Communicated by Christopher Bishop

SMEM Algorithm for Mixture Models

Naonori Ueda

Ryohei Nakano* NTT Communication Science Laboratories, Hikaridai, Seika-cho, Soraku-gun, Kyoto

619-0237 Japan

Zoubin Ghahramani

Geoffrey E. Hinton

Gatsby Computational Neuroscience Unit, University College London, London WC1N 3AR, U.K.

We present a split-and-merge expectation-maximization (SMEM) algorithm to overcome the local maxima problem in parameter estimation of finite mixture models. In the case of mixture models, local maxima often involve having too many components of a mixture model in one part of the space and too few in another, widely separated part of the space. To escape from such configurations, we repeatedly perform simultaneous split-and-merge operations using a new criterion for efficiently selecting the split-and-merge candidates. We apply the proposed algorithm to the training of gaussian mixtures and mixtures of factor analyzers using synthetic and real data and show the effectiveness of using the splitand-merge operations to improve the likelihood of both the training data and of held-out test data. We also show the practical usefulness of the proposed algorithm by applying it to image compression and pattern recognition problems.

1 Introduction ____

Mixture density models, in particular normal mixtures, have been used extensively in the field of statistical pattern recognition (MacLachlan & Basford, 1987). Recently, more sophisticated mixture density models such as mixtures of latent variable models (e.g., probabilistic PCA and factor analysis) have been proposed to approximate the underlying data manifold (Hinton, Dayan, & Revow, 1997; Tipping & Bishop, 1997; Ghahramani & Hinton, 1997). The parameters of these mixture models can be estimated using the expectation-maximization (EM) algorithm (Dempster, Laird, & Rubin, 1977) based on the maximum likelihood framework. A common and

Neural Computation 12, 2109-2128 (2000) © 2000 Massachusetts Institute of Technology

LETTER =

^{*} Present address: Nagoya Institute of Technology, Gokiso-cho, Showa-Ku, Nagoya 466–8555 Japan

serious problem associated with these EM algorithms, however, is the local maxima problem. Although this problem has been pointed out by many researchers, the best way to solve it, in practice, is still an open question.

Two of the authors have proposed the deterministic annealing EM (DAEM) algorithm (Ueda & Nakano, 1998), where a modified posterior probability parameterized by temperature is derived to avoid local maxima. However, when mixture density models are involved, local maxima arise when there are too many components of a mixture model in one part of the space and too few in another. The DAEM algorithm and other algorithms are not very effective at avoiding such local maxima because they are not able to move a component from an overpopulated region to an underpopulated region without passing through positions that give a lower like-lihood. We therefore introduce a discrete move that simultaneously merges two components in an overpopulated region and splits a component in an underpopulated region.

The idea of performing split-and-merge operations has been successfully applied to clustering (Ball & Hall, 1967) and vector quantization (Ueda & Nakano, 1994). Recently, split-and-merge operations have also been proposed for Bayesian normal mixture analysis (Richardson & Green, 1997). Since they use split-and-merge operations with a Markov chain Monte Carlo method, it is computationally much more costly than our algorithm. In addition, we introduce new split-and-merge criteria to select the split-and-merge candidates efficiently.

Although the proposed method, unlike the DAEM algorithm, is limited to mixture models, we have experimentally confirmed that our split-andmerge EM (SMEM) algorithm obtains better solutions than the DAEM algorithm. We have already given a basic idea of the SMEM algorithm (Ueda, Nakano, Ghahramani, & Hinton, 1999). This article describes the algorithm in detail and shows real applications, including image compression and pattern recognition.

2 EM Algorithm _

Before describing our SMEM algorithm, we briefly review the EM algorithm. Suppose that a set \mathcal{Z} consists of observed data \mathcal{X} and unobserved data \mathcal{Y} . $\mathcal{Z} = (\mathcal{X}, \mathcal{Y})$ and \mathcal{X} are called *complete data* and *incomplete data*, respectively. Assume that the joint probability density of \mathcal{Z} is parametrically given as $p(\mathcal{X}, \mathcal{Y}; \Theta)$, where Θ denotes parameters of the density to be estimated. The maximum likelihood estimate of Θ is a value of Θ that maximizes the incomplete data log-likelihood function:

$$\mathcal{L}(\Theta; \mathcal{X}) \stackrel{\text{def}}{=} \log p(\mathcal{X}; \Theta)$$
$$= \log \int p(\mathcal{X}, \mathcal{Y}; \Theta) d\mathcal{Y}.$$
(2.1)

The characteristic of the EM algorithm is to maximize the incomplete data

log-likelihood function by iteratively maximizing the expectation of the complete data log-likelihood function:

$$\mathcal{L}_{c}(\Theta; \mathcal{Z}) \stackrel{\text{def}}{=} \log p(\mathcal{X}, \mathcal{Y}; \Theta).$$
(2.2)

Suppose that $\Theta^{(t)}$ denotes the estimate of Θ obtained after the *t*th iteration of the algorithm. Then, at the *t* + 1th iteration, the E-step computes the expected complete data log-likelihood function denoted by $Q(\Theta|\Theta^{(t)})$ and defined by

$$Q(\Theta|\Theta^{(t)}) \stackrel{\text{def}}{=} \mathrm{E}\{\mathcal{L}_{c}(\Theta;\mathcal{Z})|\mathcal{X};\Theta^{(t)}\},\tag{2.3}$$

and the M-step finds the Θ maximizing $Q(\Theta|\Theta^{(t)})$. The convergence of the EM steps is theoretically guaranteed (Dempster, Laird, & Rubin, 1977).

