© Springer 2006

# BYY Harmony Learning on Finite Mixture: Adaptive Gradient Implementation and A Floating RPCL Mechanism

JINWEN MA\* and LE WANG

**Abstract.** In tackling the learning problem on a set of finite samples, Bayesian Ying-Yang (BYY) harmony learning has developed a new learning mechanism that makes model selection implemented either automatically during parameter learning or in help of evaluating a new class of model selection criteria. In this paper, parameter learning with automated model selection has been studied for finite mixture model via an adaptive gradient learning algorithm for BYY harmony learning on a specific bidirectional architecture (BI-architec-

evaluating a criterion incurs a large computational cost since we need to repeat the entire parameter estimating process at a number of different value of k, even though such a process is attempted to be organized in a more efficient way, e.g., embedding the checking of the criterion value within clustering as did in ISO-DATA [8].

In the field of artificial neural networks, competitive learning (CL) has been developed for clustering analysis and vector quantization [9]. However, the conventional competitive learning algorithms such as the classical competitive learning algorithm [10] and the frequency sensitive competitive learning algorithm [11], can only be considered as adaptive versions of k-means algorithm and thus are unable to solve the problem of selecting k during clustering analysis. Xu et al. [12] proposed a new kind of competitive learning algorithm, called rival penalized competitive learning (RPCL) algorithm, via such a mechanism that for each input, the winner of the units (i.e., weight vectors) is rewarded to adapt to the input, but the rival (the second winner) is penalized (or delearned) by a smaller learning rate. It was demonstrated that the RPCL algorithm has the ability of automatically allocating an appropriate number of units for a sample data set, with the other extra units being pushed far away from the sample data. Therefore, the RPCL algorithm can be used on a clustering analysis problem with an unknown k, as long as an initial number of the units is given to be larger than the number of actual component densities in the sample data set. In some extended versions of the RPCL algorithm [13], the units have been generalized to be the component densities of a finite mixture, which makes it possible for a mixture modeling problem too.

Alternatively, Bayesian Ying-Yang (BYY) harmony learning system and theory, proposed in 1995 [14] and systematically summarized in [15–17], has provided a theoretical foundation to solve this mixture modeling problem. The BYY harmony learning acts as a general statistical learning framework not only for understanding several existing major learning approaches but also for tackling the learning problem on a set of finite samples with a new learning mechanism that makes model selection implemented either automatically during parameter learning or in help of evaluating a new class of model selection criteria. Applied to the mixture modeling based problem, selection of k is made via either evaluating a criterion derived from the harmony measure or directly maximizing the harmony measure. In the latter(la)29(ttert elaborated in [17] that the latter case actually implements a RPCL mechanism that has its penalizing per sample during learning. Externally, the local maximum problem may also be solved by combining a Maximum-Likelihood(ML)-equivalent learning with the harmony via a weighted sum [15] or introducing simulated annealing like techniques that make learning gradually switched from a ML-equivalent learning to the harmony learning [15, 21]. In implementation, both cases are equivalent to replacing the WTA competition in the harmony learning and the Bayesian posterior allocation in the ML learning by some average of the both [15, 21]. Another internal mechanism, that is suggested but without further exploration in [15], is using a bidirectional architecture (BI-architecture) with its Ying path for the finite mixture and with its Yang path for regularization. In [22, 23], some batch way gradient implementations of BYY harmony learning with a BI-architecture were made on Gaussian mixture and experiments have shown that a RPCL like mechanism also occurs during parameter learning with automated model selection.

In the current paper, we study the implementation of an adaptive version of the gradient learning algorithm on finite mixture for this purpose, and analyze the learning mechanism on how it links to the RPCL learning with a floating penalizing. The simulation experiments have shown that the adaptive gradient algorithm in the Gaussian mixture setting is efficient in comparison with the batch gradient learning algorithms. Moreover, it is successfully applied to classification of the Iris data and unsupervised color image segmentation.

In the sequel, the adaptive gradient learning algorithm is derived in Section 2. Its learning mechanism is further analyzed in Section 3. Several simulation and practical experiments are conducted in Section 4 to demonstrate the algorithm for Gaussian mixture. Finally, we conclude in Section 5.

## 2. Adaptive Gradient Learning Algorithm

## 2.1. BYY HARMONY FUNCTION ON FINITE MIXTURE

We begin with a brief description of a BI-architecture of the BYY system on which the harmony learning is equivalent to the parameter learning on the finite mixture model, and leave the details to Xu [15].

