector  $[g_1, g_2, \dots, g_K]^T$  to each point in the input space. In particular, the gating net implements a parameterized function  $\mathbf{s}(\mathbf{x}, \boldsymbol{\theta}_0)$  and a normalizing function  $\mathbf{g}(\mathbf{s})$ . The function  $\mathbf{s}$  is a map from  $R^n$  to  $R^K$ , for each value of the parameter vector  $\boldsymbol{\theta}_0$ , and  $\mathbf{g}$  is a map from  $R^K$  to  $R^K$ . The particular form of  $\mathbf{g}$  used here is the softmax function

$$g_j = g_j(\mathbf{x}, \boldsymbol{\theta}_0) = \frac{e^{s_j(\mathbf{x}, \boldsymbol{\theta}_0)}}{\sum_{i=1}^K e^{s_i(\mathbf{x}, \boldsymbol{\theta}_0)}}, \quad j = 1, \dots, K. \quad (2)$$

Note that this definition implies that the  $g_j$  are positive and sum to one.

In the current paper we focus on the case in which the functions  $\mathbf{s}$  is linear (cf. Jordan & Jacobs, 1994). In this case the boundaries in the input space where  $g_i = g_{i'}$  are planar and the function  $\mathbf{g}$  can be viewed as a smoothed piecewise-planar partitioning of the input space.

Training data are assumed to be generated according to the following probability model. We assume that for a given  $\mathbf{x}$ , a label  $j$  is selected with probability  $P(j|\mathbf{x}) = g_j(\mathbf{x}, \boldsymbol{\theta}_0)$ . An output  $\mathbf{y}$  is then chosen with probability  $P(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}_j)$ . Thus, the total probability of observing  $\mathbf{y}$  from  $\mathbf{x}$  is given by the following finite mixture density

$$\begin{aligned} P(\mathbf{y}|\mathbf{x}) &= \sum_{j=1}^K P(j|\mathbf{x}) P(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}_j) \\ &= \sum_{j=1}^K g_j(\mathbf{x}, \boldsymbol{\theta}_0) P(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}_j). \end{aligned} \quad (3)$$

A training set  $\mathcal{Y} = \{(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}), t = 1, \dots, N\}$  is assumed to be generated as an independent set of draws from this mixture density. Thus, the total probability of the training set, for a specified set of input vectors  $\{\mathbf{x}^{(t)}\}_{t=1}^N$ , is given by the following likelihood function

$$\begin{aligned} L &= P(\{\mathbf{y}^{(t)}\}_1^N | \{\mathbf{x}^{(t)}\}_1^N) = \prod_{t=1}^N P(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}) \\ &= \prod_{t=1}^N \sum_{j=1}^K g_j(\mathbf{x}^{(t)}, \boldsymbol{\theta}_0) \\ &\quad \times P(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}, \boldsymbol{\theta}_j). \end{aligned} \quad (4)$$

The learning algorithms that we discuss are all maximum likelihood estimators. That is, we treat learning as the problem of finding parameters  $\boldsymbol{\theta}_0$ ,  $\boldsymbol{\theta}_j$ , and  $\Sigma_j$  to maximize  $L$ , or, more conveniently, to maximize the log likelihood  $l = \ln L$

$$l(\boldsymbol{\Theta}, \mathcal{Y}) = \sum_{t=1}^N \ln \sum_{j=1}^K g_j(\mathbf{x}^{(t)}, \boldsymbol{\theta}_0) P(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}, \boldsymbol{\theta}_j),$$

where  $\boldsymbol{\Theta} = [\boldsymbol{\theta}_0, \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_K, \Sigma_1, \Sigma_2, \dots, \Sigma_K]^T$ .

Given the probability model in eqn (3), the expected value of the output is given as follows

$$\boldsymbol{\mu} \equiv E[\mathbf{y}|\mathbf{x}] = \sum_{j=1}^K g_j(\mathbf{x}, \boldsymbol{\theta}_0) \boldsymbol{\mu}_j.$$

This motivates using the weighted output of the expert networks as the total output of the ME architecture (cf. Figure 1).

The model in eqns (3) and (1) is a finite Gaussian mixture model. It is interesting to compare this model to a related Gaussian mixture model that is widely studied in statistics; that is, the model

$$P(\mathbf{x}) = \sum_{j=1}^K \alpha_j P_j(\mathbf{x}|\boldsymbol{\theta}_j), \quad \alpha_j \geq 0, \quad \sum_{j=1}^K \alpha_j = 1. \quad (6)$$

The difference between these models is clear: the  $\alpha_j$ 's in eqn (6) are independent of the input vectors, where the  $g_j$ 's in eqn (3) are conditional on  $\mathbf{x}$  (they represent the probabilities  $P(j|\mathbf{x})$ ). Thus, model (6) represents an unconditional probability, appropriate for unsupervised learning, whereas model (3) represents a conditional probability, appropriate for supervised learning.

There is another model studied in statistics, the *switching regression* model (Quandt & Ramsey, 1972, 1978; De Veaux, 1986), that is intermediate between model (6) and model (3). The switching regression model is given as

$$P(\mathbf{y}|\mathbf{x}) = \lambda P(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}_1) + (1 - \lambda) P(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}_2), \quad (7)$$

where the  $P(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}_1)$  are univariate Gaussians and the mean of each Gaussian is assumed to be linear in  $\mathbf{x}$ . This model assumes that the data pair  $\{\mathbf{y}, \mathbf{x}\}$  is generated from a pair of linear regression models through a random switch that turns to one side with probability  $\lambda$  and to the other side with probability  $1 - \lambda$ . This model can be generalized to allow for a multinomial switch

$$P(\mathbf{y}|\mathbf{x}) = \sum_{j=1}^K \alpha_j P(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}_j), \quad (8)$$

where  $\alpha_j \geq 0$ ,  $\sum_{j=1}^K \alpha_j = 1$ , and  $P(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}_j)$  is given by eqn (1). The difference between switching regression and the ME model is that the switching regression model assumes that the setting of the switch is independent of the input vector. This

assumption does not allow for piecewise variation in the form of the regression surface; all of the regression components contribute throughout the input space. Switching regression can be viewed as one end of a continuum in which the overlap in the regression components is total; decision tree models (e.g., Breiman et al., 1984) are the other end of the continuum in which the overlap is zero. The ME model interpolates smoothly between these extremes.

## 2.2. An EM Algorithm for Training the Mixture of Experts

In many estimation problems the likelihood is a complicated nonlinear function of the parameters. Parameter estimation in such cases usually involves some sort of numerical optimization technique, typically gradient ascent. An alternative to gradient techniques, applicable in many situations, is the EM algorithm (Baum et al., 1970; Dempster et al., 1977). EM is based on the idea of solving a succession of simplified problems that are obtained by augmenting the original observed variables with a set of additional “hidden” variables. Unconditional mixture models are particularly amenable to the EM approach (Redner & Walker, 1984) and, as observed by Jordan and Jacobs (1994), the conditional mixture of experts model is also amenable to an EM approach.

Given an observed data set  $\mathcal{Y}$ , we augment  $\mathcal{Y}$  with a set of additional variables  $\mathcal{Y}_{\text{mis}}$ , called “missing” or “hidden” variables, and consider a maximum likelihood problem for a “complete-data” set  $\mathcal{Z} = \{\mathcal{Y}, \mathcal{Y}_{\text{mis}}\}$  (cf. Little & Rubin, 1987). We choose the missing variables in such a way that the resulting “complete-data log likelihood,” given by  $l_c(\Theta, \mathcal{Z}) = \ln P(\mathcal{Y}, \mathcal{Y}_{\text{mis}} | \Theta)$ , is easy to maximize with respect to  $\Theta$ . The probability model  $P(\mathcal{Y}, \mathcal{Y}_{\text{mis}} | \Theta)$  must be chosen so that its marginal distribution across  $\mathcal{Y}$ , referred to in this context as the “incomplete-data” likelihood, is the original likelihood

$$P(\mathcal{Y} | \Theta) = \int P(\mathcal{Y}, \mathcal{Y}_{\text{mis}} | \Theta) d\mathcal{Y}_{\text{mis}}. \quad (9)$$

In deriving an update to the parameters based on the complete data log likelihood, we first note that we cannot work directly with the complete data log likelihood, because this likelihood is a random function of the missing random variables  $\mathcal{Y}_{\text{mis}}$ . The idea is to average out  $\mathcal{Y}_{\text{mis}}$ , that is, to maximize the *expected* complete data log likelihood  $E_{\mathcal{Y}_{\text{mis}}}[\ln P(\mathcal{Y}, \mathcal{Y}_{\text{mis}} | \Theta)]$ . This idea motivates the EM algorithm.

The EM algorithm is an iterative algorithm consisting of two steps:

- The Expectation (E) step, which computes the following conditional expectation of the log likelihood

$$\begin{aligned} Q(\Theta | \Theta^{(k)}) &= E_{\mathcal{Y}_{\text{mis}}} \{ \ln P(\mathcal{Z} | \Theta) | \mathcal{Y}, \Theta^{(k)} \} \\ &= \int P(\mathcal{Y}_{\text{mis}} | \mathcal{Y}, \Theta^{(k)}) \ln P(\mathcal{Z} | \Theta) d\mathcal{Y}_{\text{mis}} \end{aligned} \quad (10)$$

where  $\Theta^{(k)}$  is the value of the parameter vector at iteration  $k$ .

- The Maximization (M) step, which computes

$$\Theta^{(k+1)} = \arg \max_{\Theta} Q(\Theta | \Theta^{(k)}). \quad (11)$$

The M step chooses a parameter value that increases the  $Q$  function; the expected value of the complete data log likelihood. Dempster et al. (1977) proved that an iteration of EM also increases the original log likelihood  $l$ . That is,

$$l(\Theta^{(k+1)}; \mathcal{Y}) \geq l(\Theta^{(k)}; \mathcal{Y}).$$

Thus, the likelihood  $l$  increases monotonically along the sequence of parameter estimates generated by an EM algorithm.

Although in many cases the solution to the M step can be obtained analytically, in other cases an iterative inner loop is required to optimize  $Q$ . Another possibility is to simply increase the value of  $Q$  during the M step

$$Q(\Theta^{(k+1)} | \Theta^{(k)}) > Q(\Theta^{(k)} | \Theta^{(k)}) \quad (12)$$

by some means, for example, by gradient ascent or by Newton's method. An algorithm with an M-step given by eqn (12) is referred to as a *generalized EM* (GEM) algorithm (Dempster et al., 1977).

For the ME architecture we choose the missing data to be a set of indicator random variables  $\mathcal{Y}_{\text{mis}} = \{I_j^{(t)}, j = 1, \dots, K, t = 1, \dots, N\}$  with

$$I_j^{(t)} = \begin{cases} 1 & \text{if } \mathbf{y}^{(t)} \text{ is generated from the } j^{\text{th}} \\ & \text{model given by eqn (1),} \\ 0 & \text{otherwise} \end{cases} \quad (13)$$

and

$$\sum_{j=1}^K I_j^{(t)} = 1, \text{ for each } t.$$

We assume that the distribution of the complete data is given as

$$P(\mathcal{Z}|\Theta) = \prod_{i=1}^N \prod_{j=1}^K [g_j(\mathbf{x}^{(i)}, \boldsymbol{\theta}_0) P(\mathbf{y}^{(i)}|\mathbf{x}^{(i)}, \boldsymbol{\theta}_j)]^{I_j^{(i)}}.$$

It is easy to verify that this distribution satisfies eqn (9).

From eqn (10), we also obtain

$$\begin{aligned} Q(\Theta|\Theta^{(k)}) &= E_{\mathcal{Y}_{\text{mu}}} \{\ln P(\mathcal{Z}|\Theta)|\mathcal{Y}, \Theta^{(k)}\} \\ &= \sum_{i=1}^N \sum_{j=1}^K h_j^{(k)}(t) \ln [g_j(\mathbf{x}^{(i)}, \boldsymbol{\theta}_0) P(\mathbf{y}^{(i)}|\mathbf{x}^{(i)}, \boldsymbol{\theta}_j)] \\ &= \sum_{i=1}^N \sum_{j=1}^K h_j^{(k)}(t) \ln g_j(\mathbf{x}^{(i)}) \\ &\quad + \sum_{i=1}^N h_1^{(k)}(t) \ln P(\mathbf{y}^{(i)}|\mathbf{x}^{(i)}, \boldsymbol{\theta}_1) \\ &\quad + \dots + \sum_{i=1}^N h_K^{(k)}(t) \ln P(\mathbf{y}^{(i)}|\mathbf{x}^{(i)}, \boldsymbol{\theta}_K), \end{aligned} \quad (14)$$

where

$$\begin{aligned} h_j^{(k)}(t) &= E[I_j^{(i)}|\mathcal{Y}, \Theta^{(k)}] = P(j|\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \\ &= \frac{g_j(\mathbf{x}^{(i)}, \boldsymbol{\theta}_0^{(k)}) P(\mathbf{y}^{(i)}|\mathbf{x}^{(i)}, \boldsymbol{\theta}_j^{(k)})}{\sum_{i=1}^K g_i(\mathbf{x}^{(i)}, \boldsymbol{\theta}_0^{(k)}) P(\mathbf{y}^{(i)}|\mathbf{x}^{(i)}, \boldsymbol{\theta}_i^{(k)})}, \end{aligned} \quad (15)$$

where  $P(j|\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$  denotes the probability that the pair  $\{\mathbf{x}^{(i)}, \mathbf{y}^{(i)}\}$  comes from the  $j$ th probability model. Note that we always have  $h_j^{(k)}(t) > 0$ .

With the  $Q$  function in hand, we now investigate the implementation of the M step. From eqns (14), (1) and (2), we have

$$\begin{aligned} \frac{\partial Q}{\partial \boldsymbol{\theta}_0} &= \sum_{i=1}^N \sum_{j=1}^K h_j^{(k)}(t) \frac{\partial g_j(\mathbf{x}^{(i)}, \boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}_0} / g_j(\mathbf{x}^{(i)}, \boldsymbol{\theta}_0) \\ &= \sum_{i=1}^N \sum_{j=1}^K h_j^{(k)}(t) \left[ \frac{\partial s_j}{\partial \boldsymbol{\theta}_0} - \sum_{i=1}^K g_i(\mathbf{x}^{(i)}, \boldsymbol{\theta}_0) \frac{\partial s_i}{\partial \boldsymbol{\theta}_0} \right] \\ &= \sum_{i=1}^N \sum_{j=1}^K [h_j^{(k)}(t) - g_j(\mathbf{x}^{(i)}, \boldsymbol{\theta}_0)] \frac{\partial s_j}{\partial \boldsymbol{\theta}_0}, \end{aligned} \quad (16)$$

$$\begin{aligned} \frac{\partial Q}{\partial \boldsymbol{\theta}_j} &= \sum_{i=1}^N h_j^{(k)}(t) \frac{\partial P(\mathbf{y}^{(i)}|\mathbf{x}^{(i)}, \boldsymbol{\theta}_j)}{\partial \boldsymbol{\theta}_j} / P(\mathbf{y}^{(i)}|\mathbf{x}^{(i)}, \boldsymbol{\theta}_j) \\ &= \sum_{i=1}^N h_j^{(k)}(t) \frac{\partial \mathbf{f}_j^T(\mathbf{x}^{(i)}, \boldsymbol{\theta}_j)}{\partial \boldsymbol{\theta}_j} \Sigma_j^{-1} [\mathbf{y}^{(i)} - \mathbf{f}_j(\mathbf{x}^{(i)}, \boldsymbol{\theta}_j)], \end{aligned} \quad (17)$$

and

$$\begin{aligned} \frac{\partial Q}{\partial \Sigma_j} &= \sum_{i=1}^N h_j^{(k)}(t) \frac{\partial P(\mathbf{y}^{(i)}|\mathbf{x}^{(i)}, \boldsymbol{\theta}_j)}{\partial \Sigma_j} / P(\mathbf{y}^{(i)}|\mathbf{x}^{(i)}, \boldsymbol{\theta}_j) \\ &= -\frac{1}{2} \sum_{i=1}^N h_j^{(k)}(t) \Sigma_j^{-1} \{ \Sigma_j - [\mathbf{y}^{(i)} - \mathbf{f}_j(\mathbf{x}^{(i)}, \boldsymbol{\theta}_j)] \\ &\quad \times [\mathbf{y}^{(i)} - \mathbf{f}_j(\mathbf{x}^{(i)}, \boldsymbol{\theta}_j)]^T \} \Sigma_j^{-1}. \end{aligned} \quad (18)$$

By letting

$$\left. \frac{\partial Q}{\partial \Sigma_j} \right|_{\Sigma_j = \Sigma_j^{(k+1)}} = 0,$$

we obtain the update for the covariance matrices

$$\begin{aligned} \Sigma_j^{(k+1)} &= \frac{1}{\sum_{i=1}^N h_j^{(k)}(t)} \sum_{i=1}^N h_j^{(k)}(t) [\mathbf{y}^{(i)} - \mathbf{f}_j(\mathbf{x}^{(i)}, \boldsymbol{\theta}_j)] \\ &\quad \times [\mathbf{y}^{(i)} - \mathbf{f}_j(\mathbf{x}^{(i)}, \boldsymbol{\theta}_j)]^T. \end{aligned} \quad (19)$$

Assuming that the training set  $\mathcal{Y}$  is generated by a mixture model, we note that when the sample number  $N$  is sufficiently large (relative to the dimension of  $\mathbf{y}$ ), the space spanned by the  $N$  vectors  $[\mathbf{y}^{(i)} - \mathbf{f}_j(\mathbf{x}^{(i)}, \boldsymbol{\theta}_j)]$  will be of full dimension with probability one. Recalling that  $h_j^{(k)}(t) > 0$  we observe that when the sample number  $N$  is sufficiently large the matrices  $\Sigma_j^{(k+1)}$  are therefore positive definite with probability one.