3 Split-and-Merge EM Algorithm _

3.1 Split-and-Merge Operations. We restrict ourselves here to mixture density models. The probability density function (pdf) of a mixture of *M* density models is given by

$$p(\boldsymbol{x};\Theta) = \sum_{m=1}^{M} \alpha_m p_m(\boldsymbol{x};\theta_m), \qquad (3.1)$$

where α_m is the mixing proportion of the *m*th model¹ and satisfies $\alpha_m \ge 0$ and $\sum_{m=1}^{M} \alpha_m = 1$. The $p_m(\boldsymbol{x}; \theta_m)$ is a *d*-dimensional density model corresponding to the *m*th model. Clearly, $\Theta = \{(\alpha_m, \theta_m), m = 1, ..., M\}$ is an unknown parameter set.

In the case of mixture models, the model index $m \in \{1, ..., M\}$ is unknown for an observed data x_n and therefore m corresponds to the unobserved data mentioned in section 2. Noting that the pdf of the complete data is $p(x, m; \Theta) = \alpha_m p_m(x; \theta_m)$ in this case, we have

$$Q(\Theta|\Theta^{(t)}) = \sum_{n=1}^{N} \sum_{m=1}^{M} \{\log \alpha_m p_m(\boldsymbol{x}_n; \theta_m)\} P(m|\boldsymbol{x}_n; \Theta^{(t)}).$$
(3.2)

Here, $P(m|\boldsymbol{x}_n; \Theta^{(t)})$ is a posterior probability and is computed as

$$P(m|\boldsymbol{x}_{n}; \Theta^{(t)}) = \frac{\alpha_{m}^{(t)} p_{m}(\boldsymbol{x}_{n}; \theta_{m}^{(t)})}{\sum_{l=1}^{M} \alpha_{l}^{(t)} p_{l}(\boldsymbol{x}_{n}; \theta_{l}^{(t)})} .$$
(3.3)

N is the number of observed data points.

¹ Strictly speaking, we should call it the *m*th model component or model component ω_m , but to simplify the notation and wording, we call it model *m* or the *m*th model hereafter.

Looking at equation 3.2 carefully, one can see that the *Q* function can be represented in the form of a direct sum,

$$Q(\Theta|\Theta^{(t)}) = \sum_{m=1}^{M} q_m(\theta_m|\Theta^{(t)}), \qquad (3.4)$$

where

$$q_m(\theta_m|\Theta^{(t)}) = \sum_{n=1}^N P(m|\boldsymbol{x}_n;\Theta^{(t)}) \log \alpha_m p_m(\boldsymbol{x}_n;\theta_m)$$
(3.5)

and depends on only α_m and θ_m .

Let Θ^* denote the parameter values estimated by the usual EM algorithm. Then, after the EM algorithm has converged, the *Q* function can be rewritten as

$$Q^* = q_i^* + q_j^* + q_k^* + \sum_{m, m \neq i, j, k} q_m^*.$$
(3.6)

We then try to increase the first three terms of the right-hand side of equation 3.6 by merging models *i* and *j* to produce a model i', and splitting the model *k* into two models j' and k'.

3.1.1 Initialization. To reestimate the parameters of these new models, we have to initialize the parameters corresponding to the new models using Θ^* . Intuitively natural initializations are given below. The initial parameter values for the merged model *i*' are set as linear combinations of the original ones before the merge:

$$\alpha_{i'} = \alpha_i^* + \alpha_j^* \quad \text{and} \quad \theta_{i'} = \frac{\alpha_i^* \theta_i^* + \alpha_j^* \theta_j^*}{\alpha_i^* + \alpha_i^*}.$$
(3.7)

Noting that the mixing proportion is estimated as the average of posterior over data, $\alpha_l^* = 1/N \sum_{n=1}^{N} P(l|\boldsymbol{x}_n; \Theta^*)$. One can see that the initialization of θ_i is the linear combination of θ_i^* and θ_i^* , weighted by the posteriors.

On the other hand, as for models j' and k', we set

$$\alpha_{j'} = \alpha_{k'} = \frac{\alpha_k^*}{2} \qquad \theta_{j'} = \theta_k^* + \epsilon \quad \text{and} \quad \theta_{k'} = \theta_k^* + \epsilon', \tag{3.8}$$

where ϵ or ϵ' is some small, random perturbation vector or matrix (i.e., $\|\epsilon\| \ll \|\theta_k^*\|$). In the case of mixture gaussians, covariance matrices $\Sigma_{j'}$ and $\Sigma_{k'}$ should be positive definite. In this case, we can initialize them as

$$\Sigma_{j'} = \Sigma_{k'} = \det(\Sigma_k^*)^{1/d} I_d \tag{3.9}$$



Figure 1: An example of initialization in a two-dimensional gaussian case. (a) A gaussian just before split (left) and initialized gaussians just after split (right). (b) Two gaussians just before merge (left) and an initialized gaussian just after merge (right).

instead of equation 3.8. Here, $det(\Sigma)$ denotes the determinant of matrix Σ , and I_d is the *d*-dimensional identity matrix. Figure 1 shows a simple example of the initialization steps for a two-dimensional gaussian case using equations 3.7 through 3.9.