A BYY system describes each observation  $x \in \mathcal{X} \subset \mathbb{R}^d$  and its corresponding inner representation  $y \in \mathcal{Y} \subset \mathbb{R}^m$  via the two types of Bayesian decomposition of the joint density p(x, y) = p(x)p(y|x) and q(x, y) = q(x|y)q(y), being called Yang machine and Ying machine, respectively. Given a data set  $D_x = \{x_t\}_{t=1}^N$ , the task of learning on a BYY system consists of specifying all the aspects of p(y|x), p(x), q(x|y), q(y) with a harmony learning principle implemented by maximizing the functional

$$H(p||q) = \int p(y|x)p(x)ln[q(x|y)q(y)]dxdy - lnz_q,$$
(1)

where  $z_q$  is a regularization term that implements either data smoothing or normalization as mentioned previously in the introduction section.

If p(y|x) and q(x|y) are both parametric, i.e., from a family of probability densities with a parameter  $\theta$ , the BYY system is called to have a BI-architecture for short. For the finite mixture modeling, we use the following BI-architecture of the BYY system. The inner representation y is discrete, i.e.,  $y \in \mathcal{Y} = \{1, 2, ..., k\} \subset R$  with m = 1. Then, we let  $q(y = j) = \alpha_j \ge 0$  with  $\sum_{j=1}^k \alpha_j = 1$ . Also, we ignore the regularization term  $z_q$  (i.e., set  $z_q = 1$ ) and let p(x) be directly given by the empirical density  $p_0(x)$ :

$$p_0(x) = \frac{1}{N} \sum_{t=1}^{N} \delta(x - x_t),$$
(2)

where  $x \in \mathcal{X} = R^d$  and  $\delta(\cdot)$  TLBT666.55296.2(r)296ki4.8(r)2962.1(r)296krsec, cr)ETBT9.9BT9.9

Considering the constraints on  $\alpha_j$ , we use the following so-called softmax representation:

$$\alpha_j = \frac{e^{\beta_j}}{\sum_{i=1}^k e^{\beta_i}}, \quad j = 1, \dots, k,$$
(6)

where  $-\infty < \beta_1, \ldots, \beta_k < +\infty$ .

Then, we can get the following derivatives of  $J(\Theta_k)$  per sample  $x_t$  with respect to  $\beta_j$  and  $\theta_j$ :

$$\frac{\partial J_t(\Theta_k)}{\partial \beta_j} = \sum_{i=1}^k \frac{\partial J_t(\Theta_k)}{\partial U_i(x_t)} \frac{\partial U_i(x_t)}{\partial \beta_j}$$
$$= \frac{1}{q(x_t|\Theta_k)} \sum_{i=1}^k [1 - \sum_{l=1}^k (p(l|x_t) - \delta_{il}) ln U_l(x_t)] (\delta_{ij} - \delta_{ij}) [1 - \delta_{ij}] (\delta_{ij} - \delta_{ij}] (\delta_{ij} - \delta_{ij}) [1 - \delta_{ij}] (\delta_{ij} - \delta_{ij}) [1 - \delta_{ij}] (\delta_{ij} - \delta_{ij}] (\delta_{ij} - \delta_{ij}) [1 - \delta_{ij}] (\delta_{ij} - \delta_{ij}) [1 - \delta_{ij}] (\delta_{ij} - \delta_{ij}] (\delta_{ij} - \delta_{ij}) [1 - \delta_{ij}] (\delta_{ij} - \delta_{ij}] (\delta_{ij} - \delta_{ij}) [1 - \delta_{ij}] (\delta_{ij} - \delta_{i$$

For the following experiments, we detail the adaptive gradient learning algorithm for Gaussian mixture. In this case,  $q(x|\theta_j)$  is given by

$$q(x|\theta_j) = q(x|m_j, \Sigma_j) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma_j|^{\frac{1}{2}}} e^{-\frac{1}{2}(x-m_j)^T \Sigma_j^{-1}(x-m_j)},$$
(13)

where  $m_j$  is the mean vector and  $\Sigma_j$  is the covariance matrix which is positive definite. Its derivatives take the following form:

$$\frac{\partial q(x_t|m_j, \Sigma_j)}{\partial m_j} = q(x_t|m_j, \Sigma_j) \Sigma_j^{-1}(x_t - m_j), \qquad (14)$$