Next, by letting

$$\left. \frac{\partial Q}{\partial \boldsymbol{\theta}_j} \right|_{\boldsymbol{\theta}_j = \boldsymbol{\theta}_j^{(k+1)}} = 0,$$

we obtain

$$\begin{aligned} \sum_{i=1}^N h_j^{(k)}(t) \frac{\partial \mathbf{f}_j^T(\mathbf{x}^{(i)}, \boldsymbol{\theta}_j)}{\partial \boldsymbol{\theta}_j} \\ \times (\Sigma_j^{(k)})^{-1} [\mathbf{y}^{(i)} - \mathbf{f}_j(\mathbf{x}^{(i)}, \boldsymbol{\theta}_j)] = 0, \end{aligned} \quad (20)$$

which we can solve explicitly given our assumption that the expert networks are linear

$$\boldsymbol{\theta}_j^{(k+1)} = (R_j^{(k)})^{-1} \mathbf{c}_j^{(k)} \quad (21)$$

where

$$\mathbf{c}_j^{(k)} = \sum_{i=1}^N h_j^{(k)}(t) X_i (\Sigma_j^{(k)})^{-1} \mathbf{y}^{(i)}, \quad (22)$$

$$R_j^{(k)} = \sum_{i=1}^N h_j^{(k)}(t) X_i (\Sigma_j^{(k)})^{-1} X_i^T, \quad (23)$$

and

$$X_i^T = \left\{ \begin{array}{cccc|cccc} (\mathbf{x}^{(i)})^T & 0 & \dots & 0 & 1 & 0 & \dots & 0 \\ 0 & (\mathbf{x}^{(i)})^T & 0 & \dots & 0 & 0 & 1 & 0 \\ \vdots & \vdots & & \vdots & \vdots & & & \vdots \\ 0 & \dots & \dots & 0 & (\mathbf{x}^{(i)})^T & 0 & \dots & 1 \end{array} \right\}. \quad (24)$$

Note that  $R_j^{(k)}$  is invertible with probability one when the sample size  $N$  is sufficiently large.

Finally, let us consider the update for  $\theta_0$ . Jordan and Jacobs (1994) observed that the gating network is a specific form of a generalized linear model, in particular a *multinomial logit* model (cf. McCullagh & Nelder, 1984). Multinomial logit models can be fit efficiently with a variant of Newton's method known as *iteratively reweighted least squares* (IRLS). For the purposes of the current paper, we simply write the generic form of a Newton update and refer the reader to Jordan and Jacobs (1994) for further details on IRLS. Note also that Jordan and Jacobs (1994) assume that the inner loop of IRLS fitting runs to completion. In the current paper we address only the case in which a single IRLS step is taken in the inner loop. The form of this IRLS step is also generalized to allow a learning rate parameter.<sup>2</sup>

The update for the gating network parameters is obtained as follows. Denote the gradient vector at iteration  $k$  as

$$\mathbf{e}_g^{(k)} = \sum_{i=1}^N \sum_{j=1}^K [h_j^{(k)}(t) - g_j(\mathbf{x}^{(i)}, \theta_0^{(k)})] \frac{\partial s_j}{\partial \theta_0^{(k)}}, \quad (25)$$

and the Hessian matrix at iteration  $k$  as

$$R_g^{(k)} = \sum_{i=1}^N \sum_{j=1}^K g_j(\mathbf{x}^{(i)}, \theta_0^{(k)}) [1 - g_j(\mathbf{x}^{(i)}, \theta_0^{(k)})] \frac{\partial s_j}{\partial \theta_0} \frac{\partial s_j}{\partial \theta_0^T}. \quad (26)$$

Then the generalized IRLS update is given as follows

$$\theta_0^{(k+1)} = \theta_0^{(k)} + \gamma_g (R_g^{(k)})^{-1} \mathbf{e}_g^{(k)}, \quad (27)$$

where  $\gamma_g$  is a learning rate.

In summary, the parameter update for the model eqn (3) is given as follows.

#### Algorithm 1

1. (The E step): Compute the  $h_j^{(k)}(t)$ 's by eqn (15).
2. (The M step): Compute  $\Sigma_j^{(k+1)}$  by eqn (19),

compute  $\theta_0^{(k+1)}$  by eqn (27), and also compute  $\theta_j^{(k+1)}$ ,  $j = 1, \dots, K$  by eqn (21).

Before closing this section, let us return to the switching regression model [eqn (8)]. Following the same procedure as above, we obtain the following EM algorithm for switching regression.

#### Algorithm 2

1. (The E step): Compute the  $h_j^{(k)}(t)$ 's by

$$h_j^{(k)}(t) = \frac{\alpha_j^{(k)} P(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}, \theta_j^{(k)})}{\sum_{i=1}^K \alpha_i^{(k)} P(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}, \theta_i^{(k)})}. \quad (28)$$

2. (The M step): Compute  $\Sigma_j^{(k+1)}$  by eqn (19), and let

$$\alpha_j^{(k+1)} = \frac{1}{N} \sum_{t=1}^N h_j^{(k)}(t). \quad (29)$$

Obtain  $\theta_j^{(k+1)}$ ,  $j = 1, \dots, K$  in the same manner as for model eqn (3).

We see that the EM algorithm for switching regression is simpler, because the  $\alpha_j^{(k)}$ 's are not constrained through a common parameter  $\theta_0$  as in the ME model.

### 2.3. Theoretical Convergence Results

In this section we provide a number of convergence results for the algorithm presented in the previous section. We study both the convergence and the convergence rate of the algorithm. In the Appendix we extend these results to a number of related algorithms.

We begin with a convergence theorem that establishes a relationship between the EM algorithm and gradient ascent.

**THEOREM 1.** *For the model given by eqn (3) and the learning algorithm given by Algorithm 1, we have:*

$$\begin{aligned} \theta_0^{(k+1)} - \theta_0^{(k)} &= P_g^{(k)} \frac{\partial l}{\partial \theta_0} \bigg|_{\theta_0 = \theta_0^{(k)}}, \\ \theta_j^{(k+1)} - \theta_j^{(k)} &= P_j^{(k)} \frac{\partial l}{\partial \theta_j} \bigg|_{\theta_j = \theta_j^{(k)}}, \quad j = 1, \dots, K, \\ \text{vec}[\Sigma_j^{(k+1)}] - \text{vec}[\Sigma_j^{(k)}] &= P_{\Sigma_j}^{(k)} \frac{\partial l}{\partial \text{vec}[\Sigma_j]} \bigg|_{\Sigma_j = \Sigma_j^{(k)}}, \quad j = 1, \dots, K, \end{aligned} \quad (30)$$

where  $l = \ln L$  is given by eqns (5), (3) and (1), and "vec[ $A$ ]" denotes the vector obtained by stacking the column vectors of the matrix  $A$ .

<sup>2</sup> Thus, the algorithm that we analyze in this paper is, strictly speaking, a GEM algorithm.

Moreover, assuming that the training set  $\mathcal{Y}$  is generated by the mixture model of eqns (3) and (1) and assuming that the number  $N$  is sufficiently large, we have that  $P_g^{(k)}$  is a positive definite matrix, and  $P_{\Sigma_j}^{(k)}$ ,  $P_{\Sigma_j}^{(k)}$ ,  $j = 1, \dots, K$  are positive definite matrices with probability one. Specifically, they take the following values:

- (i)  $P_g^{(k)} = \gamma_g (R_g^{(k)})^{-1}$  with  $R_g^{(k)}$  given by eqn (26).
- (ii) For  $j = 1, \dots, K$ ,  $P_j^{(k)} = (R_j^{(k)})^{-1}$  with  $R_j^{(k)}$  is given by eqn (23).
- (iii) For  $j = 1, \dots, K$ ,

$$P_{\Sigma_j}^{(k)} = \frac{2}{\sum_{t=1}^N h_j^{(k)}(t)} \Sigma_j^{(k)} \otimes \Sigma_j^{(k)} \quad (31)$$

where “ $\otimes$ ” denotes the Kronecker product. For a  $m \times n$  matrix  $A$  and  $q \times m$  matrix  $B$ , the Kronecker product  $A \otimes B$  is defined as

$$A \otimes B = \begin{Bmatrix} a_{11}B & a_{12}B & \cdots & a_{1n}B \\ a_{21}B & a_{22}B & \cdots & a_{2n}B \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1}B & a_{m2}B & \cdots & a_{mn}B \end{Bmatrix}.$$

*Proof.* From eqns (5), (3), and (1), for  $l = \ln L$ , we obtain the following derivatives:

$$\begin{aligned} \left. \frac{\partial l}{\partial \theta_0} \right|_{\theta_0 = \theta_0^{(k)}} &= \sum_{t=1}^N \sum_{j=1}^K \\ &\times \left\{ \frac{g_j(\mathbf{x}^{(t)}, \theta_0^{(k)}) P(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}, \theta_j^{(k)})}{\sum_{i=1}^K g_i(\mathbf{x}^{(t)}, \theta_0^{(k)}) P(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}, \theta_i^{(k)})} \right\} \\ &\times \frac{\partial g_j(\mathbf{x}^{(t)}, \theta_0) / \partial \theta_0 |_{\theta_0 = \theta_0^{(k)}}}{g_j(\mathbf{x}^{(t)}, \theta_0^{(k)})} \\ &= \sum_{t=1}^N \sum_{j=1}^K [h_j^{(k)}(t) - g_j(\mathbf{x}^{(t)}, \theta_0^{(k)})] \frac{\partial s_j}{\partial \theta_0} \Big|_{\theta_0 = \theta_0^{(k)}}, \end{aligned} \quad (32)$$

$$\begin{aligned} \left. \frac{\partial l}{\partial \theta_j} \right|_{\theta_j = \theta_j^{(k)}} &= \sum_{t=1}^N \left\{ \frac{g_j(\mathbf{x}^{(t)}, \theta_0^{(k)}) P(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}, \theta_j^{(k)})}{\sum_{i=1}^K g_i(\mathbf{x}^{(t)}, \theta_0^{(k)}) P(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}, \theta_i^{(k)})} \right\} \\ &\times \frac{\partial P(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}, \theta_j) / \partial \theta_j |_{\theta_j = \theta_j^{(k)}}}{P(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}, \theta_j^{(k)})} \\ &= \sum_{t=1}^N h_j^{(k)}(t) \frac{\partial \mathbf{f}_j^T(\mathbf{x}^{(t)}, \theta_j)}{\partial \theta_j} \Big|_{\theta_j = \theta_j^{(k)}} \\ &\times (\Sigma_j^{(k)})^{-1} [\mathbf{y}^{(t)} - \mathbf{f}_j(\mathbf{x}^{(t)}, \theta_j^{(k)})], \end{aligned} \quad (33)$$

$$\begin{aligned} \left. \frac{\partial l}{\partial \Sigma_j} \right|_{\Sigma_j = \Sigma_j^{(k)}} &= \sum_{t=1}^N \left\{ \frac{g_j(\mathbf{x}^{(t)}, \theta_0^{(k)}) P(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}, \theta_j^{(k)})}{\sum_{i=1}^K g_i(\mathbf{x}^{(t)}, \theta_0^{(k)}) P(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}, \theta_i^{(k)})} \right\} \\ &\times \frac{\partial P(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}, \theta_j) / \partial \Sigma_j |_{\Sigma_j = \Sigma_j^{(k)}}}{P(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}, \theta_j^{(k)})} \\ &= -\frac{1}{2} \sum_{t=1}^N h_j^{(k)}(t) (\Sigma_j^{(k)})^{-1} \\ &\times \{ \Sigma_j^{(k)} - [\mathbf{y}^{(t)} - \mathbf{f}_j(\mathbf{x}^{(t)}, \theta_j^{(k)})] \\ &\times [\mathbf{y}^{(t)} - \mathbf{f}_j(\mathbf{x}^{(t)}, \theta_j^{(k)})]^T \} (\Sigma_j^{(k)}) \end{aligned} \quad (34)$$

We now prove points (i), (ii) and (iii).

(i) Comparing eqn (32) with eqn (27) it follows that  $P_g^{(k)} = \gamma_g (R_g^{(k)})^{-1}$ . To show that  $P_g^{(k)}$  is positive definite, we show that  $R_g^{(k)}$  is positive definite.<sup>3</sup> For an arbitrary vector  $\mathbf{u}$ , from eqn (26) we have

$$\begin{aligned} \mathbf{u}^T R_g^{(k)} \mathbf{u} &= \sum_{t=1}^N \sum_{j=1}^K g_j(\mathbf{x}^{(t)}, \theta_0^{(k)}) [1 - g_j(\mathbf{x}^{(t)}, \theta_0^{(k)})] \mathbf{u}^T \frac{\partial s_j}{\partial \theta_0} \frac{\partial s_j}{\partial \theta_0^T} \mathbf{u} \\ &= \sum_{t=1}^N \sum_{j=1}^K g_j(\mathbf{x}^{(t)}, \theta_0^{(k)}) [1 - g_j(\mathbf{x}^{(t)}, \theta_0^{(k)})] \mathbf{v}^T \mathbf{v} \geq 0, \end{aligned}$$

because  $g_j(\mathbf{x}^{(t)}, \theta_0^{(k)}) [1 - g_j(\mathbf{x}^{(t)}, \theta_0^{(k)})] > 0$ . Equality holds in the above equation only when  $\mathbf{v} = [\partial s_j / \partial \theta_0^T] \mathbf{u} = 0$  for any  $\mathbf{u}$ , which is impossible. Thus, we have established that  $R_g^{(k)}$  [and thus also  $(R_g^{(k)})^{-1}$ ] is positive definite.

(ii) Let  $\mathbf{c}_j^{(k)}$  and  $R_j^{(k)}$  be given by eqns (22) and (23). From eqn (33), we obtain

$$[\mathbf{c}_j^{(k)} - (R_j^{(k)}) \theta_j^{(k)}] = \left. \frac{\partial l}{\partial \theta_j} \right|_{\theta_j = \theta_j^{(k)}}.$$

Furthermore, it follows from eqn (21) that

$$\begin{aligned} \theta_j^{(k+1)} &= \theta_j^{(k)} + (R_j^{(k)})^{-1} \mathbf{c}_j^{(k)} - \theta_j^{(k)} \\ &= \theta_j^{(k)} + (R_j^{(k)})^{-1} [\mathbf{c}_j^{(k)} - (R_j^{(k)}) \theta_j^{(k)}]. \end{aligned} \quad (35)$$

That is, we have

$$\theta_j^{(k+1)} = \theta_j^{(k)} + (R_j^{(k)})^{-1} \left. \frac{\partial l}{\partial \theta_j} \right|_{\theta_j = \theta_j^{(k)}}$$

and  $P_j^{(k)} = (R_j^{(k)})^{-1}$ .