3.1.2 *Partial EM Steps.* The parameter reestimation for m' = i', j', and k' can be done by using EM steps, but instead of equation 3.3, we use the following modified posterior probability:

$$P(m'|\boldsymbol{x}; \Theta^{(t)}) = \frac{\alpha_{m'}^{(t)} p_{m'}(\boldsymbol{x}; \theta_{m'}^{(t)})}{\sum_{l=i', j', k'} \alpha_l^{(t)} p_l(\boldsymbol{x}; \theta_l^{(t)})}$$

$$\times \sum_{m=i,j,k} P(m|\boldsymbol{x}; \Theta^*), \quad \text{for } m' = i', j', k'.$$
(3.10)

Using equation 3.10, the sum of posterior probabilities for models i', j', and k' becomes equal to the sum of posterior probabilities for models i, j, and k just before split-and-merge. That is,

$$\sum_{m'=i',j',k'} P(m'|\boldsymbol{x};\Theta^{(t)}) = \sum_{m=i,j,k} P(m|\boldsymbol{x};\Theta^*)$$
(3.11)

always holds during the reestimation process. By this, we can reestimate the parameters for models i', j', and k' consistently without affecting the other models. We call these EM steps *partial* EM steps. These partial EM steps make the total algorithm efficient.

3.2 SMEM Algorithm. After the partial EM steps, the usual EM steps, called the *full EM steps*, are performed as a postprocessing operation. After these steps, if Q is improved, then we accept the new estimate and repeat the above after setting the new parameters to Θ^* . Otherwise we reject it, go back to Θ^* , and try another candidate. We summarize these as the following SMEM algorithm:

- 1. Perform the usual EM updates from some initial parameter value Θ until convergence. Let Θ^* and Q^* denote the estimated parameters and corresponding Q function value after the EM algorithm has converged, respectively.
- 2. Sort the split-and-merge candidates by computing split-and-merge criteria (described in the next section) based on Θ^* . Let $\{i, j, k\}_c$ denote the *c*th candidate.
- 3. For $c = 1, ..., C_{\text{max}}$, perform the following: After making the initial parameter settings based on Θ^* , perform the partial EM steps for $\{i, j, k\}_c$ and then perform the full EM steps until convergence. Let Θ^{**} be the obtained parameters and Q^{**} be the corresponding Q function value after the full EM has converged. If $Q^{**} > Q^*$, then set $Q^* \leftarrow Q^{**}$, $\Theta^* \leftarrow \Theta^{**}$ and go to step 2.
- 4. Halt with Θ^* as the final parameters.

Note that when a certain split-and-merge candidate that improves the Q function value is found in step 3, the other successive candidates are ignored. There is no guarantee therefore that the split-and-merge candidates that are chosen will give the largest possible improvement in Q. This is not a major problem, however, because the split-and-merge operations are performed repeatedly. If there were no heuristics for ordering potential split-and-merge operations, we would have to consider them all $C_{\text{max}} = M(M-1)$

(M - 2)/2, but experimentally we have confirmed that $C_{\text{max}} \simeq 5$ may be enough because the split-and-merge criteria do work well.

The SMEM algorithm monotonically increases the *Q* function value, and if the *Q* function value does not increase for all $c = 1, ..., C_{max}$, then the algorithm stops. Since the full EM steps equivalent to the original EM steps are performed after the convergence of the partial EM steps, it is clear that the SMEM algorithm maintains the global convergence properties of the EM algorithm.

In the SMEM algorithm, the split-and-merge operations are simultaneously performed so that the total number of mixture components is unchanged. In general, the Q function value increases as the number of parameters (or the number of mixture components) increases. Thus, in order to check whether Q is improved by the rearrangement of model components at step 3, it is necessary to keep the number of model components unchanged.

Intuitively, a simultaneous split-and-merge can be viewed as a way of tunneling through low-likelihood barriers, thereby eliminating many poor local optima. In this respect, it has some similarities with simulated annealing, but the moves that are considered are long range and very specific to the particular problems that arise when fitting mixture models.

3.3 Split-and-Merge Criteria. Each of the split-and-merge candidates can be evaluated by its *Q* function value after step 3 of the SMEM algorithm mentioned in section 3.2. However, since there are so many candidates, some reasonable criteria for ordering the split-and-merge candidates should be used to accelerate the SMEM algorithm.

3.3.1 *Merge Criterion.* In general, when there are many data points, each of which has almost equal posterior probability given by equation 3.3 for any two components, it can be thought that these two components might be merged. To evaluate this numerically, we define the following merge criterion:²

$$J_{merge}(i, j; \Theta^*) = \mathbf{P}_i(\Theta^*)^T \mathbf{P}_i(\Theta^*), \qquad (3.12)$$

where $\mathbf{P}_i(\Theta^*) = (P(i|\boldsymbol{x}_1; \Theta^*), \dots, P(i|\boldsymbol{x}_N; \Theta^*))^T \in \mathcal{R}^N$ is an *N*-dimensional vector consisting of the posterior probabilities for the *i*th model. *T* denotes

$$J_{merge}(i, j; \Theta^*) = \frac{\mathbf{P}_i(\Theta^*)^T \mathbf{P}_j(\Theta^*)}{\|\mathbf{P}_i(\Theta^*)\| \|\mathbf{P}_j(\Theta^*)\|}$$

can be used. In our experiments, however, we used equation 3.12 for simplicity, and the merge criterion did work well, as shown later.