$$\frac{\partial q(x_t|m_j, \Sigma_j)}{\partial \Sigma_j} = \frac{1}{2} q(x_t|m_j, \Sigma_j) [\Sigma_j^{-1} (x_t - m_j) (x_t - m_j)^T - I] \Sigma_j^{-1},$$
(15)

where *I* is the *d*-dim identity matrix. Substituting the derivative of  $q(x_t|m_j, \Sigma_j)$  with respect  $m_j$  into Equation (11), we then have the following adaptive update rule of  $m_j$ :

$$\Delta m_j = \eta p(j|x_t) \lambda_j(t) \Sigma_j^{-1}(x_t - m_j).$$
<sup>(16)</sup>

However, since  $\Sigma_j$  is constrained to be positive definite, we cannot guarantee its positive definite in the adaptive iteration if the derivative of  $q(x_t|m_j, \Sigma_j)$  with respect to  $\Sigma_j$  is substituted into Equation (11) directly. Instead, we use the decomposition technique suggested in [23], that is, letting  $\Sigma_j = B_j B_j^T$ , where  $B_j$  is a nonsingular square matrix, we have the following adaptive update rule of  $B_j$ :

$$\Delta \operatorname{vec} B_j = \frac{\eta}{2} p(j|x_t) \lambda_j(t) \frac{\partial (B_j B_j^T)}{\partial B_j} \operatorname{vec} [\Sigma_j^{-1} (x_t - m_j) (x_t - m_j)^T \Sigma_j^{-1} - \Sigma_j^{-1}],$$
(17)

where vec[A] denotes the vector obtained by stacking the column vectors of the matrix A, and

$$\frac{\partial (BB^T)}{\partial B} = I_{d \times d} \bigotimes B_{d \times d}^T + E_{d^2 \times d^2} \cdot B_{d \times d}^T \bigotimes I_{d \times d},$$

where  $\bigotimes$  denotes the Kronecker product (or tensor product), and

$$E_{d^{2}\times d^{2}} = \frac{\partial B^{T}}{\partial B} = (\Gamma_{ij})_{d^{2}\times d^{2}} = \begin{pmatrix} \Gamma_{11} & \Gamma_{12} & \cdots & \Gamma_{1d} \\ \Gamma_{21} & \Gamma_{22} & \cdots & \Gamma_{2d} \\ \cdots & \cdots & \cdots & \cdots \\ \Gamma_{d1} & \Gamma_{d2} & \cdots & \Gamma_{dd} \end{pmatrix}_{d^{2}\times d^{2}},$$

where  $\Gamma_{ij}$  is a  $d \times d$  matrix whose (j, i)th element is just 1, with all the other elements being zero. Combining Equations (10), (16), and (17) together, we get the specific adaptive gradient learning algorithm for Gaussian mixture in which  $\alpha_j$  and  $\Sigma_j$  are adaptively updated accordingly.

## 3. Analysis of the Adaptive Gradient Rule

In this section, we analyze the learning mechanism implied in the adaptive gradient learning rule given by Equations (10) and (11). First, we give some properties of  $\lambda_i(t)$  via the following theorem:

 $x_t$ 

#### THEOREM 1.

(i)  $p(j_1|x_t) \leq p(j_2|x_t) \leq \cdots \leq p(j_k|x_t)$  { $j_1, j_2, \dots, j_k$ } { $\{1, 2, \dots, k\}$   $\lambda_{j_1}(t) \leq \lambda_{j_2}(t) \leq \cdots \leq \lambda_{j_k}(t)$  $T(t) = e^{-(1+H(x_t))}$   $p(j|x_t) > (\leq)T(t)$   $\lambda_j(t) > (\leq)0$ ,

$$H(x_t) = H(p(1|x_t), \dots, p(k|x_t)) = -\sum_{i=1}^k p(i|x_t) ln p(i|x_t);$$
(18)

(ii)  $j_c = argmaxp(j|x_t)$   $\lambda_{j_c}(t) > 0$ 

(iii) 
$$p(j|x_t) > \frac{1}{e}$$
,  $\lambda_j(t) > 0$   $p(j|x_t) \leq \frac{1}{ke}$ ,  $\lambda_j(t) < 0$ ,  $e$ 

See the Appendix for the proof.