<sup>3</sup> A matrix  $A$  is positive definite iff  $A^{-1}$  is positive definite.

We now prove that  $R_j^{(k)}$  is positive definite. For an arbitrary vector  $\mathbf{u}$ , from eqn (21) we have

$$\begin{aligned}\mathbf{u}^T R_j^{(k)} \mathbf{u} &= \sum_{t=1}^N h_j^{(k)}(t) \mathbf{u}^T X_t (\Sigma_j^{(k)})^{-1} X_t^T \mathbf{u} \\ &= \sum_{t=1}^N h_j^{(k)}(t) \mathbf{v}^T (\Sigma_j^{(k)})^{-1} \mathbf{v} \geq 0\end{aligned}$$

where  $\mathbf{v} = X_t^T \mathbf{u}$ . From the note immediately following eqn (15), we know that  $\Sigma_j^{(k)}$  given by eqn (19) is positive definite and invertible for each  $k$  with probability one. Thus, with probability one, the equality of the above equation holds only when  $\mathbf{v} = 0$  for any  $\mathbf{u}$ , which is impossible. Thus, we have established that  $R_j^{(k)}$  [and thus also  $(R_j^{(k)})^{-1}$ ] is positive definite with probability one.

(iii) We consider eqn (19) for updating  $\Sigma_j^{(k)}$ . This equation can be expanded as

$$\begin{aligned}\Sigma_j^{(k+1)} &= \Sigma_j^{(k)} + \frac{1}{\sum_{t=1}^N h_j^{(k)}(t)} \sum_{t=1}^N h_j^{(k)}(t) [\mathbf{y}^{(t)} - \mathbf{f}_j(\mathbf{x}^{(t)}, \boldsymbol{\theta}_j)] \\ &\quad \times [\mathbf{y}^{(t)} - \mathbf{f}_j(\mathbf{x}^{(t)}, \boldsymbol{\theta}_j)]^T - \Sigma_j^{(k)} \\ &= \Sigma_j^{(k)} + \frac{2\Sigma_j^{(k)}}{\sum_{t=1}^N h_j^{(k)}(t)} V_{\Sigma_j} \Sigma_j^{(k)},\end{aligned}\quad (36)$$

where

$$\begin{aligned}V_{\Sigma_j} &= -\frac{1}{2} \sum_{t=1}^N h_j^{(k)}(t) (\Sigma_j^{(k)})^{-1} \{ \Sigma_j^{(k)} - [\mathbf{y}^{(t)} - \mathbf{f}_j(\mathbf{x}^{(t)}, \boldsymbol{\theta}_j^{(k)})] \\ &\quad \times [\mathbf{y}^{(t)} - \mathbf{f}_j(\mathbf{x}^{(t)}, \boldsymbol{\theta}_j^{(k)})]^T \} (\Sigma_j^{(k)})^{-1}.\end{aligned}$$

It follows from eqn (34); that

$$V_{\Sigma_j} = \left. \frac{\partial l}{\partial \Sigma_j} \right|_{\Sigma_j = \Sigma_j^{(k)}}.$$

That is, we have

$$\Sigma_j^{(k+1)} = \frac{2\Sigma_j^{(k)}}{\sum_{t=1}^N h_j^{(k)}(t)} \left. \frac{\partial l}{\partial \text{vec} \Sigma_j} \right|_{\Sigma_j = \Sigma_j^{(k)}} \Sigma_j^{(k)}.$$

Utilizing the identity  $\text{vec}[ABC] = (C^T \otimes A) \text{vec}[B]$ , we obtain

$$\text{vec}[\Sigma_j^{(k+1)}] = \frac{2}{\sum_{t=1}^N h_j^{(k)}(t)} (\Sigma_j^{(k)} \otimes \Sigma_j^{(k)}) \left. \frac{\partial l}{\partial \Sigma_j} \right|_{\Sigma_j = \Sigma_j^{(k)}}.$$

Thus,

$$P_{\Sigma_j}^{(k)} = \frac{2}{\sum_{t=1}^N h_j^{(k)}(t)} (\Sigma_j^{(k)} \otimes \Sigma_j^{(k)}).$$

Moreover, for an arbitrary matrix  $U$ , we have

$$\text{vec}[U]^T (\Sigma_j^{(k)} \otimes \Sigma_j^{(k)}) \text{vec}[U] = \text{tr}(\Sigma_j^{(k)} U \Sigma_j^{(k)} U^T)$$

$$= \text{tr}((\Sigma_j^{(k)} U)^T (\Sigma_j^{(k)} U)) = \text{vec}[\Sigma_j^{(k)} U]^T \text{vec}[\Sigma_j^{(k)} U] \geq 0$$

where the equality holds only when  $\Sigma_j^{(k)} U = 0$ , which is impossible with probability one because  $U$  is arbitrary, and  $\Sigma_j^{(k)}$  is, as indicated above, positive definite with probability one. Thus, we have established that  $P_{\Sigma_j}^{(k)}$  is positive definite with probability one. ■

Theorem 1 can be used to establish a relationship between the step taken by the EM algorithm and the direction of steepest ascent. Recall that for a positive matrix  $B$ , we have

$$\frac{\partial l^T}{\partial \boldsymbol{\theta}} B \frac{\partial l}{\partial \boldsymbol{\theta}} > 0.$$

This implies the following corollary.

**COROLLARY 1.** Assume that the training set  $\{(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}), t = 1, \dots, N\}$  comes from the mixture model of eqns (3) and (1) and that  $N$  is sufficiently large. With probability one, the search direction of the EM algorithm has a positive projection on the gradient ascent searching direction of  $l = \ln L$ .

That is, the EM algorithm can be viewed as a modified gradient ascent algorithm for maximizing  $l = \ln L$ . From Theorem 1,  $B$  changes with the iteration step  $k$ , thus, the EM algorithm can also be regarded as a variable metric gradient ascent algorithm. This algorithm searches in an uphill direction, so if the learning rate is appropriate, the searching process will converge to a local maximum or a saddle point of the likelihood  $l = \ln L$ .

Similar results have been obtained for unsupervised mixture models by Xu and Jordan (1993) and for Hidden Markov models by Baum and Sell (1968). See Xu and Jordan (1993) for further discussion of the relationships between these theorems.

We now establish the convergence of the parameters  $\boldsymbol{\Theta}^{(k)}$ . We also provide convergence rates for both  $l(\boldsymbol{\Theta}^{(k)})$  and  $\boldsymbol{\Theta}^{(k)}$ .

**THEOREM 2.** Assume that the training set  $\mathcal{Y}$  is generated by the mixture model of eqns (3) and (1)



and that  $N$  is sufficiently large. Assume further that  $\Sigma_1, \dots, \Sigma_K$  are diagonal and let  $\mathbf{v}_{\Sigma_j}$  be a vector consisting of the diagonal elements of  $\Sigma_j$ .

Let us denote

$$\Theta = [\theta_0^T, \theta_1^T, \dots, \theta_K^T, \mathbf{v}_{\Sigma_1}, \dots, \mathbf{v}_{\Sigma_K}]^T$$

$$P = \text{diag}[P_g^{(k)}, P_1, \dots, P_K, P_{\Sigma_1}, \dots, P_{\Sigma_K}],$$

and

$$H(\Theta) = \frac{\partial^2 l(\Theta)}{\partial \Theta \partial \Theta^T}.$$

Assume that on a given domain  $D_\Theta$

- (i)  $\partial^2 g_j(\theta_q)/\partial \theta_q \partial \theta_q^T, \partial^2 f_{ji}(\theta_j)/\partial \theta_j \partial \theta_j^T, j = 1, \dots, K, i = 1, \dots, m$  exist and are continuous;
- (ii) the Hessian matrix  $H(\Theta)$  is negative definite;
- (iii)  $\Theta^*$  is a local maximum of  $l(\Theta)$ , and  $\Theta^* \in D_\Theta$ .

Then with probability one,

(1) Letting  $-M, -m$  (here  $M > m > 0$ ) be the minimum and maximum eigenvalues of the negative definite matrix  $(P^{1/2})^T H(\Theta) (P^{1/2})$  (or equivalently the minimum and maximum eigenvalues of  $PH(\Theta)$ , because we have  $PH\mathbf{e} = \lambda\mathbf{e}$  from  $(P^{1/2})^T HP^{1/2}\mathbf{e} = \lambda\mathbf{e}$ ), we have

$$l(\Theta^*) - l(\Theta^{(k)}) \leq r^k [l(\Theta^*) - l(\Theta_0)], \quad (37)$$

$$\|P^{-1/2}(\Theta^{(k)} - \Theta^*)\| \leq |r|^{k/2} \sqrt{\frac{2}{m} [l(\Theta^*) - l(\Theta_0)]}, \quad (38)$$

where

$$r = 1 - \left(1 - \frac{M}{2}\right) \frac{m^2}{M} < 1.$$

We also have  $0 < |r| < 1$  when  $M < 2$ .

(2) For any initial point  $\Theta_0 \in D_\Theta$ ,  $\lim_{k \rightarrow \infty} \Theta^{(k)} = \Theta^*$  when  $M < 2$ .

*Proof.* As indicated earlier, when the training set  $\{(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}), t = 1, \dots, N\}$  is generated from the mixture model of eqns (3) and (1) and  $N$  is sufficiently large,  $\Sigma_j^{(k)}$  remains positive definite during the learning process. Thus, under the condition (i), it follows from eqns (5), (3) and (1) that  $H(\Theta)$  exists and remains continuous on  $D_\Theta$ . Expanding the log likelihood in a Taylor expansion, we have

$$l(\Theta) - l(\Theta^*) = (\Theta - \Theta^*)^T \frac{\partial l(\Theta)}{\partial \Theta} \bigg|_{\Theta=\Theta^*} + \frac{1}{2} (\Theta - \Theta^*)^T H(\Theta^* + \xi(\Theta - \Theta^*)) (\Theta - \Theta^*)$$

with  $0 < \xi < 1$ . Because

$$\frac{\partial l(\Theta)}{\partial \Theta} \bigg|_{\Theta=\Theta^*} = 0,$$

we have

$$l(\Theta) - l(\Theta^*) = \frac{1}{2} (\Theta - \Theta^*)^T H(\Theta^* + \xi(\Theta - \Theta^*)) (\Theta - \Theta^*). \quad (39)$$

From Theorem 1 we know that  $P$  is positive definite. Furthermore, from condition (ii),  $H(\Theta)$  is negative definite on  $D_\Theta$ . This implies that  $P^{1/2}$  exists and  $(P^{1/2})^T H(\Theta) (P^{1/2})$  is negative definite on  $D_\Theta$ . Utilizing the Rayleigh quotient we obtain that for any  $\mathbf{u}$

$$-M \|\mathbf{u}\|^2 \leq \mathbf{u}^T (P^{1/2})^T H(\Theta) (P^{1/2}) \mathbf{u} \leq -m \|\mathbf{u}\|^2. \quad (40)$$

Substituting eqn (40) into eqn (39), we obtain

$$l(\Theta) - l(\Theta^*) = \frac{1}{2} (\Theta - \Theta^*)^T (P^{-1/2})^T (P^{1/2})^T \times H(\Theta^* + \xi(\Theta - \Theta^*)) P^{1/2} P^{-1/2} (\Theta - \Theta^*) \quad (41)$$

$$l(\Theta) - l(\Theta^*) \geq -\frac{M}{2} \|P^{-1/2}(\Theta - \Theta^*)\|^2. \quad (42)$$

Moreover, we have

$$\begin{aligned} & -m \|P^{-1/2}(\Theta - \Theta^*)\|^2 \\ & \geq \left| (\Theta - \Theta^*)^T \left[ \frac{\partial l(\Theta)}{\partial \Theta} - \frac{\partial l(\Theta)}{\partial \Theta} \bigg|_{\Theta=\Theta^*} \right] \right| \\ & \geq -\left\| P^{1/2} \frac{\partial l(\Theta)}{\partial \Theta} \right\| \|P^{-1/2}(\Theta - \Theta^*)\|. \end{aligned}$$

Thus,

$$\|P^{-1/2}(\Theta - \Theta^*)\| \leq \frac{1}{m} \left\| P^{1/2} \frac{\partial l(\Theta)}{\partial \Theta} \right\|.$$

Together with eqn (42), we obtain

$$-\left\| P^{1/2} \frac{\partial l(\Theta)}{\partial \Theta} \right\|^2 \leq \frac{2m^2}{M} [l(\Theta) - l(\Theta^*)]. \quad (43)$$

On the other hand, we also have

$$l(\Theta) - l(\Theta^{(k)}) = (\Theta - \Theta^{(k)})^T \frac{\partial l(\Theta)}{\partial \Theta} \bigg|_{\Theta=\Theta^{(k)}} + \frac{1}{2} (\Theta - \Theta^{(k)})^T H(\Theta^{(k)} + \xi'(\Theta - \Theta^{(k)})) (\Theta - \Theta^{(k)})$$

with  $0 < \xi' < 1$ . By Theorem 1, we know that for the EM algorithm,

$$\Theta^{(k+1)} - \Theta^{(k)} = P \frac{\partial l(\Theta)}{\partial \Theta} \bigg|_{\Theta=\Theta^{(k)}}$$

Utilizing this result in the above equation, we obtain

$$\begin{aligned} l(\Theta^{(k+1)}) - l(\Theta^{(k)}) &= \left\| P^{1/2} \frac{\partial l(\Theta)}{\partial \Theta} \bigg|_{\Theta=\Theta^{(k)}} \right\|^2 \\ &\quad + \frac{1}{2} \left( P^{1/2} \frac{\partial l(\Theta)}{\partial \Theta} \bigg|_{\Theta=\Theta^{(k)}} \right)^T (P^{1/2})^T \\ &\quad \times H(\Theta^{(k)} + \xi'(\Theta - \Theta^{(k)})) P^{1/2} \\ &\quad \times \left( P^{1/2} \frac{\partial l(\Theta)}{\partial \Theta} \bigg|_{\Theta=\Theta^{(k)}} \right) \\ &\geq \left( 1 - \frac{M}{2} \right) \left\| P^{1/2} \frac{\partial l(\Theta)}{\partial \Theta} \bigg|_{\Theta=\Theta^{(k)}} \right\|^2. \end{aligned} \quad (44)$$

Combining eqn (44) and eqn (43), we obtain

$$l(\Theta^{(k+1)}) - l(\Theta^{(k)}) \geq - \left( 1 - \frac{M}{2} \right) \frac{2m^2}{M} [l(\Theta^{(k)}) - l(\Theta^*)]$$

and furthermore

$$\begin{aligned} l(\Theta^{(k+1)}) - l(\Theta^*) &\geq \left[ 1 - \left( 1 - \frac{M}{2} \right) \frac{2m^2}{M} \right] [l(\Theta^{(k)}) - l(\Theta^*)] \\ &\geq \left[ 1 - \left( 1 - \frac{M}{2} \right) \frac{2m^{2^k}}{M} \right] [l(\Theta_0) - l(\Theta^*)]. \end{aligned} \quad (45)$$

Let  $r = [1 - (1 - M/2) 2m^2/M]$ . Multiplying both sides of the above equation by negative one, we obtain eqn (37). In addition, it is easy to verify that  $0 < |r| < 1$  when  $M < 2$  (recall that  $M > m$ ). Furthermore, it follows from eqn (40) and eqn (41) that we have

$$l(\Theta^{(k)}) - l(\Theta^*) \leq -\frac{m}{2} \|P^{-1/2}(\Theta^{(k)} - \Theta^*)\|^2$$

which, by eqn (45), becomes

$$-\frac{m}{2} \|P^{-1/2}(\Theta^{(k)} - \Theta^*)\|^2 \geq r^k [l(\Theta_0) - l(\Theta^*)]$$

$$\|P^{-1/2}(\Theta^{(k)} - \Theta^*)\| \leq |r|^{k/2} \sqrt{\frac{2}{m} [l(\Theta^*) - l(\Theta_0)]}$$

which is just eqn (38). In addition, when  $M < 2$ ,  $|r| < 1$ , we have  $\lim_{k \rightarrow \infty} \Theta^{(k)} = \Theta^*$  because  $P$  is positive definite. ■

We see from this theorem that the EM algorithm converges linearly. Moreover, the speed of convergence depends on the difference between  $M$  and  $m$ : the smaller the difference, the faster the convergence.