 $^{^2\,}$ This merging criterion favors merges with larger classes. To avoid this bias, a modified criterion,

the transpose operation. $\|\cdot\|$ denotes the Euclidean vector norm. Clearly, two components ω_i and ω_j with large $J_{merge}(i, j; \Theta^*)$ are good candidates for a merge.

3.3.2 *Split Criterion*. As a split criterion (J_{split}), we define the local Kullback divergence as

$$J_{split}(k; \Theta^*) = \int f_k(\boldsymbol{x}; \Theta^*) \log \frac{f_k(\boldsymbol{x}; \Theta^*)}{p_k(\boldsymbol{x}; \theta_k^*)} d\boldsymbol{x}, \qquad (3.13)$$

which is the distance between two distributions: the local data density $f_k(x)$ around the *k*th model and the density of the *k*th model specified by the current parameter estimate Θ^* . The local data density is defined as

$$f_k(\boldsymbol{x}; \Theta^*) = \frac{\sum_{n=1}^N \delta(\boldsymbol{x} - \boldsymbol{x}_n) P(k | \boldsymbol{x}_n; \Theta^*)}{\sum_{n=1}^N P(k | \boldsymbol{x}_n; \Theta^*)}.$$
(3.14)

This is a modified empirical distribution weighted by the posterior probability so that the data around the *k*th model are focused on. Note that when the weights are equal, that is, $P(k|x; \Theta^*) = 1/M$, equation 3.14 is the usual empirical distribution:

$$p_k(\boldsymbol{x}; \Theta^*) = \frac{1}{N} \sum_{n=1}^N \delta(\boldsymbol{x} - \boldsymbol{x}_n).$$
(3.15)

Since it can be thought that the model with the largest $J_{split}(k; \Theta^*)$ has the worst estimate of the local density, we should try to split it.

The split criterion defined by equation 3.13 can be viewed as a likelihood ratio test. That is, $f_k(\boldsymbol{x}; \Theta^*)/p_k(\boldsymbol{x}; \theta_k^*)$ can be interpreted as the likelihood ratio test statistic. In this sense, our split criterion is similar to a cluster validity test formula proposed by Wolfe (1970).

3.3.3 Sorting Candidates. Using J_{merge} and J_{split} , we sort the split-andmerge candidates as follows. First, the merge candidates are sorted based on J_{merge} . Then, for each sorted merge candidate $\{i, j\}_c$, the split candidates, excluding $\{i, j\}_c$, are sorted as $\{k\}_c$. By combining these results and renumbering them, we obtain $\{i, j, k\}_c$, c = 1, ..., M(M - 1)(M - 2)/2.

4 Application to Density Estimation by Mixture of Gaussians ____

4.1 Synthetic Data. We apply the proposed algorithm to a density estimation problem using a mixture of gaussians. In this case, $p_m(x; \theta_m)$ on the right-hand side of equation 3.1 becomes

$$p_m(x; \theta_m) = (2\pi)^{-d/2} \det(\Sigma_m)^{-1/2}$$

$$\times \exp\left\{-\frac{1}{2}(\boldsymbol{x}-\boldsymbol{\mu}_m)^T\boldsymbol{\Sigma}_m^{-1}(\boldsymbol{x}-\boldsymbol{\mu}_m)\right\}.$$
(4.1)

Clearly, θ_m corresponds to mean vector $\boldsymbol{\mu}_m$ and covariance matrix $\boldsymbol{\Sigma}_m$.

In the case of density estimation by a mixture of gaussians, as is well known, the global maxima of the likelihood correspond to singular solutions with zero covariances. Therefore, the SMEM algorithm may converge to the singular solutions. To prevent this, we used a Bayesian regularization method (Ormoneit & Tresp, 1996). That is, the update equation for the covariance matrix is given by

$$\Sigma_m^{(t+1)} = \frac{\sum_{\boldsymbol{x} \in \mathcal{X}} (\boldsymbol{x} - \boldsymbol{\mu}_m^{(t+1)}) (\boldsymbol{x} - \boldsymbol{\mu}_m^{(t+1)})^T P(m | \boldsymbol{x}; \Theta^{(t)}) + \lambda I_d}{\sum_{\boldsymbol{x} \in \mathcal{X}} P(m | \boldsymbol{x}; \Theta^{(t)}) + 1},$$

$$m = 1, \dots, M.$$
(4.2)

Here I_d is the *d*-dimensional unit matrix, and λ is a regularization constant determined by some validation data. In the experiment we set $\lambda = 0.1$.

First, we used the two-dimensional synthetic data in Figure 2 to demonstrate visually the usefulness of the split-and-merge operations. The initial mean vectors and covariance matrices were, as shown in Figure 2b, set to near the means of all of the data and unit matrices, respectively. The usual EM algorithm converged to the local maximum solution shown in Figure 2c, whereas the SMEM algorithm converged to the superior solution shown in Figure 2f, very close to the true one. The split of the first gaussian shown in Figure 2d appeared to be redundant, but as shown in Figure 2e they are successfully merged, and the original two gaussians were improved. This indicates that the split-and-merge operations not only appropriately assign the number of gaussians in a local data space, but can also improve the gaussian parameters themselves.