That is,  $\lambda_j(t)$  is floating with  $p(j|x_t)$  in the same descent order at each sample  $x_t$ and there exists a threshold value T(t) such that if  $p(j|x_t) > (\leq)T(t)$ ,  $\lambda_j(t) > (\leq)0$ . Obviously, T(t) varies with the sample  $x_t$  via the Shannon entropy of the posterior probabilities  $p(j|x_t)$  of the components in the finite mixture. As  $H(x_t)$  is high, i.e., the belonging component of  $x_t$  is obscure, T(t) becomes low; otherwise, as  $H(x_t)$ is low, i.e., the belonging component of  $x_t$  is clear, T(t) becomes high.

Moreover, when T(t) is low, there are generally several components with  $\lambda_j(t) > 0$ ; otherwise, when T(t) is high, there are a few components or just a single component with  $\lambda_j(t) > 0$ . If component j is the maximum one at the sample  $x_t$ ,  $\lambda_j(t)$  must be positive. On the other hand, if component j is relatively very small,  $\lambda_j(t)$  becomes negative.

We consider the incremental of component  $U_j(x_t)$  by the update of the adaptive gradient leaning rule and have

$$\Delta U_j(x_t) = q(x_t|\theta_j) \sum_{i=1}^k \frac{\partial \alpha_j}{\partial \beta_i} \Delta \beta_i + \alpha_j \left[ \frac{\partial q(x_t|\theta_j)}{\partial \theta_j} \right]^T \Delta \theta_j.$$
(19)

It follows from Equations (10) and (11) that

$$\Delta U_{j}(x_{t}) = \eta U_{j}(x_{t}) \Big[ (\lambda_{j}(t)p(j|x_{t}) - \sum_{i=1}^{k} \alpha_{i}\lambda_{i}(t)p(i|x_{t})) + \Big(\sum_{i=1}^{k} \alpha_{i}^{2} - \alpha_{j}\Big) \\ + \frac{\alpha_{j}\lambda_{j}(t)}{q(x_{t}|\Theta_{k})q(x_{t}|\theta_{j})} \Big| \Big| \frac{\partial q(x_{t}|\theta_{j})}{\partial \theta_{j}} \Big| \Big|^{2} \Big].$$

$$(20)$$

For analysis, we let

$$\gamma_j(t) = \gamma_{j1}(t) + \gamma_{j2}(t),$$
(21)

where

$$\gamma_{j1}(t) = \lambda_j(t) p(j|x_t) - \sum_{i=1}^k \alpha_i \lambda_i(t) p(i|x_t);$$
(22)

$$\gamma_{j2}(t) = \sum_{i=1}^{k} \alpha_i^2 - \alpha_j.$$
 (23)

We then have

$$\Delta U_j(x_t) = \eta U_j(x_t)\gamma_{j1}(t) + \eta U_j(x_t)\gamma_{j2}(t) + \eta \frac{\alpha_j^2 \lambda_j(t)}{q(x_t|\Theta_k)} \left\| \frac{\partial q(x_t|\theta_j)}{\partial \theta_j} \right\|^2.$$
(24)

According to Equation (22) and (i), it can be found that  $\gamma_{j1}(t)$  also varies with  $p(i|x_t)$  and there is another threshold value T'(t) such that if  $p(i|x_t) > (\leq)T'(t)$ ,  $\gamma_{j1}(t) > (\leq)0$ . Since  $\sum_{i=1}^{k} \alpha_i \lambda_i(t) p(i|x_t)$  is a weighted sum of these  $\lambda_i(t) p(i|x_t)$ ,  $\gamma_{j1}(t)$  becomes positive if  $p(j|x_t)$  is large enough. Therefore,  $\gamma_{j1}(t)$ behaves similarly as  $\lambda_i(t)$ , but their threshold values for  $p(j|x_t)$  may be different.

On the other hand,  $\gamma_{j2}(t)$  is only related to the mixing proportions  $\alpha_1, \ldots, \alpha_k$ . It follows from Equation (23) that if  $\alpha_j$  is larger than a threshold value  $\hat{\alpha} = \sum_{i=1}^k \alpha_i^2$ ,  $\gamma_{j2}(t)$  will be negative; otherwise, it will be positive.