### 3. THEORETICAL ANALYSIS OF AN EM ALGORITHM FOR THE HIERARCHICAL MIXTURE OF EXPERTS ARCHITECTURE

#### 3.1. An EM Algorithm for Training the Hierarchical Architecture

The ME architecture can be viewed as an architecture for splitting the input space into regions in which different local functions are fit. The hierarchical mixture of experts (HME) architecture generalizes this idea to a nested model in which regions in the input space are split recursively into subregions (Jordan & Jacobs, 1992). The resulting tree-structured architecture can be viewed as a multi-resolution function approximator in which smoothed piecewise functions are fit at a variety of levels of resolution.

As shown in Figure 2, the HME architecture is a tree. In this tree, each terminal node is an expert network, and each nonterminal node is a root of a subtree that itself corresponds to an HME architecture. At every nonterminal node in the tree there is a gating network that is responsible for the topmost split of the HME architecture rooted at that node. All of the expert networks and the gating networks in the architecture have the same input vector  $x \in R^n$ . In the remainder of this section, as in Jordan and Jacobs (1994), we consider the case in which the expert

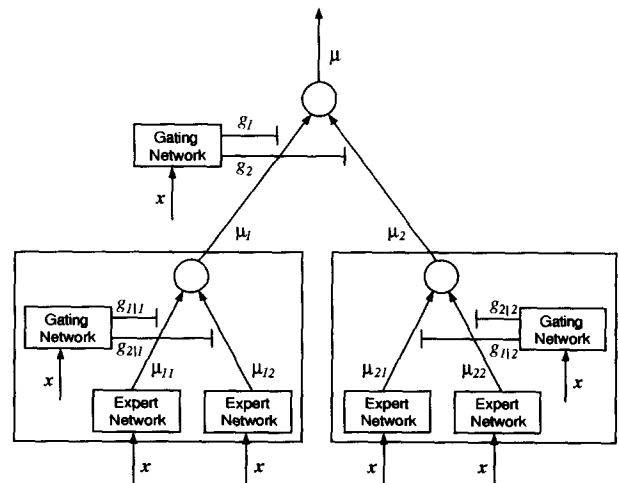


FIGURE 2. A two-level hierarchical mixture of experts. To form a deeper tree, each expert is expanded recursively into a gating network and a set of subexperts.

networks and the gating networks are generalized linear models. Furthermore, for simplicity, we consider only the case in which the probability model for the experts is Gaussian.

Let us denote a node at depth  $r$  by  $v_{i_0 i_1 \dots i_r}$ . This node is the  $i_r$ th daughter of the node  $v_{i_0 i_1 \dots i_{r-1}}$ . The root node of the tree is  $v_{i_0}$ . The number of branches emitted from  $v_{i_0 \dots i_r}$  is denoted by  $K_{i_0 i_1 \dots i_r}$ . For simplicity, we can omit  $i_0$  and write  $v_{i_1 \dots i_r}$  and  $K_{i_1 \dots i_r}$ . In addition, the output of the subtree rooted at  $v_{i_1 \dots i_r}$  is denoted

$$y_{i_1 \dots i_r} = \sum_{i_{r+1}=1}^{K_{i_1 \dots i_r}} g_{i_1 \dots i_r i_{r+1}}(\mathbf{x}, \boldsymbol{\theta}_{i_1 \dots i_r}^g) y_{i_1 \dots i_r i_{r+1}}$$

where  $g_{i_1 \dots i_r i_{r+1}}$  is the gating coefficient generated by the gating network attached at  $v_{i_1 \dots i_r}$ . This coefficient satisfies

$$\sum_{i_{r+1}=1}^{K_{i_1 \dots i_r}} g_{i_1 \dots i_r i_{r+1}}(\mathbf{x}, \boldsymbol{\theta}_{i_1 \dots i_r}^g) = 1,$$

for any  $\mathbf{x}$ , where  $\boldsymbol{\theta}_{i_1 \dots i_r}^g$  is the parameter vector of the gating network.

Given a training set  $\mathcal{Y} = \{(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}), t = 1, \dots, N\}$ , we want to maximize the likelihood function [cf. eqn (5)], that is,

$$L = P(\{\mathbf{y}^{(t)}\}_1^N | \{\mathbf{x}^{(t)}\}_1^N) = \prod_{t=1}^N P(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}).$$

Expanding the probability model, we have

$$P(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}) = \sum_{i_1=1}^{K_0} g_{i_1}(\mathbf{x}^{(t)}, \boldsymbol{\theta}_{i_0}^g) P(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}, v_{i_1})$$

where

$$P(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}, v_{i_1}) = \begin{cases} \sum_{i_2=1}^{K_{i_1}} g_{i_1 i_2}(\mathbf{x}^{(t)}, \boldsymbol{\theta}_{i_1}^g) P(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}, v_{i_1 i_2}), & \text{if } v_{i_1} \text{ is a nonterminal node,} \\ \frac{1}{(2\pi)^{m/2} |\Sigma_{i_1}|^{1/2}} \times e^{\{-1/2[\mathbf{y}^{(t)} - \mathbf{f}_{i_1}(\mathbf{x}^{(t)}, \boldsymbol{\theta}_{i_1})]^T \Sigma_{i_1}^{-1} [\mathbf{y}^{(t)} - \mathbf{f}_{i_1}(\mathbf{x}^{(t)}, \boldsymbol{\theta}_{i_1})]\}}, & \text{if } v_{i_1} \text{ is a terminal node,} \end{cases} \quad (46)$$

and recursively

$$P(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}, v_{i_1 \dots i_r}) = \begin{cases} \sum_{i_{r+1}=1}^{K_{i_1 \dots i_r}} g_{i_1 \dots i_r i_{r+1}}(\mathbf{x}, \boldsymbol{\theta}_{i_1 \dots i_r}^g) P(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}, v_{i_1 \dots i_r i_{r+1}}), & \text{if } v_{i_1 \dots i_r} \text{ is a nonterminal node,} \\ \frac{1}{(2\pi)^{m/2} |\Sigma_{i_1 \dots i_r}|^{1/2}} \times e^{\{-1/2[\mathbf{y}^{(t)} - \mathbf{f}_{i_1 \dots i_r}(\mathbf{x}^{(t)}, \boldsymbol{\theta}_{i_1 \dots i_r})]^T \Sigma_{i_1 \dots i_r}^{-1} [\mathbf{y}^{(t)} - \mathbf{f}_{i_1 \dots i_r}(\mathbf{x}^{(t)}, \boldsymbol{\theta}_{i_1 \dots i_r})]\}}, & \text{if } v_{i_1 \dots i_r} \text{ is a terminal node} \end{cases} \quad (47)$$

where  $\mathbf{f}_{i_1 \dots i_r}(\cdot, \cdot)$  is the parameterized function implemented by the expert network at terminal node  $v_{i_1 \dots i_r}$ , and  $\boldsymbol{\theta}_{i_1 \dots i_r}$ ,  $\Sigma_{i_1 \dots i_r}$  are the parameters of this expert network;  $P(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}, v_{i_1 \dots i_r})$  is the probability that  $\mathbf{y}^{(t)}$  is generated from the probability model rooted at  $v_{i_1 \dots i_r}$  when  $\mathbf{x}^{(t)}$  is the input.

To derive an EM algorithm for the HME architecture, we attach a set of indicator random variables to each nonterminal node  $v_{i_1 \dots i_r}$

$$I_{i_1 \dots i_r}^{(t)} = \begin{cases} 1, & \text{if } \mathbf{y}^{(t)} \text{ is generated by the subtree rooted at } v_{i_1 \dots i_r}, \\ 0, & \text{otherwise.} \end{cases}$$

The missing data  $\mathcal{Y}_{\text{mis}}$  consists of all of the indicator variables attached to the nonterminal nodes throughout the tree. In addition, we denote by  $\mathcal{Y}_{i_1 \dots i_r}^{\text{mis}}$  the set consisting of the indicator variables attached to the nonterminal nodes in the subtree rooted at  $v_{i_1 \dots i_r}$ .

We define the distribution of the complete data  $\mathcal{Z} = \{\mathcal{Y}, \mathcal{Y}_{\text{mis}}\}$  as

$$P(\mathcal{Z} | \boldsymbol{\Theta}) = \prod_{t=1}^N \prod_{i_1=1}^{K_0} [g_{i_1}(\mathbf{x}^{(t)}, \boldsymbol{\theta}_{i_0}^g) P(\mathcal{Z}_{i_1} | \mathbf{x}^{(t)}, v_{i_1})]^{I_{i_1}^{(t)}} \quad (48)$$

where

$$P(\mathcal{Z}_{i_1} | \mathbf{x}^{(t)}, v_{i_1}) = \begin{cases} \prod_{i_2=1}^{K_{i_1}} g_{i_1 i_2}(\mathbf{x}^{(t)}, \boldsymbol{\theta}_{i_1}^g) P(\mathcal{Z}_{i_1 i_2} | \mathbf{x}^{(t)}, v_{i_1 i_2})]^{I_{i_1 i_2}^{(t)}}, & \text{if } v_{i_1} \text{ is a nonterminal node,} \\ \frac{1}{(2\pi)^{m/2} |\Sigma_{i_1}|^{1/2}} \times e^{\{-1/2[\mathbf{y}^{(t)} - \mathbf{f}_{i_1}(\mathbf{x}^{(t)}, \boldsymbol{\theta}_{i_1})]^T \Sigma_{i_1}^{-1} [\mathbf{y}^{(t)} - \mathbf{f}_{i_1}(\mathbf{x}^{(t)}, \boldsymbol{\theta}_{i_1})]\}}, & \text{if } v_{i_1} \text{ is a terminal node} \end{cases} \quad (49)$$

and recursively

$$P(\mathcal{Z}_{i_1 \dots i_r} | \mathbf{x}^{(t)}, v_{i_1 \dots i_r}) = \begin{cases} \prod_{i_{r+1}=1}^{K_{i_1 \dots i_r}} [g_{i_1 \dots i_r+1}(\mathbf{x}, \boldsymbol{\theta}_{i_1 \dots i_r}^g)]^{I_{i_1 \dots i_r+1}^{(t)}} \\ \times P(\mathcal{Z}_{i_1 \dots i_r+1} | \mathbf{x}^{(t)}, v_{i_1 \dots i_r+1})]^{I_{i_1 \dots i_r+1}^{(t)}}, \\ \text{if } v_{i_1 \dots i_r} \text{ is a nonterminal node,} \\ \\ 1 \\ (2\pi)^{m/2} |\Sigma_{i_1 \dots i_r}|^{1/2} \\ \times e^{\{-1/2[\mathbf{y}^{(t)} - \mathbf{f}_{i_1 \dots i_r}(\mathbf{x}^{(t)}, \boldsymbol{\theta}_{i_1 \dots i_r})]^T \Sigma_{i_1 \dots i_r}^{-1} [\mathbf{y}^{(t)} - \mathbf{f}_{i_1 \dots i_r}(\mathbf{x}^{(t)}, \boldsymbol{\theta}_{i_1 \dots i_r})]\}}, \\ \text{if } v_{i_1 \dots i_r} \text{ is a terminal node.} \end{cases} \quad (50)$$

It is not difficult to verify that this distribution satisfies eqn (9) as required.

We now compute the  $Q$  function as required by the E step of the EM algorithm. From eqn (10), we obtain

$$Q(\boldsymbol{\Theta} | \boldsymbol{\Theta}^{(k)}) = \sum_{t=1}^N \sum_{i_1=1}^{K_0} h_{i_1}^{(k)}(t) \left\{ \ln[g_{i_1}(\mathbf{x}^{(t)}, \boldsymbol{\theta}_{i_1}^g)] + \ln F_{i_1} \right\}$$

where

$$h_{i_1}^{(k)}(t) = E[I_{i_1}^{(t)} | \mathcal{Y}, \boldsymbol{\Theta}^{(k)}] = \frac{g_{i_1}(\mathbf{x}^{(t)}, \boldsymbol{\theta}_{i_1}^{g(k)}) P^{(k)}(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}, v_{i_1})}{\sum_{i_1=1}^{K_0} g_{i_1}(\mathbf{x}^{(t)}, \boldsymbol{\theta}_{i_1}^{g(k)}) P^{(k)}(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}, v_{i_1})}, \quad (51)$$

$$\ln F_{i_1} = \begin{cases} \sum_{i_2=1}^{K_{i_1}} h_{i_1 i_2}^{(k)}(t) \{ \ln[g_{i_1 i_2}(\mathbf{x}^{(t)}, \boldsymbol{\theta}_{i_1 i_2}^g)] + \ln F_{i_1 i_2} \}, \\ \text{if } v_{i_1} \text{ is a nonterminal node,} \\ \\ \ln P(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}, v_{i_1}), \\ \text{if } v_{i_1} \text{ is a terminal node} \end{cases}$$

$$h_{i_1 i_2}^{(k)}(t) = E[I_{i_1 i_2}^{(t)} | \mathcal{Y}, \boldsymbol{\Theta}^{(k)}, I_{i_1}^{(t)} = 1] = \frac{g_{i_1 i_2}(\mathbf{x}^{(t)}, \boldsymbol{\theta}_{i_1 i_2}^{g(k)}) P^{(k)}(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}, v_{i_1 i_2})}{\sum_{i_2=1}^{K_{i_1}} g_{i_1 i_2}(\mathbf{x}^{(t)}, \boldsymbol{\theta}_{i_1 i_2}^{g(k)}) P^{(k)}(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}, v_{i_1 i_2})}, \quad (52)$$

and recursively

$$\ln F_{i_1 \dots i_r} = \begin{cases} \sum_{i_{r+1}=1}^{K_{i_1 \dots i_r}} h_{i_1 \dots i_r i_{r+1}}^{(k)}(t) \{ \ln[g_{i_1 \dots i_r i_{r+1}}(\mathbf{x}, \boldsymbol{\theta}_{i_1 \dots i_r i_{r+1}}^g)] \\ + \ln F_{i_1 \dots i_r i_{r+1}} \}, \\ \text{if } v_{i_1 \dots i_r} \text{ is a nonterminal node,} \\ \\ \ln P(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}, v_{i_1 \dots i_r}), \\ \text{if } v_{i_1 \dots i_r} \text{ is a terminal node} \end{cases}$$

$$h_{i_1 \dots i_r i_{r+1}}^{(k)}(t) = E[I_{i_1 \dots i_r i_{r+1}}^{(t)} | \mathcal{Y}, \boldsymbol{\Theta}^{(k)}, I_{i_1 i_2}^{(t)} = 1, \dots, I_{i_1 \dots i_r}^{(t)} = 1] = \frac{g_{i_1 \dots i_r i_{r+1}}(\mathbf{x}^{(t)}, \boldsymbol{\theta}_{i_1 \dots i_r i_{r+1}}^{g(k)}) P^{(k)}(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}, v_{i_1 \dots i_r i_{r+1}})}{\sum_{i_{r+1}=1}^{K_{i_1 \dots i_r}} g_{i_1 \dots i_r i_{r+1}}(\mathbf{x}^{(t)}, \boldsymbol{\theta}_{i_1 \dots i_r i_{r+1}}^{g(k)}) P^{(k)}(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}, v_{i_1 \dots i_r i_{r+1}})}, \quad (53)$$

where  $\boldsymbol{\theta}_{i_1 \dots i_r}^{g(k)}$  is the estimate  $\boldsymbol{\theta}_{i_1 \dots i_r}^g$  at iteration  $k$ .  $P(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}, v_{i_1 \dots i_r})$  is given by eqns (46) and (47) and  $P^{(k)}(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}, v_{i_1 \dots i_r})$  means that the probability is determined with all the parameters in the subtree rooted at  $v_{i_1 \dots i_r}$  being fixed at the estimates obtained in iteration  $k$ .