4.2 Real Data. Next, we tested the proposed algorithm using 20-dimensional real data (facial images processed into feature vectors) where the local maxima made the optimization difficult (see Ueda & Nakano, 1998, for the details.) The data size was 103 for training and 103 for test. We ran three algorithms (EM, DAEM, and SMEM) for 10 different initializations using the *k*-means clustering algorithm. We set M = 5 and used a diagonal covariance for each gaussian. Table 1 shows the summary statistics (mean, standard deviation (std), maximum, and minimum) of log-likelihood values per sample size obtained by each of the EM, DAEM, and SMEM algorithms for 10 different initializations. As shown in Table 1, even the worst solution found by the SMEM algorithm was better than the best solutions found by the other algorithms on both the training and test data. Moreover, the log-likelihoods achieved by the SMEM algorithm.



Figure 2: Results by the EM and SMEM algorithms for a two-dimensional gaussian mixture density estimation problem. (a) Contours of true gaussians, (b) initial density, (c) result by the EM algorithm, (d)–(e) examples of split and merge operations by the SMEM algorithm, and (f) the final result.

		Initial value	EM	DAEM	SMEM
Training	mean	-159.1	-148.2	-147.9	-145.1
	std	1.77	0.24	0.04	0.08
	max	-157.3	-147.7	-147.8	-145.0
	min	-163.2	-148.6	-147.9	-145.2
Test	mean	-168.2	-159.7	-159.8	-155.9
	std	2.80	1.00	0.37	0.09
	max	-165.5	-158.0	-159.6	-155.9
	min	-174.2	-160.8	-159.8	-156.0

Table 1: Log-Likelihood/Sample Size.

Figure 3 shows log-likelihood value trajectories accepted in step 3 of the SMEM algorithm during the estimation process. The dotted lines in Figure 3 denote the starting points of step 2. Note that it is due to the initialization in step 3 that the log-likelihood decreases just after the splitand-merge. Comparing the convergence points at step 3 marked by the " \circ " symbol in Figure 3, one can see that the successive split-and-merge operations improved the log-likelihood for both the training and test data,



Figure 3: Trajectories of log-likelihood. The upper (lower) result corresponds to the training (test) data.

as we expected. Table 2 compares the number of iterations executed by the three algorithms. Note that in the SMEM algorithm, the number includes not only partial and full EM steps for accepted operations, but also EM-steps for rejected ones. From Table 2, the SMEM algorithm was about 8.7 times slower than the original EM algorithm. The average rank of the accepted split-and-merge candidates was 1.8 (std = 0.9), which indicates that the proposed split-and-merge criteria worked very well.

5 Application to Dimensionality Reduction Using Mixture of Factor Analyzers

5.1 Factor Analyzers. A single factor analyzer (FA) (Anderson, 1984) assumes that an observed *p*-dimensional variable *x* is generated as a linear transformation of some lower *q*-dimensional latent variable $z \sim \mathcal{N}(0, \mathbf{I})$

	EM	DAEM	SMEM
mean	47	147	409
std	16	39	84
max	65	189	616
min	37	103	265

Table 2: The Number of EM-Steps.

plus additive gaussian noise $v \sim \mathcal{N}(\mathbf{0}, \Psi)$. Ψ is a diagonal matrix. That is, the generative model can be written as

$$\boldsymbol{x} = \boldsymbol{W}\boldsymbol{z} + \boldsymbol{v} + \boldsymbol{\mu}. \tag{5.1}$$

Here, $\mathbf{W} \in \mathcal{R}^{p \times q}$ is a transformation matrix and is called a *factor loading matrix*. $\boldsymbol{\mu}$ is a mean vector. Then, from a simple calculation, the pdf of the observed data by an FA model can be obtained by

$$p(\boldsymbol{x}; \Theta) = \mathcal{N}(\boldsymbol{\mu}, \mathbf{W}\mathbf{W}^T + \Psi).$$
(5.2)

Noting that **W** can be represented by linearly independent column vectors as $\mathbf{W} = [w_1, ..., w_q]$, where $w_i \in \mathbb{R}^{p \times 1}$, we can rewrite equation 5.1 as

$$\boldsymbol{x} = \sum_{i=1}^{q} z_i \boldsymbol{w}_i + \boldsymbol{v} + \boldsymbol{\mu}, \tag{5.3}$$

where z_i is the *i*th component of z. Clearly, equation 5.3 shows that w_1, \ldots, w_q form the basis in a latent space. Hence, FA can be interpreted as a dimensionality reduction model extracting a linear manifold (affine subspace) $\mathcal{M} = L^{(q)} + \mu$ underlying the given observed data space, where $L^{(q)}$ denotes the linear subspace of \mathcal{R}^p spanned by the basis w_1, \ldots, w_q .