With the above preparations, we can uncover the competitive learning mechanism among the mixture components  $U_1, \ldots, U_k$  by the adaptive gradient learning rule. When a sample  $x_t$  is inputted, the component parameters  $\beta_j, \alpha_j, \theta_j$  are modified according to Equations (6), (10), and (11), respectively. Accordingly, component  $U_j$  is modified and its gain<sup>1</sup> consists of three parts:  $\eta U_j(x_t)\gamma_{j1}(t), \eta U_j(x_t)\gamma_{j2}(t)$ , and  $\eta \frac{\alpha_j^2 \lambda_j(t)}{q(x_t |\Theta_j)} || \frac{\partial q(x_t |\theta_j)}{\partial \theta_j} ||^2$ .

Since both the descent queuing orders of  $\gamma_{j1}(t)$  in  $\{\gamma_{j1}(t)\}_{j=1}^{k}$  and  $\lambda_{j}(t)$  in  $\{\lambda_{j}(t)\}_{j=1}^{k}$  are that of  $p(j|x_t)$  in  $\{p(j|x_t)\}_{j=1}^{k}$ , the first and last parts of the gain can be both considered as a result of the competition among the components  $U_1, \ldots, U_k$  at the sample  $x_t$  through the posterior probabilities  $p(1|x_t), \ldots, p(k|x_t)$ . By each of these two parts, all the posterior probabilities can be put in a descent order and there is a threshold value T(t) such that if  $p(j|x_t) > T(t)$ , the part of the gain is positive and  $U_j$  becomes a winner and is rewarded to increase; otherwise, if  $p(j|x_t) < T(t)$ , the part of the gain is negative and  $U_j$  becomes a loser and is penalized to decreased. Because  $U_j(x_t) = q(x_t|\Theta_k)p(j|x_t)$ , the competition of the posterior probabilities is equivalent to that of the components themselves.

<sup>&</sup>lt;sup>1</sup>Since  $\eta$  is very small, we neglect its quadratic or higher order terms in the incremental of  $U_j$  at the sample  $x_t$ .

By the middle part of the gain, we further find a proportion conscience mechanism implied in the adaptive gradient learning rule. Clearly,  $\hat{\alpha} = \sum_{i=1}^{k} \alpha^2$ 

algorithm is applied to classification of the Iris data and unsupervised color image segmentation.

## 4.1. SIMULATION EXPERIMENTS

4.1.1.

We begin with a brief description of the four sets of sample data used for our simulation experiments. We conducted four Monte Carlo experiments in which samples were drawn from a mixture of four or three bivariate Gaussians distributions on the plane coordinate system (i.e., d = 2). As shown in Figure 1, each data set of samples are generated with a certain degree of overlap among the clusters (Gaussians) in the mixture. They are four typical sets of sample data from Gaussian mixtures. The clusters in  $S_1$  are sphere-shaped, with the equal number of samples. But those in  $S_2$  are ellipse-shaped, with different numbers of samples.



## BYY HARMONY LEARNING ON FINITE MIXTURE

Moreover,  $S_3$  consists of three very flat clusters and  $S_4$  has a small number of samples, with the same structure of  $S_2$ . The detailed parameters for these four sets of sample data are given in Table I where  $m_i$ ,  $\Sigma_i = [\sigma_{jk}^i]_{2\times 2}$ ,  $\alpha_i$  and  $N_i$  denote the mean vector, covariance matrix, mixing proportion, and the number of samples of the *i*th Gaussian(cluster), respectively.

# 4.1.2.

In our simulation experiments and application experiments in the sequential subsections, k is always selected to be larger than  $k^*$ , i.e., the number of actual Gaussians in a sample or real data set. As mentioned in Section 1, we do not know the true number  $k^*$  of actual Gaussians in a data set. Generally, we can overestimate it with some additional information in the data set. If this can be done, we can let the overestimate of  $k^*$  be k so that  $k > k^*$ . Certainly, we can always choose a large number as k such that it is surely larger than  $k^*$  of the data set. But this may not be desirable since a large  $k(>>k^*)$  will increase the training time and the risk of model selection error (see the following simulation analysis). Fortunately, an efficient procedure via the RPCL algorithm was suggested in [12] to get an appropriate k for a data set. We can start with a small number or a guess for k and perform the RPCL algorithm on the data set. By checking whether there are spare units (i.e., far-away weight vectors) left, we would know if  $k > k^*$ . If yes,