Proceeding now to the M step of the EM algorithm, we obtain parameter updates by optimizing the  $Q$  function. If the node  $v_{i_1 \dots i_r}$  is a terminal node, by setting the partial derivative of  $Q$  with respect to  $\Sigma_{i_1 \dots i_r}$  equal to zero, we obtain an update for the covariance matrices

$$\Sigma_{i_1 \dots i_r}^{(k+1)} = \sum_{t=1}^N h_{i_1}^{(k)}(t) h_{i_1 i_2}^{(k)}(t) \dots h_{i_1 \dots i_r}^{(k)}(t) \times [\mathbf{y}^{(t)} - \mathbf{f}_{i_1 \dots i_r}(\mathbf{x}^{(t)}, \boldsymbol{\theta}_{i_1 \dots i_r})] \times \frac{[\mathbf{y}^{(t)} - \mathbf{f}_{i_1 \dots i_r}(\mathbf{x}^{(t)}, \boldsymbol{\theta}_{i_1 \dots i_r})]^T}{\sum_{t=1}^N h_{i_1}^{(k)}(t) h_{i_1 i_2}^{(k)}(t) \dots h_{i_1 \dots i_r}^{(k)}(t)} \quad (54)$$

To obtain an update for the parameters of the expert networks we differentiate  $Q$  with respect to  $\boldsymbol{\theta}_{i_1 \dots i_r}$  and find that we must solve the following equation

$$\sum_{t=1}^N h_{i_1}^{(k)}(t) h_{i_1 i_2}^{(k)}(t) \dots h_{i_1 \dots i_r}^{(k)}(t) \frac{\partial \mathbf{f}_{i_1 \dots i_r}^T(\mathbf{x}^{(t)}, \boldsymbol{\theta}_{i_1 \dots i_r})}{\partial \boldsymbol{\theta}_{i_1 \dots i_r}} \times \Sigma_{i_1 \dots i_r}^{-1} [\mathbf{y}^{(t)} - \mathbf{f}_{i_1 \dots i_r}(\mathbf{x}^{(t)}, \boldsymbol{\theta}_{i_1 \dots i_r})] = 0. \quad (55)$$

In the case of linear expert networks, this equation is a weighted least squares equation, which can be solved as

$$\boldsymbol{\theta}_{i_1 \dots i_r}^{(k+1)} = (\mathbf{R}_{i_1 \dots i_r}^{(k)})^{-1} \mathbf{c}_{i_1 \dots i_r}^{(k)},$$

where

$$\mathbf{c}_{i_1 \dots i_r}^{(k)} = \sum_{t=1}^N h_{i_1}^{(k)}(t) h_{i_1 i_2}^{(k)}(t) \dots h_{i_1 \dots i_r}^{(k)}(t) X_t (\Sigma_{i_1 \dots i_r}^{(k)})^{-1} \mathbf{y}^{(t)},$$

and

$$R_{i_1 \dots i_r}^{(k)} = \sum_{t=1}^N h_{i_1}^{(k)}(t) h_{i_2}^{(k)}(t) \dots h_{i_r}^{(k)}(t) X_t (\Sigma_{i_1 \dots i_r}^{(k)})^{-1} X_t^T. \quad (56)$$

Finally, for any nonterminal node  $v_{i_1 \dots i_r}$ , setting the partial derivative of  $Q$  with respect to  $\theta_{i_1 \dots i_r}^{g(k+1)}$  equal to zero, we have that  $\theta_{i_1 \dots i_r}^{g(k+1)}$  is the solution of the following nonlinear system

$$\sum_{t=1}^N h_{i_1}^{(k)}(t) h_{i_2}^{(k)}(t) \dots h_{i_r}^{(k)}(t) \sum_{i_{r+1}=1}^{K_{i_1 \dots i_r}} [h_{i_1 \dots i_{r+1}}^{(k)}(t) - g_{i_1 \dots i_{r+1}}(\mathbf{x}^{(t)}, \theta_{i_1 \dots i_r}^g)] \frac{\partial s_{i_1 \dots i_{r+1}}}{\partial \theta_{i_1 \dots i_r}^g} = 0. \quad (57)$$

As in the case of the one-level ME architecture [cf. eqn (25) and eqn (27)], we obtain the following Newton step for updating the gating network parameters

$$\theta_{i_1 \dots i_r}^{g(k+1)} = \theta_{i_1 \dots i_r}^{g(k)} + \gamma_{i_1 \dots i_r} (R_{i_1 \dots i_r}^{g(k)})^{-1} e_{i_1 \dots i_r}^{g(k)}, \quad (58)$$

where

$$\begin{aligned} R_{i_1 \dots i_r}^{g(k)} &= \sum_{t=1}^N h_{i_1}^{(k)}(t) h_{i_2}^{(k)}(t) \dots h_{i_r}^{(k)}(t) \\ &\times \sum_{i_{r+1}=1}^{K_{i_1 \dots i_r}} g_{i_1 \dots i_{r+1}}(\mathbf{x}^{(t)}, \theta_{i_1 \dots i_r}^g) \\ &\times [1 - g_{i_1 \dots i_r}(\mathbf{x}^{(t)}, \theta_{i_1 \dots i_r}^g)] \frac{\partial s_{i_1 \dots i_{r+1}}}{\partial \theta_{i_1 \dots i_r}^{g(k)}} \frac{\partial s_{i_1 \dots i_{r+1}}}{\partial \theta_{i_1 \dots i_r}^{g(k)T}} \end{aligned} \quad (59)$$

and

$$\begin{aligned} e_{i_1 \dots i_r}^{g(k)} &= \sum_{t=1}^N h_{i_1}^{(k)}(t) h_{i_2}^{(k)}(t) \dots h_{i_r}^{(k)}(t) \\ &\times \sum_{i_{r+1}=1}^{K_{i_1 \dots i_r}} [h_{i_1 \dots i_{r+1}}^{(k)}(t) - g_{i_1 \dots i_{r+1}}(\mathbf{x}^{(t)}, \theta_{i_1 \dots i_r}^g)] \\ &\times \frac{\partial s_{i_1 \dots i_{r+1}}}{\partial \theta_{i_1 \dots i_r}^{g(k)}}. \end{aligned} \quad (60)$$

As in the case of the one-level architecture [cf. eqn (27)] the algorithm in eqn (58) is essentially the same as the IRLS algorithm suggested by Jordan and Jacobs (1994), although we have introduced a learning rate parameter and we restrict the update to a single step.

In summary, the EM algorithm for the HME architecture is given as follows.

### Algorithm 3

1. (The E step): Compute the  $h_{i_1}^{(k)}(t)$ ,  $h_{i_2}^{(k)}(t), \dots, h_{i_r}^{(k)}(t)$  by eqns (51), (52) and (53).

These computations can be performed via a single recursive pass upward in the tree.

2. (The M step): Compute the  $\Sigma_{i_1 \dots i_r}^{(k+1)}$  by eqn (54), compute the  $\theta_{i_1 \dots i_r}^{g(k+1)}$  by eqn (58), and compute the  $\theta_{i_1 \dots i_r}^{(k+1)}$  by eqn (56).

We can think of the E step as assigning credit (posterior probability) to various branches of the tree for each data point and the M step as solving weighted least squares problems in which the weights are given by the posterior probabilities assigned in the E step. The updates for the gating networks simply cache away the posteriors.

### 3.2. Theoretical Convergence Results

Much of the analysis developed in Section 2.2 can be extended to cover the EM algorithm for the HME architecture. In this subsection we extend Theorem 1 and Theorem 2 to cover the hierarchical case. The results are given as Theorems 3 and 4, respectively.

We first compute the derivatives of the  $Q$  function

$$\begin{aligned} \frac{\partial Q}{\partial \theta_{i_1 \dots i_r}^g} &= \sum_{t=1}^N h_{i_1}^{(k)}(t) h_{i_2}^{(k)}(t) \dots h_{i_r}^{(k)}(t) \\ &\times \sum_{i_{r+1}=1}^{K_{i_1 \dots i_r}} [h_{i_1 \dots i_{r+1}}^{(k)}(t) - g_{i_1 \dots i_{r+1}}(\mathbf{x}^{(t)}, \theta_{i_1 \dots i_r}^g)] \\ &\times \frac{\partial s_{i_1 \dots i_{r+1}}}{\partial \theta_{i_1 \dots i_r}^g}, \end{aligned} \quad (61)$$

$$\begin{aligned} \frac{\partial Q}{\partial \theta_{i_1 \dots i_r}} &= \sum_{t=1}^N h_{i_1}^{(k)}(t) h_{i_2}^{(k)}(t) \dots h_{i_r}^{(k)}(t) \\ &\times \frac{\partial \mathbf{f}_{i_1 \dots i_r}^T(\mathbf{x}^{(t)}, \theta_{i_1 \dots i_r})}{\partial \theta_{i_1 \dots i_r}} \\ &\times \Sigma_{i_1 \dots i_r}^{-1} [\mathbf{y}^{(t)} - \mathbf{f}_{i_1 \dots i_r}(\mathbf{x}^{(t)}, \theta_{i_1 \dots i_r})], \end{aligned} \quad (62)$$

$$\begin{aligned} \frac{\partial Q}{\partial \Sigma_{i_1 \dots i_r}} &= -\frac{1}{2} \sum_{t=1}^N h_{i_1}^{(k)}(t) h_{i_2}^{(k)}(t) \dots h_{i_r}^{(k)}(t) \Sigma_{i_1 \dots i_r}^{-1} \\ &\times \{\Sigma_{i_1 \dots i_r} - [\mathbf{y}^{(t)} - \mathbf{f}_{i_1 \dots i_r}(\mathbf{x}^{(t)}, \theta_{i_1 \dots i_r})] \\ &\times [\mathbf{y}^{(t)} - \mathbf{f}_{i_1 \dots i_r}(\mathbf{x}^{(t)}, \theta_{i_1 \dots i_r})]^T\} \Sigma_{i_1 \dots i_r}^{-1} \end{aligned} \quad (63)$$

and the derivatives of the log likelihood

$$\begin{aligned} \frac{\partial l}{\partial \theta_{i_1 \dots i_r}^g} \bigg|_{\theta_{i_1 \dots i_r}^g = \theta_{i_1 \dots i_r}^{g(k)}} &= \sum_{t=1}^N h_{i_1}^{(k)}(t) h_{i_2}^{(k)}(t) \dots h_{i_r}^{(k)}(t) \\ &\times \sum_{i_{r+1}=1}^{K_{i_1 \dots i_r}} [h_{i_1 \dots i_{r+1}}^{(k)}(t) - g_{i_1 \dots i_{r+1}}(\mathbf{x}^{(t)}, \theta_{i_1 \dots i_r}^{g(k)})] \\ &\times \frac{\partial s_{i_1 \dots i_{r+1}}}{\partial \theta_{i_1 \dots i_r}^g} \bigg|_{\theta_{i_1 \dots i_r}^g = \theta_{i_1 \dots i_r}^{g(k)}} \end{aligned} \quad (64)$$

$$\begin{aligned} \frac{\partial Q}{\partial \boldsymbol{\theta}_{i_1 \dots i_r}} \bigg|_{\boldsymbol{\theta}_{i_1 \dots i_r} = \boldsymbol{\theta}_{i_1 \dots i_r}^{(k)}} &= \sum_{t=1}^N h_{i_1}^{(k)}(t) h_{i_2}^{(k)}(t) \dots h_{i_r}^{(k)}(t) \\ &\times \frac{\partial \mathbf{f}_{i_1 \dots i_r}^T(\mathbf{x}^{(t)}, \boldsymbol{\theta}_{i_1 \dots i_r})}{\partial \boldsymbol{\theta}_{i_1 \dots i_r}} \bigg|_{\boldsymbol{\theta}_{i_1 \dots i_r} = \boldsymbol{\theta}_{i_1 \dots i_r}^{(k)}} \\ &\times \Sigma_{i_1 \dots i_r}^{-1} [\mathbf{y}^{(t)} - \mathbf{f}_{i_1 \dots i_r}(\mathbf{x}^{(t)}, \boldsymbol{\theta}_{i_1 \dots i_r}^{(k)})] \quad (65) \end{aligned}$$

$$\begin{aligned} \frac{\partial Q}{\partial \Sigma_{i_1 \dots i_r}} \bigg|_{\Sigma_{i_1 \dots i_r} = \Sigma_{i_1 \dots i_r}^{(k)}} &= -\frac{1}{2} \sum_{t=1}^N h_{i_1}^{(k)}(t) h_{i_2}^{(k)}(t) \dots h_{i_r}^{(k)}(t) (\Sigma_{i_1 \dots i_r}^{(k)})^{-1} \\ &\times \{ \Sigma_{i_1 \dots i_r}^{(k)} - [\mathbf{y}^{(t)} - \mathbf{f}_{i_1 \dots i_r}(\mathbf{x}^{(t)}, \boldsymbol{\theta}_{i_1 \dots i_r}^{(k)})] \\ &\times [\mathbf{y}^{(t)} - \mathbf{f}_{i_1 \dots i_r}(\mathbf{x}^{(t)}, \boldsymbol{\theta}_{i_1 \dots i_r}^{(k)})]^T \} \\ &\times (\Sigma_{i_1 \dots i_r}^{(k)}) \quad (66) \end{aligned}$$

Based on these two sets of derivatives, we follow the same line of thought as in the proof of Theorem 1 to establish the following theorem for the HME architecture.

**THEOREM 3.** *For the HME architecture of eqns (46) and (47), with the parameter updates given by Algorithm 3, we have that for every node  $v_{i_1 \dots i_r}$*

$$\begin{aligned} \boldsymbol{\theta}_{i_1 \dots i_r}^{g(k+1)} - \boldsymbol{\theta}_{i_1 \dots i_r}^{g(k)} &= P_{i_1 \dots i_r}^{g(k)} \frac{\partial l}{\partial \boldsymbol{\theta}_{i_1 \dots i_r}^{g(k)}} \bigg|_{\boldsymbol{\theta}_{i_1 \dots i_r} = \boldsymbol{\theta}_{i_1 \dots i_r}^{g(k)}}, \\ \boldsymbol{\theta}_{i_1 \dots i_r}^{(k+1)} - \boldsymbol{\theta}_{i_1 \dots i_r}^{(k)} &= P_{i_1 \dots i_r}^{(k)} \frac{\partial l}{\partial \boldsymbol{\theta}_{i_1 \dots i_r}^{(k)}} \bigg|_{\boldsymbol{\theta}_{i_1 \dots i_r} = \boldsymbol{\theta}_{i_1 \dots i_r}^{(k)}}, \\ \text{vec}[\Sigma_{i_1 \dots i_r}^{(k+1)}] - \text{vec}[\Sigma_{i_1 \dots i_r}^{(k)}] &= P_{i_1 \dots i_r}^{(k)} \frac{\partial l}{\partial \text{vec}[\Sigma_{i_1 \dots i_r}^{(k)}]} \bigg|_{\Sigma_{i_1 \dots i_r} = \Sigma_{i_1 \dots i_r}^{(k)}}, \quad (67) \end{aligned}$$

where  $l = \ln L$  is defined by eqns (5), (46), and (47).

Moreover, assuming that the training set  $\mathcal{Y}$  is generated by the HME model of eqns (46) and (47), and that the number  $N$  is sufficiently large, then  $P_{i_1 \dots i_r}^{g(k)}$  is a positive definite matrix, and  $P_{i_1 \dots i_r}^{(k)}$ ,  $P_{\Sigma_{i_1 \dots i_r}}^{(k)}$  are positive definite matrices with probability one. Specifically, they take the following values:

- (i)  $P_{i_1 \dots i_r}^{g(k)} = \gamma_{i_1 \dots i_r}^g (R_{i_1 \dots i_r}^{g(k)})^{-1}$  with  $R_{i_1 \dots i_r}^{g(k)}$  given by eqn (58).
- (ii)  $P_{i_1 \dots i_r}^{(k)} = (R_{i_1 \dots i_r}^{(k)})^{-1}$  with  $R_{i_1 \dots i_r}^{(k)}$  given by eqn (56).
- (iii) For  $P_{\Sigma_{i_1 \dots i_r}}^{(k)}$ , we have

$$P_{\Sigma_{i_1 \dots i_r}}^{(k)} = \frac{2}{\Sigma_{t=1}^N h_{i_1}^{(k)}(t) h_{i_2}^{(k)}(t) \dots h_{i_r}^{(k)}(t)} \Sigma_{i_1 \dots i_r}^{(k)} \otimes \Sigma_{i_1 \dots i_r}^{(k)} \quad (68)$$

where “ $\otimes$ ” denotes the Kronecker product as defined in Theorem 3.