5.2 Mixture of Factor Analyzers. A mixture of factor analyzers (MFA), proposed by Ghahramani and Hinton (1997), is an extension of single FA. That is, MFA is defined as the combination of *M* mixture of FAs and can be thought of as a reduced dimensional mixture of gaussians. The MFA model extracts *q*-dimensional locally linear manifolds $\mathcal{M}_m = L_m^{(q)} + \boldsymbol{\mu}_m$ for $m = 1, \ldots, M$ underlying the given high-dimensional data. More intuitively, the MFA model can perform clustering and dimensionality reduction simultaneously. Since the MFA model extracts a globally nonlinear manifold, it is more flexible than the single FA model.

The pdf of the observed data by *M* mixtures of FAs is given by

$$p(\boldsymbol{x}; \Theta) = \sum_{m=1}^{M} \alpha_m \mathcal{N}(\boldsymbol{\mu}_m, \mathbf{W}_m \mathbf{W}_m^T + \boldsymbol{\Psi}_m).$$
(5.4)

(See Ghahramani & Hinton, 1997, for details.) Note that equation 5.4 is a natural extension of equation 5.2. The unknown parameters $\Theta = \{\alpha_m, \mu_m, W_m \mid m = 1, ..., M\}$ of the MFA model can be estimated by the EM algorithm. In this case, the complete data log-likelihood becomes

$$\mathcal{L}_{c} = \log \prod_{n=1}^{N} \prod_{m=1}^{M} p(\boldsymbol{x}_{n}, \boldsymbol{z}_{n}, m; \boldsymbol{\Theta}_{m})^{u_{m}}$$
$$= \sum_{n=1}^{N} \sum_{m=1}^{M} u_{m} \log p(\boldsymbol{x}_{n}, \boldsymbol{z}_{n}, m; \boldsymbol{\Theta}_{m}).$$
(5.5)



Figure 4: Results by the EM and SMEM algorithms for a one-dimensional manifold extraction problem. (a) True manifold and generated data, (b) initial estimate, (c) result by the EM algorithm, (d)–(e) examples of split and merge operations, and (f) the final result.

Here, u_m is a mixture indicator variable, where if x_n is generated by the *m*th model, $u_m = 1$; otherwise, $u_m = 0$. Since \mathcal{L}_c can be represented in the form of a direct sum, the *Q* function is also decomposable, and therefore the SMEM algorithm is straightforwardly applicable to the parameter estimation of the MFA model.

5.3 Demonstration. Figure 4 shows results of extracting a one-dimensional manifold from three-dimensional data (noisy shrinking spiral) using the EM and SMEM algorithms. The data points in Figure 4 were generated by

 $(X_1, X_2, X_3) = ((13 - 0.5t) \cos t, -(13 - 0.5t) \sin t, t) + additive noise,$ (5.6)

where $t \in [0, 4\pi]$. In this case, each factor loading matrix \mathbf{W}_m becomes a three-dimensional column vector corresponding to each thick line in Figure 4. The center position and the direction of each thick line are μ_m and \mathbf{W}_m , respectively. In addition, the length of each thick line is $2\|\mathbf{W}_m\|$.

Although the EM algorithm converged to poor local maxima as shown in Figure 4(c), the SMEM algorithm successfully extracted the data manifold shown in Figure 4(f).

Table 3 compares average log-likelihoods per data point over 10 different initializations. The log-likelihood values were drastically improved on both the training and test data by the SMEM algorithm.

Table 3: Log-Likelihood/Sample Size.

	EM	SMEM
Training	-7.68 (0.151)	-7.26 (0.017)
Test	-7.75 (0.171)	-7.33 (0.032)

5.4 Practical Applications.

5.4.1 Image Compression. An MFA model is available for block transform image coding. In this method, as in usual block transform coding approaches such as Karhunen-Lòeve transformation or principal component analysis (PCA), an image is subdivided into nonoverlapping blocks of $b \times b$ pixels. Typically, b = 8 or b = 16 is employed. Each block is regarded as a $d(=b \times b)$ -dimensional vector x. Let \mathcal{X} be a set of obtained x values. Then, using \mathcal{X} , an image is transformed by the following steps in the compression algorithm:

- 1. Set the desired dimensionality *q* and the number of mixture components *M*.
- 2. Estimate μ_m and W_m , for m = 1, ..., M by fitting an MFA model to \mathcal{X} .
- 3. For each $x \in \mathcal{X}$, compute

$$\hat{\boldsymbol{x}}^{(m)} = \boldsymbol{W}_m (\boldsymbol{W}_m^T \boldsymbol{W}_m)^{-1} \boldsymbol{W}_m^T (\boldsymbol{x} - \boldsymbol{\mu}_m) + \boldsymbol{\mu}_m,$$

for $m = 1, \dots, M.$ (5.7)

Then assign $\hat{x}^{(m^*)}$ as a reconstructed vector that minimizes the squared error $\|\hat{x}^{(m)} - x\|^2$ by its reconstruction.

4. Transform each $\hat{x}^{(m^*)}$ into a block image (reconstructed block image).

The derivation of equation 5.7 is as follows. The least-squares reconstructed vector \hat{x} is, as shown in Figure 5, obtained by orthogonally projecting $x - \mu$ onto $L_m^{(q)}$ and adding μ_m . That is,

$$\hat{\boldsymbol{x}} = P_m(\boldsymbol{x} - \boldsymbol{\mu}_m) + \boldsymbol{\mu}_m, \tag{5.8}$$

where P_m is the projection matrix of $L_m^{(q)}$ and is computed from

$$P_m = \mathbf{W}_m (\mathbf{W}_m^T \mathbf{W}_m)^{-1} \mathbf{W}_m^T.$$
(5.9)

Substituting equation 5.9 into equation 5.8, we obtain equation 5.7.