The sample set	Gaussian	m <sub>i</sub>	$\sigma^i_{11}$	$\sigma^i_{12}$	$\sigma_{22}^i$	$lpha_i$	$N_i$
$\mathcal{S}_1$	Gaussian 1	(2.50, 0)	0.50	0	0.50	0.25	400
(N = 1600)	Gaussian 2	(0, 2.50)	0.50	0	0.50	0.25	400
	Gaussian 3	(-2.50, 0)	0.50	0	0.50	0.25	400
	Gaussian 4	(0, -2.50)	0.50	0	0.50	0.25	400
$\mathcal{S}_2$	Gaussian 1	(2.50, 0)	0.45	-0.25	0.55	0.34	544
(N = 1600)	Gaussian 2	(0, 2.50)	0.65	0.20	0.25	0.28	448
	Gaussian 3	(-2.50, 0)	1	0.10	0.35	0.22	352
	Gaussian 4	(0, -2.50)	0.30	0.15	0.80	0.16	256
$\mathcal{S}_3$	Gaussian 1	(2.50, 0)	0.10	-0.20	1.25	0.50	600
(N = 1200)	Gaussian 2	(0, 2.50)	1.25	0.35	0.15	0.30	360
· · · · ·	Gaussian 3	(-1, -1)	1	-0.80	0.75	0.20	240
$\mathcal{S}_4$	Gaussian 1	(2.50, 0)	0.28	-0.20	0.32	0.34	68
(N = 200)	Gaussian 2	(0, 2.50)	0.34	0.20	0.22	0.28	56
	Gaussian 3	(-2.50, 0)	0.50	0.04	0.12	0.22	44
	Gaussian 4	(0, -2.50)	0.10	0.05	0.50	0.16	32
		(., 100)		2.00			

The parameters in the actual Gaussian mixtures to the three sets of sample data.

then take this value; if not, increase k by adding a positive increment and perform the RPCL algorithm again. The same procedure can be repeated until  $k > k^*$ . In this way, we can get a moderately large number as k of the algorithm for the data set.

The other parameters in the adaptive gradient learning algorithm can be initialized randomly within certain intervals under the constraints. Actually, the initial value of  $\beta_j$  can be freely selected from some interval. But it has been found by the experiments that when the initial values of these  $\beta_j$  are equal or close and thus the initial values of the corresponding  $\alpha_j$  are equal or close, the adaptive gradient learning algorithm converges more efficiently. Thus, we give the same initial value or close values to these  $\beta_j$ . As for  $m_j$ , we can select its initial value in the field of the sample data or some sample point. However, the initial matrix of  $B_j$  should be nonsingular. For convenience, we can let the initial matrix of each  $B_j$  be the identity matrix or around it. The learning rate  $\eta$  is given by  $\eta(n) =$  $\eta_0/n$  with  $\eta_0 = 0.1$  for the simulation experiments. The algorithm is stopped when  $|J(\Theta_k^{\text{new}}) - J(\Theta_k^{\text{odd}})| < 10^{-7}$ .

# 4.1.3.

We implemented the adaptive gradient learning algorithm on these four typical sets of sample data. The experimental results on  $S_1$  and  $S_2$  are given in Figures 2 and 3, respectively, with case k=8 and  $k^*=4$ . We observe that four Gaussians are finally located accurately in either the sphere-shaped and symmetric structure case or the ellipse-shaped and asymmetric structure case, while the mixing proportions of the other four Gaussians were reduced to below 0.0001, i.e, these Gaussians are extra and can be discarded. Therefore, the correct number of the clusters have been detected on these data sets. Moreover, the experiment has been made on  $S_3$  with  $k=8, k^*=3$ . As shown in Figure 4, clusters are far from spherical shapes (actually they are very flat). Again, three Gaussians are located accurately, while the mixing proportions of the other five extra Gaussians become less than 0.001. That is, the correct number of the clusters can still be detected on such a special data set. Furthermore, the adaptive gradient learning algorithm was also implemented on  $S_4$ with  $k = 8, k^* = 4$ . As shown in Figure 5, even each cluster has a small number of samples, the correct number of clusters can still be detected, with the mixing proportions of other four extra Gaussians reduced below 0.0002. Therefore, the adaptive gradient learning algorithm can make model selection automatically during the parameter learning in each of these sample data sets.

The further experiments of the adaptive gradient algorithm had been also made on these sets of sample data in the similar cases. Actually, a failure on the correct number detection rarely happen when we initially set  $k^* \le k \le 3k^*$ . However, the adaptive gradient learning may lead to a wrong detection when  $k > 3k^*$ . On the other hand, an extra Gaussian may be stable with any shape as its mixing

## ING ON FINITE MIXTURE



result of the adaptive gradient algorithm on the sample set  $S_1$  (after us and the following three figures, the contour lines of each Gaussian are less than  $e^{-3}$ (peak).