From Theorem 5, we can again reach Corollary 1. Again, we see that the EM algorithm for training the HME architecture is a type of variable metric gradient ascent algorithm for maximizing  $l = \ln L$ .

Finally, we can also generalize Theorem 2 as follows.

**THEOREM 4.** *Assume that the training set  $\mathcal{Y}$  is generated by the HME model of eqns (46) and (47) and that the number  $N$  is sufficiently large. Assume that  $\Sigma_{i_1 \dots i_r}$  is diagonal and  $\mathbf{v}_{\Sigma_{i_1 \dots i_r}}$  is a vector consisting of the diagonal elements of  $\Sigma_{i_1 \dots i_r}$ .*

Let  $\boldsymbol{\Theta}$  be a vector produced by cascading every vector  $\mathbf{v}_{i_1 \dots i_r} = [(\boldsymbol{\theta}_{i_1 \dots i_r}^g)^T, \boldsymbol{\theta}_{i_1 \dots i_r}^T, \mathbf{v}_{\Sigma_{i_1 \dots i_r}}^T]^T$  of every node  $v_{i_1 \dots i_r}$  in the HME architecture. Let  $P$  be a diagonal block matrix with each diagonal item being a positive diagonal block matrix  $D_{i_1 \dots i_r}^b = \text{diag}[P_{i_1 \dots i_r}^g, P_{i_1 \dots i_r}, P_{\Sigma_{i_1 \dots i_r}}]$ . The items of  $\boldsymbol{\Theta}$  and  $P$  are arranged in such a way that  $D_{i_1 \dots i_r}^b$  in  $P$  corresponds to  $\mathbf{v}_{i_1 \dots i_r}$  in  $\boldsymbol{\Theta}$ .

Furthermore, assume that on a given domain  $D_{\boldsymbol{\Theta}}$

- (i) The parameterized functions of all the expert networks and the gating networks have second-order continuous derivatives.
- (ii) The Hessian matrix

$$H(\boldsymbol{\Theta}) = \frac{\partial^2 l(\boldsymbol{\Theta})}{\partial \boldsymbol{\Theta} \partial \boldsymbol{\Theta}^T}$$

is negative definite;

- (iii)  $\boldsymbol{\Theta}^*$  is a local maximum of  $l(\boldsymbol{\Theta})$ , and  $\boldsymbol{\Theta}^* \in D_{\boldsymbol{\Theta}}$ .

Then we have the same conclusion as given in Theorem 2. That is,

(1) Letting  $-M, -m$  (here  $M > m > 0$ ) be the minimum and maximum eigenvalues of the negatively definite matrix  $(P^{1/2})^T H(\boldsymbol{\Theta}) (P^{1/2})$  [or equivalently the minimum and maximum eigenvalues of  $PH(\boldsymbol{\Theta})$ ], because we have  $P\mathbf{H}\mathbf{e} = \lambda\mathbf{e}$  from  $(P^{1/2})^T H P^{1/2} \mathbf{e} = \lambda\mathbf{e}$ , we have

$$\begin{aligned} l(\boldsymbol{\Theta}^*) - l(\boldsymbol{\Theta}^{(k)}) &\leq r^k [l(\boldsymbol{\Theta}^*) - l(\boldsymbol{\Theta}_0)], \\ \|P^{-1/2}(\boldsymbol{\Theta}^{(k)} - \boldsymbol{\Theta}^*)\| &\leq |r|^{k/2} \sqrt{\frac{2}{m} [l(\boldsymbol{\Theta}^*) - l(\boldsymbol{\Theta}_0)]}, \end{aligned}$$

where  $r = 1 - (1 - M/2)m^2/M < 1$ . We also have  $0 < |r| < 1$  when  $M < 2$ .

(2) For any initial point  $\Theta_0 \in D_\Theta$ ,  $\lim_{k \rightarrow \infty} \Theta^{(k)} = \Theta^*$  when  $M < 2$ .

We should point out that the similarity in the conclusions of Theorem 2 and Theorem 4 does not mean that the EM algorithm for the hierarchical architecture has the same convergence rate as that for the one-level architecture. In the two cases the matrix  $P$  is different, and thus  $M, m, r$  are also different. This results in different convergence rates. Indeed, in practice, the hierarchical architecture is usually faster. The similarity in the conclusions of the two theorems does mean, however, that the convergence rates for both the algorithms are of the same (linear) order.

#### 4. VARIANTS OF THE EM ALGORITHM AND SIMULATIONS

##### 4.1. Variants of the EM Algorithm

For convenience, we denote an EM update of the parameter vector as follows

$$\Theta^{(k+1)} = M_p(\Theta^{(k)}). \quad (69)$$

From Theorems 1 and 2 and Theorems 3 and 4, we see that this update is actually a line search method along an ascent direction of  $l = \ln L$

$$\Theta^{(k+1)} = \Theta^{(k)} + P_\Theta^{(k)} \frac{\partial l}{\partial \Theta} \bigg|_{\Theta=\Theta^{(k)}}, \quad (70)$$

with  $P_\Theta^{(k)}$  being a positive definite matrix evaluated at  $\Theta^{(k)}$ . Moreover this update has a linear convergence rate. This link between the EM algorithm and conventional gradient-based optimization techniques suggests the possibility of using acceleration techniques for improving the convergence. In the sequel we suggest two such acceleration techniques.

##### 4.1.1. Modified Line Search

Equation (70) can be replaced by a modified line search

$$\begin{aligned} \Theta^{(k+1)} &= \Theta^{(k)} + \lambda_k \mathbf{d}_k, \\ \mathbf{d}_k &= P_\Theta^{(k)} \frac{\partial l}{\partial \Theta} \bigg|_{\Theta=\Theta^{(k)}} = M_p(\Theta^{(k)}) - \Theta^{(k)}, \end{aligned} \quad (71)$$

where  $\lambda_k$  is a stepsize that is optimized by maximizing  $l(\Theta^{(k)} + \lambda_k \mathbf{d}_k)$  with respect to  $\lambda_k$  via a one-dimensional search (e.g., Fibonacci search).

The implementation of a one-dimensional optimization method at every parameter update is typically

expensive. One often uses an inaccurate line search by decreasing (e.g.,  $\lambda_k \rightarrow r\lambda_k$ ) or increasing (e.g.,  $\lambda_k \rightarrow 1/r\lambda_k$ ) the stepsize heuristically according to a stopping rule. One frequently used stopping rule is the so-called Goldstein test (Luenberger, 1984). The Goldstein test is implemented as

$$\begin{aligned} l(\lambda_k) &\leq l(0) + \varepsilon l'(0)\lambda_k, \\ l(\lambda_k) &> l(0) + (1 - \varepsilon)l'(0)\lambda_k, \\ l'(0) &= \mathbf{d}_k^T \frac{\partial l}{\partial \Theta} \bigg|_{\Theta=\Theta^{(k)}}, \end{aligned} \quad (72)$$

where  $0 < \varepsilon < 1$  is a specified error bound.

Interestingly, if we rewrite the update as

$$\begin{aligned} \Theta^{(k+1)} &= \Theta^{(k)} + \lambda_k [M_p(\Theta^{(k)}) - \Theta^{(k)}] \\ &= (1 - \lambda_k)\Theta^{(k)} + \lambda_k M_p(\Theta^{(k)}), \end{aligned}$$

we find that eqn (71) is identical to the speedup technique for the EM algorithm studied by Peters and Walker (1978a, b), Meilijson (1989), and Redner and Walker (1984). These authors reported a significant speedup for an appropriately selected  $\lambda_k$  even without the Goldstein test.

##### 4.1.2. A Speed-up Formula Based on Local Linearization

Using a first-order Taylor expansion of  $M_p(\Theta^{(k)})$  around  $\Theta^{(k-1)}$ , we have, approximately,

$$\begin{aligned} M_p(\Theta^{(k)}) &= M_p(\Theta^{(k-1)}) + B(\Theta^{(k)} - \Theta^{(k-1)}) \\ \Theta^{(k+1)} &= \Theta^{(k)} + B(\Theta^{(k)} - \Theta^{(k-1)}) \end{aligned}$$

or

$$\Delta \Theta_k = B \Delta \Theta_{k-1}, \quad (73)$$

where

$$B = \frac{\partial M_p}{\partial \Theta} \bigg|_{\Theta=\Theta^{(k)}}$$

and  $\Delta \Theta_k = \Theta^{(k+1)} - \Theta^{(k)}$ .

For the matrix  $B$ , we have the following characteristic equation

$$\begin{aligned} \det(\lambda I - B) &= \lambda^n + \mu_1 \lambda^{n-1} + \dots + \mu_{n-1} \lambda + \mu_n \\ &= \lambda^n + \sum_{j=1}^n \mu_j \lambda^{n-j} = 0 \end{aligned}$$

where

$$\mu_j = (-1)^j \sum_{1 \leq i_1 < i_2 < \dots < i_j \leq n} \lambda_{i_1} \lambda_{i_2} \dots \lambda_{i_j}$$

and  $\lambda_i$  are the eigenvalues of  $B$ . It follows from the Cayley–Hamilton theorem that

$$B^n + \sum_{j=1}^n \mu_j B^{n-j} = 0.$$

Multiplying by  $\Delta \theta_{k-n}$  ( $k \geq n$ ), we have

$$B^n \Delta \theta_{k-n} + \sum_{j=1}^n \mu_j B^{n-j} \Delta \theta_{k-n} = 0, \quad (74)$$

From eqn (73), we obtain

$$\Delta \theta_{k-j} = B \Delta \theta_{k-j-1} = \dots = B^{n-j} \Delta \theta_{k-n}, \quad j = 0, 1, \dots, n.$$

Substituting into eqn (74), we have

$$\Delta \theta_k + \sum_{j=1}^n \mu_j \Delta \theta_{k-j} = 0. \quad (75)$$

Assuming that in comparison with the first  $l$  eigenvalues, the remaining eigenvalues can be neglected, we have approximately

$$\lambda^n + \sum_{j=1}^l \mu_j \lambda^{n-j} = 0, \quad B^n + \sum_{j=1}^l \mu_j B^{n-j} = 0$$

and correspondingly

$$\Delta \theta_i + \sum_{j=1}^l \mu_j \Delta \theta_{i-j} = 0, \quad i = k, k+1, \dots \quad (76)$$

The approximation becomes exact when the last  $n-l$  eigenvalues are zero.

By minimizing

$$\left\| \Delta \theta_k + \sum_{j=1}^l \mu_j \Delta \theta_{k-j} \right\|^2,$$

we obtain the following linear equation for  $\boldsymbol{\mu} = [\mu_1, \dots, \mu_l]^T$

$$S \boldsymbol{\mu} = \mathbf{s}_0, \quad S = [s_{ij}]_{l \times l}, \quad \mathbf{s}_0 = [-s_{01}, -s_{02}, \dots, -s_{0l}]^T, \quad (77)$$

$$s_{ij} = s_{ji} = (\Delta \theta_{k-i})^T (\Delta \theta_{k-j}).$$

We make use of the following identity

$$\begin{aligned} \sum_{i=k}^{\infty} \sum_{j=0}^l \mu_j \Delta \theta_{i-j} &= \sum_{j=0}^l \mu_j \Delta \theta_{k-j} + \sum_{i=k+1}^{\infty} \Delta \theta_i \\ &+ \sum_{i=k+1}^{\infty} \sum_{j=1}^l \mu_j \Delta \theta_{i-j} \end{aligned}$$

where  $\mu_0 = 1$ . From eqn (76), the left-hand side of this expression is the sum of zeros and thus zero. The first term on the right-hand side is also zero. We also have

$$\sum_{i=k+1}^{\infty} \Delta \theta_{i-j} = -\theta^{(k+1-j)} + \theta^*$$

(where  $\theta^* = \lim_{k \rightarrow \infty} \theta^{(k)}$ ). Using this identity together with eqn (76), we have

$$-\theta^{(k+1)} + \theta^* + \sum_{j=1}^l \mu_j (-\theta^{(k+1-j)} + \theta^*) = 0. \quad (78)$$

On the other hand, we can also show

$$\sum_{j=0}^l \mu_j \theta^{(k+1-j)} = \theta^{(k+1)} \sum_{j=0}^l \mu_j - \sum_{i=0}^{l-1} \left[ \sum_{j=i+1}^l \mu_j \right] \Delta \theta_{k-i}, \quad (79)$$

This can be easily shown by exchanging the order of summation in

$$\sum_{i=0}^{l-1} \left[ \sum_{j=i+1}^l \mu_j \right] \Delta \theta_{k-i},$$

yielding

$$\sum_{j=1}^l \mu_j \left[ \sum_{i=0}^{l-1} \Delta \theta_{k-i} \right] = \sum_{j=1}^l \mu_j (\theta^{(k+1-j)} - \theta^{(k+1)}).$$

Letting  $\theta_{k+1}^*$  replace  $\theta^*$  in eqn (78) and utilizing eqn (79), we finally obtain

$$\theta_{k+1}^* = \theta^{(k+1)} - \frac{\sum_{i=0}^{l-1} \left[ \sum_{j=i+1}^l \mu_j \right] \Delta \theta_{k-i}}{\sum_{j=0}^l \mu_j}. \quad (80)$$

This formula can be used for speeding up the EM algorithm in two ways.

- Given an initial  $\theta^{(0)}$ , we compute  $\theta^{(1)}, \dots, \theta^{(l+1)}$  via the EM update and then from  $\theta^{(j)}$ ,  $j = 1, \dots, l$ , we solve for  $\mu_1, \dots, \mu_l$  utilizing eqn (77). This yields a new  $\theta_{l+1}^*$  via eqn (80). We then let  $\theta^{(0)} = \theta_{l+1}^*$  and repeat the cycle.
- Instead of starting a new cycle after obtaining  $\theta_{l+1}^*$ , we simply let  $\theta^{(l+1)} = \theta_{l+1}^*$ , and use the EM update [eqn (69)] to obtain a new  $\theta^{(l+2)}$ , then we use  $\theta^{(0)}, \theta^{(1)}, \dots, \theta^{(l+2)}$  to obtain  $\theta_{l+2}^*$ . Similarly, after obtaining  $\theta_k^*$ ,  $k > l$ , we let  $\theta^{(k)} = \theta_k^*$  and use the EM update to get a new  $\theta^{(k+1)}$ , and then use  $\theta^{(k-1)}, \dots, \theta^{(k+1)}$  to obtain  $\theta_{k+1}^*$ .



Specifically, when  $l = 1$ , we have

$$\theta_{k+1}^* = \theta^{(k)} + \frac{\Delta\theta_k}{1 + \mu_1}, \quad \mu_1 = -\frac{(\Delta\theta_k)^T(\Delta\theta_{k-1})}{(\Delta\theta_{k-1})^T(\Delta\theta_{k-1})}. \quad (81)$$

In this case, the extra computation required by the acceleration technique is quite small.

#### 4.2. Simulations

We conducted two sets of computer simulations to compare the performance of the EM algorithm with the two variants described in the previous section. The training data for each simulation consisted of 1000 data points generated from the piecewise linear function  $y = a_1x + a_2 + n_t$ ,  $x \in [x_L, x_U]$  and  $y = a'_1x + a'_2 + n_t$ ,  $x \in [x'_L, x'_U]$ , where  $n_t$  is a Gaussian random variable with zero mean and variance  $\sigma^2 = 0.3$ . Training data were sampled from

the first function with probability 0.4 and from the second function with probability 0.6.