Figure 6 shows an example of compressed image by MFA with q = 4 and M = 10. Figure 6a shows a sample image used for the training data \mathcal{X} . We



Figure 5: Illustration of image reconstruction.

should use test images independent of the training image, but for simplicity we used the same image for the training and test. For comparison, we also tried the usual PCA-based image compression method (see Figure 6b). In the case of image compression by PCA, **W** is an orthogonal matrix composed of *q* column vectors (eigenvectors) corresponding to the *q* largest eigenvalues of the sample autocorrelation matrix:

$$S = \frac{1}{|\mathcal{X}|} \sum_{\boldsymbol{x} \in \mathcal{X}} \boldsymbol{x} \boldsymbol{x}^{T}.$$
(5.10)

Here, $|\mathcal{X}|$ denotes the number of components of \mathcal{X} . Since $\mathbf{W}\mathbf{W}^T = I_q$ (i.e., a set of column vectors { w_1, \ldots, w_q } is orthogonal), the reconstructed vector by PCA is given by

$$\hat{\boldsymbol{x}} = \boldsymbol{W}\boldsymbol{W}^T\boldsymbol{x}.$$
(5.11)

The major difference between the PCA and MFA approaches is that the former finds a global subspace from \mathcal{X} , while the latter simultaneously clusters the data and finds an optimum subspace for each cluster. Therefore, it is natural that the results by the MFA approaches (Figures 6c and 6d were much richer than that by the PCA approach—Figure 6b). Note that a mixture of PCA (MPCA) (Tipping & Bishop, 1997) is a much better model than PCA, but we have not compared it here because our purpose was to compare the SMEM algorithm with the EM algorithm for the MFA-based image compression.



(a) Original image



(c) MFA with EM



(b) PCA



(d) MFA with SMEM

Figure 6: An example of image reconstruction. (a) Original image, (b) Result by PCA, (c) Result by MFA model trained by the usual EM algorithm, and (d) Result by MFA model trained by the SMEM algorithm.

Comparing Figure 6c to Figure 6d, one can see that the quality of the reconstructed image using the SMEM algorithm is better than that using the EM algorithm. The mean squared errors per block of these constructed images were 15.8×10^3 , 10.1×10^3 , and 7.3×10^3 for Figures 6b, 6c, and 6d, respectively.

5.4.2 Application to Pattern Recognition. The MFA model is also applicable to pattern recognition tasks (Hinton, et al., 1997) since once an MFA model is fitted to each class, we can compute the posterior probability for each data point. More specifically, the pdf of class ω_i by the MFA model is

given by

$$p_i(\boldsymbol{x}; \Theta_i) = \sum_{m=1}^{M} P_{im} \mathcal{N}(\boldsymbol{\mu}_{im}, \mathbf{W}_{im} \mathbf{W}_{im}^T + \Psi_{im}), \qquad (5.12)$$

where Θ_i is a set of MFA model parameters for class ω_i . That is, $\Theta_i = \{P_{im}, \mu_{im}, \mathbf{W}_{im}, \Psi_{im} \mid m = 1, ..., M\}$. P_{im} is the mixing proportion of the *m*th model for class ω_i .

Using equation 5.12, the posterior probability of class ω_i given x is computed from

$$P(\omega_i | \boldsymbol{x}) = \frac{P_i p_i(\boldsymbol{x}; \Theta_i)}{\sum_{j=1}^{C} P_j p_j(\boldsymbol{x}; \Theta_j)}.$$
(5.13)

Here, *C* is the number of classes and P_i is the prior of class ω_i and is estimated from

$$P_i = \frac{N_i}{N},\tag{5.14}$$

where N_i is the number of training samples of class ω_i and N is the total sample size of all classes. Then, the optimum class i^* for x based on the Bayes decision rule is written as follows:

$$i^{*} = \arg \max_{i} P(\omega_{i} | \boldsymbol{x})$$

= $\arg \max_{i} N_{i} \sum_{m=1}^{M} P_{im} \mathcal{N}(\boldsymbol{\mu}_{im}, \mathbf{W}_{im} \mathbf{W}_{im}^{T} + \Psi_{im}).$ (5.15)

We compared the MFA model with another classification method based on clustering and dimensionality reduction, called the multiple subspace method (MSS) proposed by (Sugiyama and Ariki (1998). In order to define the MSS method we will first describe the simpler subspace (SS) method for classification known as CLAFIC (Oja, 1983). In the CLAFIC method, a linear subspace is extracted from the training data for each class, and then the distance between an input vector \boldsymbol{x} to be classified and its projected vector onto the linear subspace is computed for each class. Next, the input vector is classified as class ω_i^* with the minimum distance (see Oja, 1983, for details).

In the SS method, a single subspace is extracted for each class. On the other hand, in the MSS method, the data are first divided into several clusters by using the *k*-means algorithm. Then, for each cluster, the CLAFIC method is performed. This MSS method can be expected to improve the classification performance of the CLAFIC method, but due to the absence of a probability density model, it does not perform Bayes classification.