We studied a modular architecture with  $K = 2$  expert networks. The experts were linear; that is,  $f_j(\mathbf{x}^{(i)}, \theta_j)$  were linear functions  $[x, 1]^T \theta_j$ . For the gating net, we have

$$s_j = [x, 1]^T \theta_0,$$

where  $g_j$  is given by eqn (2). For simplicity, we updated  $\theta_{0j}$  by gradient ascent

$$\theta_{0j}^{(k+1)} = \theta_{0j}^{(k)} + r_g \frac{\partial Q}{\partial \theta_{0j}}. \quad (82)$$

The learning rate parameter was set to  $r_g = 0.05$  for the first data set and  $r_g = 0.002$  for the second data set. We used eqn (19) and eqn (21) to update the parameters  $\theta_j^{(k+1)}$  and  $\sigma_j^{(k+1)}$ , respectively.

The initial values of  $\theta_{0j}^{(0)}$ ,  $\theta_j^{(0)}$  and  $\sigma_j^{(0)}$  were picked

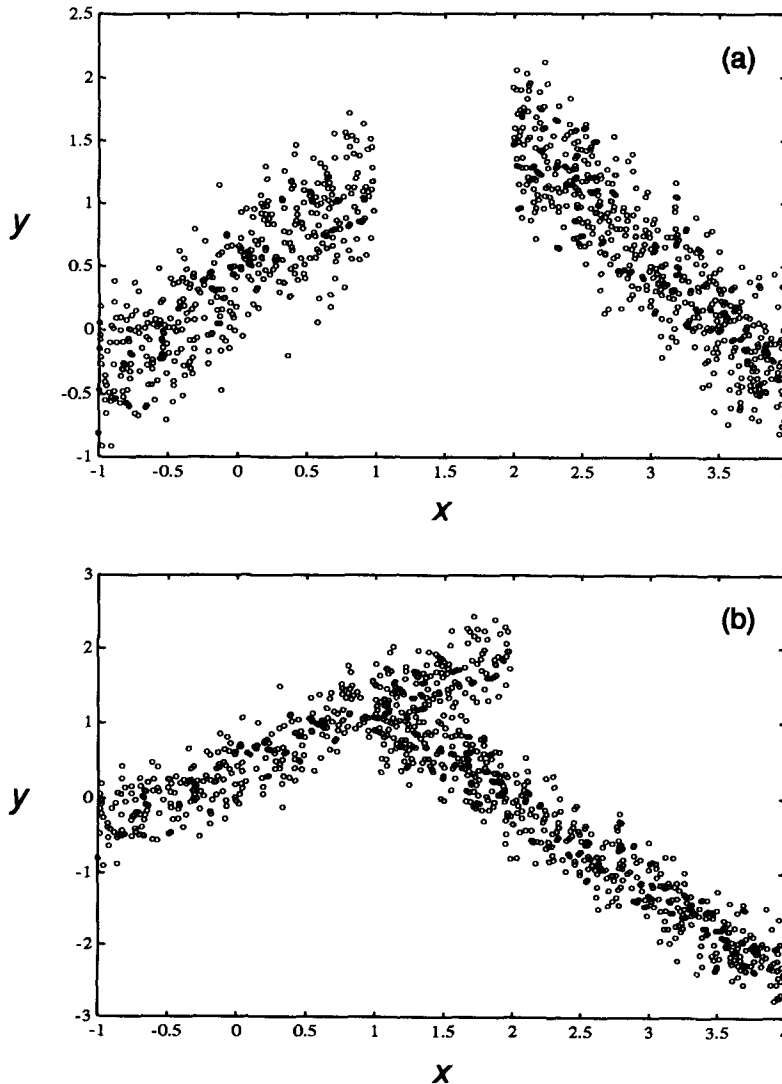


FIGURE 3. The data sets for the simulation experiments.

randomly. To compare the performance of the algorithms, we let each algorithm start from the same set of initial values.

The first data set (Figure 3a) was generated using the following parameter values:  $a_1 = 0.8$ ,  $a_2 = 0.4$ ,  $x_L = -1.0$ ,  $x_U = 1.0$ ,  $a'_1 = -1.0$ ,  $a'_2 = 3.6$ ,  $x'_L = 2.0$ ,  $x'_U = 4.0$ . The performance of the original algorithm,

the modified line search variant with  $\lambda_k = 1.1$ , the modified line search variant with  $\lambda_k = 0.5$ , and the algorithm based on local linearization are shown in Figures 4–7, respectively. As seen in Figures 4a and 5a, the log likelihood converged after 19 steps using both the original algorithm and the modified line search variant with  $\lambda_k = 1.1$ . When a smaller value

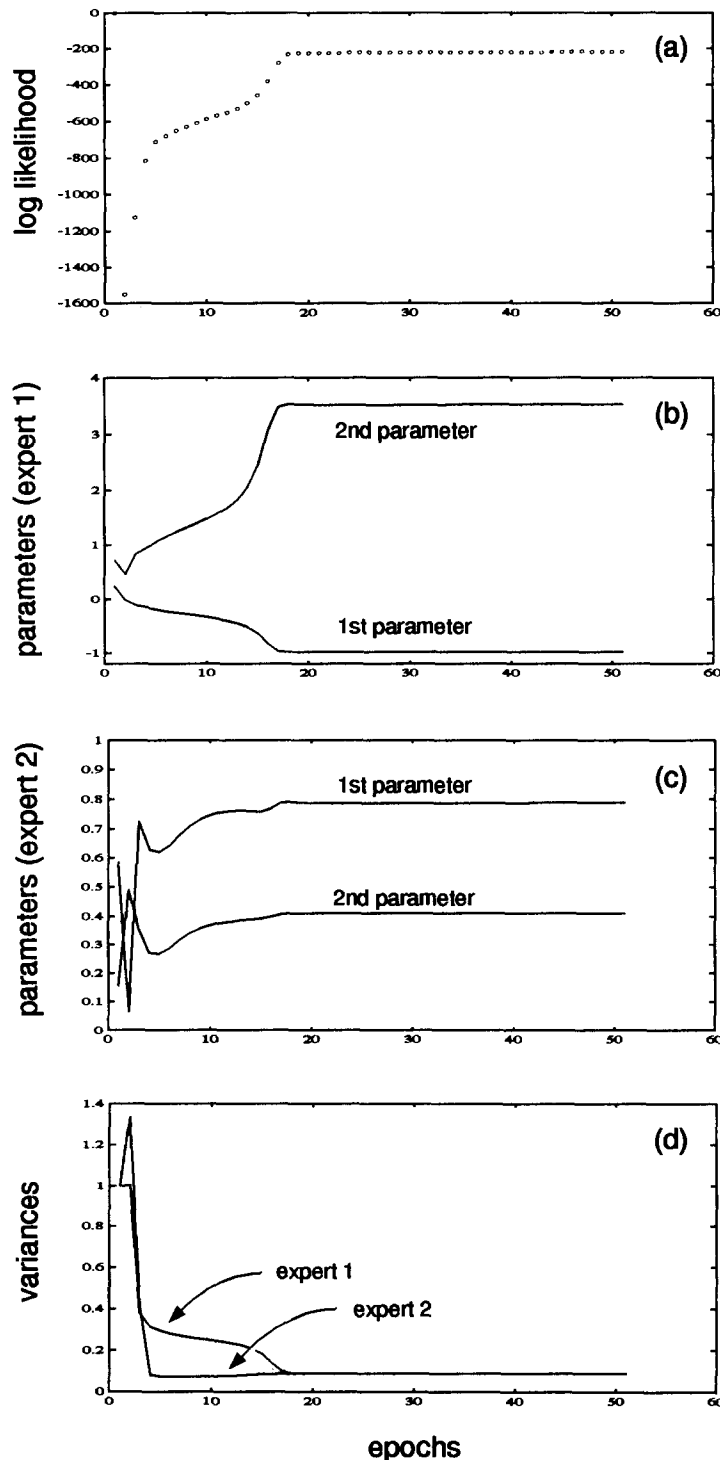


FIGURE 4. The performance of the original EM algorithm: (a) the evolution of the log likelihood; (b) the evolution of the parameters for expert network 1; (c) the evolution of the parameters for expert network 2; (d) the evolution of the variances.

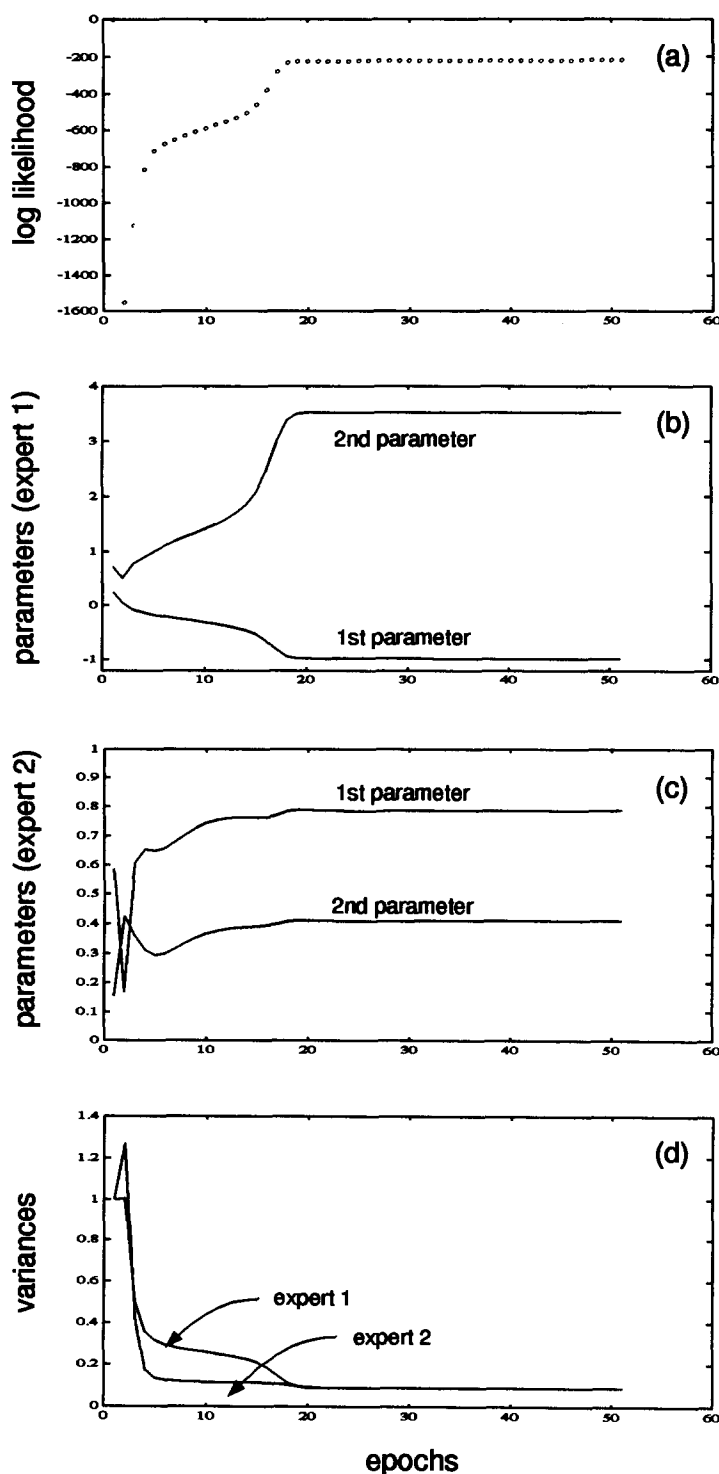


FIGURE 5. The performance of the linesearch variant with  $\lambda_k = 1.1$ .

was used ( $\lambda_k = 0.5$ ), the algorithm converged after 24 steps (Figure 6a). Trying other values of  $\lambda_k$ , we verified that  $\lambda_k < 1$  slows down the convergence, whereas  $\lambda_k > 1$  may speed up the convergence (cf. Redner & Walker, 1984). We found, however, that the outcome was quite sensitive to the selection of the value of  $\lambda_k$ . For example, setting  $\lambda_k = 1.2$  led the algorithm to diverge. Allowing  $\lambda_k$  to be determined

by the Goldstein test [eqn (72)] yielded results similar to the original algorithm, but required more computer time. Finally, Figure 7a shows that the algorithm based on local linearization yielded substantially improved convergence—the log likelihood converged after only eight steps.

Figure 4b and c shows the evolution of the parameters for the first expert net and the second

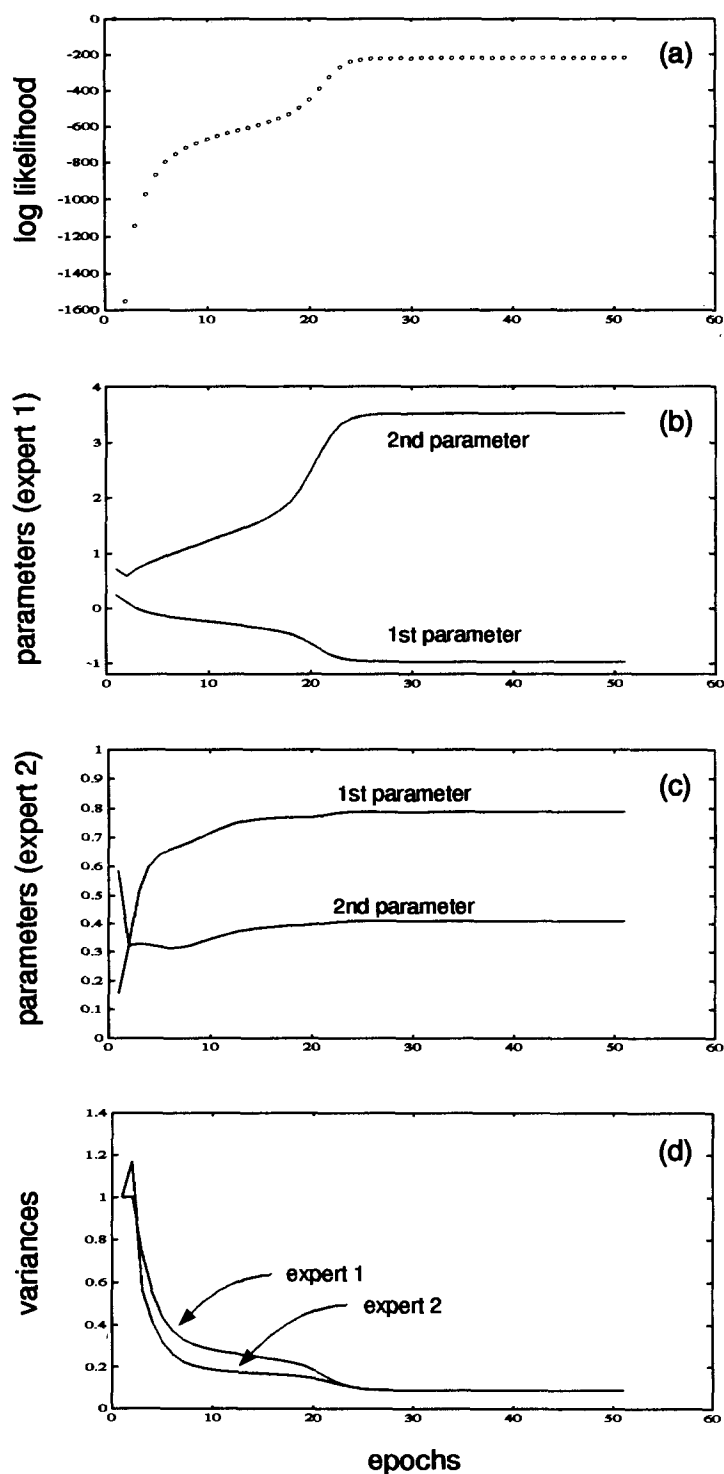


FIGURE 6. The performance of the linesearch variant with  $\lambda_k = 0.5$ .

expert net, respectively. Comparison of these figures to Figure 5b and c shows that the original algorithm and the modified line search variant with  $\lambda_k = 1.1$  behaved almost identically:  $\theta_1$  converged to the correct solution after about 18 steps in either case. Figure 6b and c shows the slowdown obtained by using  $\lambda_k = 0.5$ . Figure 7b and c shows the improved performance obtained using the local linearization

algorithm. In this case, the weight vectors converged to the correct values within seven steps.

Panel d in each of the figures shows the evolution of the estimated variances  $\sigma_1^2$  and  $\sigma_2^2$ . The results were similar to those for the expert net parameters. Again, the algorithm based on local linearization yielded significantly faster convergence than the other algorithms.

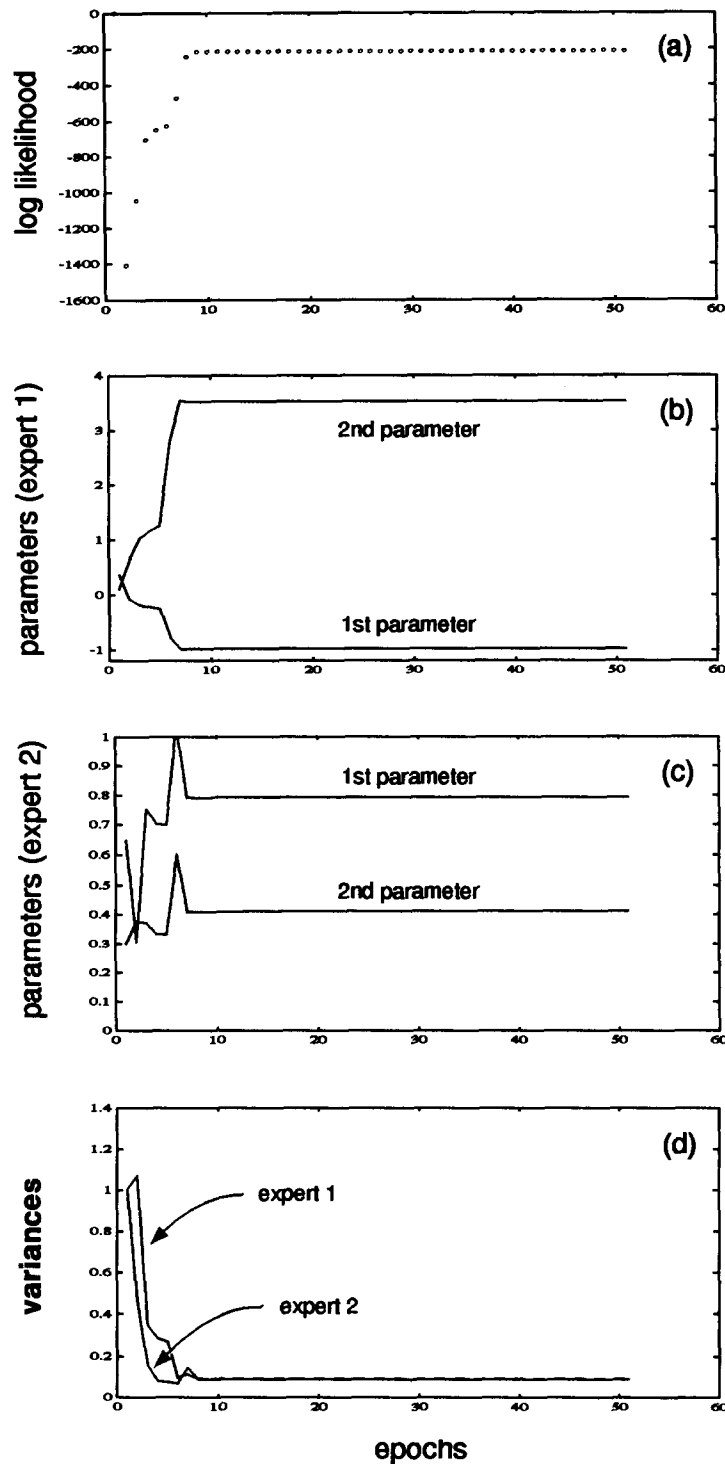


FIGURE 7. The performance of the algorithm based on local linearization.

A second simulation was run using the following parameter values (Figure 3b):  $a_1 = 0.8$ ,  $a_2 = 0.4$ ,  $x_L = -1.0$ ,  $x_U = 2.0$ ,  $a'_1 = -1.2$ ,  $a'_2 = 2.4$ ,  $x'_L = 1.0$ ,  $x'_U = 4.0$ . The results obtained in this simulation were similar to those obtained in the first simulation. The EM algorithm converged in 11 steps and the local linearization algorithm converged in six steps.

The results from a number of other simulation

experiments confirmed the results reported here. In general the algorithm based on local linearization provided significantly faster convergence than the original EM algorithm. The modified line search variant did not appear to converge faster (if the parameter  $\lambda_k$  was fixed). We also tested gradient ascent in these experiments and found that convergence was generally one to two orders of magnitude

slower than convergence of the EM algorithm and its variants. Moreover, convergence of gradient ascent was rather sensitive to the learning rate and the initial values of the parameters.

## 5. CONCLUDING REMARKS

Finite mixture models have become increasingly popular as models for unsupervised learning, partly because they occupy an interesting niche between parametric and nonparametric approaches to statistical estimation. Mixture-based approaches are parametric in that particular parametric forms must be chosen for the component densities, but they can also be regarded as nonparametric by allowing the number of components of the mixture to grow. The advantage of this niche in statistical theory is that these models have much of the flexibility of nonparametric approaches, but retain some of the analytical advantages of parametric approaches (McLachlan & Basford, 1988). Similar remarks can be made in the case of supervised learning: the ME architecture and the HME architecture provide flexible models for general nonlinear regression while retaining a strong flavor of parametric statistics. The latter model, in particular, compares favorably to decision tree models in this regard (Jordan & Jacobs, 1994).

In the current paper we have contributed to the theory of mixture-based supervised learning. We have analyzed an EM algorithm for ME and HME architectures and provided theorems on the convergence of this algorithm. In particular, we have shown that a learning algorithm can be regarded as a variable metric algorithm with its metric matrix  $P$  being positive definite, so that the searching direction of the algorithm always has a positive projection on the gradient of the log likelihood. We have shown that the algorithm converges linearly, with a rate determined by the difference between the minimal and maximal eigenvalues of a negative definite matrix.

Similar results to those obtained here can also be obtained for the case of the unsupervised learning of finite mixtures (Xu & Jordan, 1994).

## REFERENCES

- Baum, L. E., & Sell, G. R. (1968). Growth transformation for functions on manifolds. *Pacific Journal of Mathematics*, **27**, 211–227.
- Baum, L. E., Petrie, T., Soules, G., & Weiss, N. (1970). A maximization technique occurring in the statistical analysis of probabilistic functions of Markov chains. *The Annals of Mathematical Statistics*, **41**, 164–171.
- Breiman, L., Friedman, J. H., Olshen, R. A., & Stone, C. J. (1984). *Classification and regression trees*. Belmont, CA: Wadsworth International Group.
- De Veaux, R. D. (1986). *Parameter estimation for a mixture of linear regressions*. Doctoral Dissertation and Tech. Rept.

- No. 247, Department of Statistics, Stanford University, Stanford, CA.
- Dempster, A. P., Laird, N. M., & Rubin, D. B. (1977). Maximum-likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society*, **B39**, 1–38.
- Friedman, J. H. (1991). Multivariate adaptive regression splines. *The Annals of Statistics*, **19**, 1–141.
- Jacobs, R. A., Jordan, M. I., Nowlan, S. J., & Hinton, G. E. (1991). Adaptive mixtures of local experts. *Neural Computation*, **3**, 79–87.
- Jordan, M. I., & Jacobs, R. A. (1992). Hierarchies of adaptive experts. In J. E. Moody, S. Hanson, & R. P. Lippmann (Eds.), *Advances in neural information processing system 4* (pp. 985–992). San Mateo: Morgan Kaufmann.
- Jordan, M. I., & Jacobs, R. A. (1994). Hierarchical mixtures of experts and the EM algorithm. *Neural Computation*, **6**, 181–214.
- Little, R. J. A., & Rubin, D. B. (1987). *Statistical analysis with missing data*. New York: John Wiley.
- McCullagh, P., & Nelder, J. A. (1983). *Generalized linear models*. London: Chapman and Hall.
- McLachlan, G. J., & Basford, K. E. (1988). *Mixture models: Inference and application to clustering*. New York: Dekker.
- Meilijson, I. (1989). A fast improvement to the EM algorithm on its own terms. *Journal of the Royal Statistical Society*, **B51**, 127–138.
- Peters, B. C., & Walker, H. F. (1978a). An iterative procedure for obtaining maximum-likelihood estimates of the parameters for a mixture of normal distributions. *SIAM Journal of Applied Mathematics*, **35**, 362–378.
- Peters, B. C., & Walker, H. F. (1978b). The numerical evaluation of the maximum-likelihood estimates of a subset of mixture proportions. *SIAM Journal of Applied Mathematics*, **35**, 447–452.
- Quandt, R. E., & Ramsey, J. B. (1972). A new approach to estimating switching regressions. *Journal of the American Statistical Society*, **67**, 306–310.
- Quandt, R. E., & Ramsey, J. B. (1978). Estimating mixtures of normal distribution and switching regressions. *Journal of the American Statistical Society*, **73**, 730–738.
- Redner, R. A., & Walker, H. F. (1984). Mixture densities, maximum likelihood, and the EM algorithm. *SIAM Review*, **26**, 195–239.
- Wu, C. F. J. (1983). On the convergence properties of the EM algorithm. *The Annals of Statistics*, **11**, 95–103.
- Xu, L., & Jordan, M. I. (1993). Theoretical and experimental studies of the EM algorithm for unsupervised learning based on finite Gaussian mixtures (Tech. Rep. 9302), MIT Computational Cognitive Science, MIT, Cambridge, MA.
- Xu, L., & Jordan, M. I. (1994). On convergence properties of the EM algorithm for Gaussian mixtures. *Neural Computation* (in press).

## APPENDIX

The theoretical results presented in the main text show that the EM algorithm for the ME and HME architectures converges linearly with a rate determined by the condition number of a particular matrix. These results were obtained for a special case in which the expert networks are linear with a Gaussian probability model and the gating networks are multinomial logit models. In this section we discuss extensions of these results to other architectures.

We first note that Theorems 2 and 4 make no specific reference to the particular probability models utilized in specifying the architecture. The results on convergence rate in these theorems require only that the matrix  $P$  be positive definite. These theorems apply directly to other architectures if the corresponding  $P$  matrices can be shown to be positive definite. We therefore need only consider generalizations of Theorem 1, the theorem that established the positive definiteness of  $P$  for the generalized linear ME

architectures. An analogous generalization of Theorem 3 for the HME architectures can also be obtained.

Let us consider the case in which the function implemented by each expert network ( $\mathbf{f}_j(\mathbf{x}, \boldsymbol{\theta}_j)$ ) is nonlinear in the parameters. We consider two possible updates for the parameters: (1) a gradient algorithm

$$\boldsymbol{\theta}_j^{(k+1)} = \boldsymbol{\theta}_j^{(k)} + \gamma_j \mathbf{e}_j^{(k)}, \quad (\text{A.1})$$

where

$$\mathbf{e}_j^{(k)} = \sum_{i=1}^N h_j^{(k)}(t) \frac{\partial \mathbf{f}_j^T(\mathbf{x}^{(i)}, \boldsymbol{\theta}_j^{(k)})}{\partial \boldsymbol{\theta}_j^{(k)}} (\Sigma_j^{(k)})^{-1} [\mathbf{y}^{(i)} - \mathbf{f}_j(\mathbf{x}^{(i)}, \boldsymbol{\theta}_j^{(k)})], \quad (\text{A.2})$$

and (2) a Newton algorithm:

$$\boldsymbol{\theta}_j^{(k+1)} = \boldsymbol{\theta}_j^{(k)} + \gamma_j (R_j^{(k)})^{-1} \mathbf{e}_j^{(k)}, \quad (\text{A.3})$$

where

$$R_j^{(k)} = \sum_{i=1}^N h_j^{(k)}(t) \frac{\partial \mathbf{f}_j^T(\mathbf{x}^{(i)}, \boldsymbol{\theta}_j^{(k)})}{\partial \boldsymbol{\theta}_j^{(k)}} (\Sigma_j^{(k)})^{-1} \frac{\partial \mathbf{f}_j(\mathbf{x}^{(i)}, \boldsymbol{\theta}_j^{(k)})}{\partial (\boldsymbol{\theta}_j^{(k)})^T}, \quad (\text{A.4})$$

where  $\gamma_j > 0$  is a learning rate.

These updates are covered by the following extension of Theorem 1.

**THEOREM 1A.** For the model given by eqn (3) and the updates given by eqn (84) or eqn (86), we have:

$$\boldsymbol{\theta}_j^{(k+1)} - \boldsymbol{\theta}_j^{(k)} = P_j^{(k)} \frac{\partial l}{\partial \boldsymbol{\theta}_j} \bigg|_{\boldsymbol{\theta}_j = \boldsymbol{\theta}_j^{(k)}}, \quad j = 1, \dots, K,$$

where  $P_j^{(k)}$  is positive definite.

*Proof.* For the gradient descent algorithm, we have  $P_j^{(k)} = \gamma_j I_K$ , which is obviously positive definite because  $\gamma_j > 0$ . For the Newton

algorithm, we have that  $P_j^{(k)} = \gamma_j (R_j^{(k)})^{-1}$ . We now show that this matrix is positive definite. For an arbitrary vector  $\mathbf{u}$ , we have

$$\begin{aligned} \mathbf{u}^T R_j^{(k)} \mathbf{u} &= \sum_{i=1}^N h_j^{(k)}(t) \mathbf{u}^T \frac{\partial \mathbf{f}_j^T(\mathbf{x}^{(i)}, \boldsymbol{\theta}_j^{(k)})}{\partial \boldsymbol{\theta}_j^{(k)}} (\Sigma_j^{(k)})^{-1} \frac{\partial \mathbf{f}_j(\mathbf{x}^{(i)}, \boldsymbol{\theta}_j^{(k)})}{\partial (\boldsymbol{\theta}_j^{(k)})^T} \mathbf{u} \\ &= \sum_{i=1}^N h_j^{(k)}(t) \mathbf{v}^T (\Sigma_j^{(k)})^{-1} \mathbf{v} \\ &\geq 0. \end{aligned}$$

Equality holds only when

$$\mathbf{v}^T = \mathbf{u}^T \frac{\partial \mathbf{f}_j^T(\mathbf{x}^{(i)}, \boldsymbol{\theta}_j^{(k)})}{\partial \boldsymbol{\theta}_j^{(k)}} = 0,$$

because  $\Sigma_j^{(k)}$  is positive definite with probability one. This is impossible for any  $\mathbf{u}$ . So with probability one,  $R_j^{(k)}$  [and thus  $(R_j^{(k)})^{-1}$  also] is positive definite. ■

Note that the Newton update [eqn (86)] is particularly appropriate for the case in which the experts are generalized linear models (McCullagh & Nelder, 1983); that is, the case in which  $\mathbf{f}_j(\mathbf{x}^{(i)}, \boldsymbol{\theta}_j) = [f_{j1}(\mathbf{x}^{(i)}, \boldsymbol{\theta}_j), \dots, f_{jd_y}(\mathbf{x}^{(i)}, \boldsymbol{\theta}_j)]$  ( $d_y$  is the dimension of  $\mathbf{y}$ ) with

$$f_{ji}(\mathbf{x}^{(i)}, \boldsymbol{\theta}_j) = F_{ji}([\boldsymbol{\theta}_{j,1}, \dots, \boldsymbol{\theta}_{j,m}]^T \mathbf{x}^{(i)} + \boldsymbol{\theta}_{j,m+1}),$$

where  $F_{ji}(\cdot)$  is a continuous univariate nonlinear function known as the *link* function. In this case the Newton algorithm reduces to the IRLS algorithm. The extension to generalized linear models also allows probability models for the generalized exponential family (cf. Jordan & Jacobs, 1994) and Theorem 1A is applicable to this case as well.

We can also consider the case in which the gating network is nonlinear in the parameters. Both the Newton update (IRLS update) and the gradient update are applicable in this case. Theorem 1 already established that the Newton update for the gating network involves a positive definite  $P_g$  matrix. As in Theorem 1A, the result for the gradient update is immediate.