Figure 7: Recognition rates for hand-written digit data. (a) Results by MSS methods (M = 1 corresponds to the SS method). (b) Results by an MFA-based method trained by the EM algorithm. (c) Result by an MFA-based method trained by the SMEM algorithm.

We tried a digit recognition task (10 digits (classes)) using three methods: the MSS method, the MFA-based method with the EM algorithm, and the MFA-based method with the SMEM algorithm. The data were created using (degenerate) Glucksman's features (16-dimensional data) by NTT labs (Ishii, 1989). The data size was 200 per class for training and 200 per class for test.

The recognition accuracy values for the training and test data obtained by these methods are given in Figure 7. In Figure 7, q denotes the dimensionality of the latent space and M is the number of components in mixture. Note that the SS (CLAFIC) method corresponds to M = 1 in the MSS method shown in Figure 7a. Clearly, the MFA-based method with the SMEM algorithm consistently outperformed both the MSS method and the MFA-based method with the EM algorithm. The recognition accuracy by the 3-nearest neighbor (3NN) classifier was 88.3%. It is interesting that the MFA approach by the SMEM algorithm could outperform the nearest-neighbor approach when q = 3 and M = 5 (91.9%). This suggests that the intrinsic dimensionality of the data might be three or so.

6 Conclusion

We have shown how simultaneous split-and-merge operations can be used to move components of a mixture model from regions in a space in which there are too many components to regions in which there are too few. Such moves cannot be accomplished by methods that continuously move components through intermediate locations because the likelihoods are lower at these locations. A simultaneous split-and-merge can be viewed as a way of tunneling through low-likelihood barriers, thereby eliminating many nonglobal optima.

Note that the SMEM algorithm is applicable to a wide variety of mixture models, as long as decomposition 7 holds. To make the split-and-merge method more efficient, we have introduced criteria for deciding which splits and merges to consider and have shown that these criteria work well for low-dimensional synthetic data sets and higher-dimensional real data sets. Our SMEM algorithm consistently outperforms the standard EM algorithm, and therefore it can be very useful in practice.

In the SMEM algorithm, the split-and-merge operations are used to improve the parameter estimates within the maximum likelihood framework. However, by introducing probability measures over model, we could also use the split-and-merge operations to determine the appropriate number of components within the Bayesian framework. This extension is now in progress.

References _

- Anderson, T. W. (1984). *An introduction to multivariate statistical analysis* (2nd ed). New York: Wiley.
- Ball, G. H., & Hall, D. J. (1967). A clustering technique for summarizing multivariate data. *Behavioral Science*, 12, 153–155.
- Dempster, A. P., Laird, N. M., & Rubin, D. B. (1977). Maximum likelihood from incomplete data via the EM algorithm. *Journal of Royal Statistical Society B*, 39, 1–38.
- Ghahramani, Z., & Hinton, G. E. (1997). The EM algorithm for mixtures of factor analyzers (Tech. Report No. CRG-TR-96-1). Toronto: University of Toronto. Available online at: http://www.gatsby.ucl.ac.uk/~zoubin/papers/tr-96-1.ps.gz.
- Hinton, G. E., Dayan, P., & Revow, M. (1997). Modeling the minifolds of images of handwritten digits. *IEEE Trans. PAMI*, 8(1), 65–74.
- Ishii, K. (1989). Design of a recognition dictionary using artificially distorted characters. Systems and Computers in Japan, 21(9), 669–677.
- MacLachlan, G., & Basford, K. (1987). Mixture models: Inference and applications to clustering. New York: Marcel Dekker.
- Oja, E. (1983). Subspace methods of pattern recognition. Letchworth: Research Studies Press Ltd.
- Ormoneit, D., & Tresp, V. (1996). Improved gaussian mixture density estimates using Bayesian penalty terms and network averaging. In D. Touretzky, M. Moser, & M. Hasselmo (Eds.), *Neural information processing systems*, 8 (pp. 542–548). Cambridge, MA: MIT Press.
- Richardson, S., & Green, P. J. (1997). On Bayesian analysis of mixtures with an unknown number of components. J. R. Statist. Soc. B, 59(4), 731–792.
- Sugiyama, Y., & Ariki, Y. (1998). Automatic classification of TV sports news videos by multiple subspace method. *Trans. Institute of Electronics, Information* and Communication Engineers, 9 (pp. 2112–2119) (in Japanese).
- Tipping, M. E., & Bishop, C. M. (1997). Mixtures of probabilistic principal component analysers (Tech. Rep. No. NCRG-97-3). Birmingham, UK: Aston University.

- Ueda, N., & Nakano, R. (1994). A new competitive learning approach based on an equidistortion principle for designing optimal vector quantizers. *Neural Networks*, 7(8), 1211–1227.
- Ueda, N., & Nakano, R. (1998). Deterministic annealing EM algorithm. Neural Networks, 11(2), 271–282.
- Ueda, N., Nakano, R., Ghahramani, Z., & Hinton, G. E. (1999). SMEM algorithm for mixture models. In M. S. Keans, S. A. Solla, & D. A. Cohn (Eds.), *Neural information processing systems*, 11. Cambridge, MA: MIT Press.
- Wolfe, J. H. (1970). Pattern clustering by multivariate mixture analysis. *Multivariate Behavioral Research*, *5*, 329–350.

Received March 22, 1999; accepted September 20, 1999.