\_\_\_\_

# JINWEN MA AND LE WANG



The average errors of the estimated parameters by the adaptive gradient learning and EM algorithms on the four data sets.

S	$\mathcal{S}_1$	$\mathcal{S}_2$	$\mathcal{S}_3$	$\mathcal{S}_4$
AGL algorithm	0.027539	0.025610	0.034658	0.016092
EM algorithm	0.019434	0.015928	0.029762	0.008623

proportion attenuates to zero. Sometimes, it degenerates into a 1-dim Gaussian taking a shape of a line.

## 4.1.4.

In addition to the automated model selection, we further compared the converged values of parameters (discarding the extra Gaussians) with those parameters in the true Gaussian mixture where the samples come from. The performance metric of average error is adopted as the absolute deviation between the estimated parameters and the true parameters in each case. In Table II, we give the average errors of the adaptive gradient learning (AGL) algorithm. For comparison, we also give the average errors of the EM algorithms on these four sets of sample data under  $k = k^*$  in Table II.

From Table II, we can find that the adaptive gradient learning algorithm usually converges to a good estimate of the true parameters with a low average error. However, the EM algorithm slightly outperform the adaptive gradient learning algorithm on the parameter estimation. That is, there is a slight deviation between the BYY harmony learning estimate and the maximum likelihood estimate (obtained from the EM algorithm). It is clear that this deviation is caused by the floating RPCL mechanism of the adaptive gradient learning algorithm that each (2) The adaptive gradient learning algorithm converges to the correct solution more frequently than the batch gradient algorithms do in the cases that the overlap among the Gaussians in a sample data set becomes strong or k is much greater than  $k^*$ . That is, the adaptive gradient algorithm can search the global maximum of  $J(\Theta_k)$  more efficiently.

As a result, the adaptive gradient learning algorithm is more efficient that the bath gradient learning algorithms for the practical applications, which is really demonstrated by the practical applications in the following two subsections. As compared with two typical methods of the maximum certainty partitioning using a number of kernel functions [25] and the variational Bayesian learning [26] for parameter learning and model selection on Gaussian mixture, our adaptive learning algorithm can lead to a similar or better result, but have much less computation.

## 4.2. CLASSIFICATION OF THE IRIS DATA

We further applied the adaptive gradient learning algorithm to classification (or recognition) of the Iris data [27] which is a typical real data set for testing a classification algorithm. Actually, it consists of 150 samples of three classes where each class contains 50 samples and each sample or datum is four-dimensional and consists of measures of the plants morphology. Although there is the class index for each sample, we could not use the class indexes of these samples in the adaptive gradient learning algorithm since it learns in an unsupervised way. However, the class indexes were used to check the classification result of the adaptive gradient learning algorithm on the Iris data.

The adaptive gradient learning algorithm was applied to the Iris data

## BYY HARMONY LEARNING ON FINITE MIXTURE

#### 4.3. UNSUPERVISED COLOR IMAGE SEGMENTATION

We finally applied the adaptive gradient learning algorithm to unsupervised color image segmentation that has been considered as a promising and challenging area in image processing [29]. Here, we utilized each Gaussian in the adaptive gradient learning algorithm to represent an object (including the background) in a color image; however, the number of actual objects is not given in advance. We set k to be larger than the true number  $k^*$  of the actual objects and the pixels in the image are partitioned according to the maximum posteriori probabilities among  $p(j|x_l)$ .

We used three typical color images of house, jellies, and peppers in the segmentation experiments. Each pixel in the image was represented by a three-dimensional real point. In our experiments, we first regularized all the three coordinates of the pixels in each color image via dividing them by 32 so that the regularized coordinates were within an interval of [0,8]. We then ran the adaptive gradient learning algorithm on the data sets of these three color images, respectively, by letting k = 6 (We could guess that the number of objects in a color image is less than 6) with the other initial parameters set by  $\beta_j = 1 + 0.1 \times \text{rand}(1)$ ,  $B_j = (1.1 + 0.1 \times \text{rand}(1)) \times I_3$ ,  $\eta_n = 0.0008/n$ . Again, we set a low threshold value T = 0.04 such that when some  $\alpha_j < T$ , we cancel the *j*th Gaussian in the mixture for the following learning iterations. The algorithm was also stopped when  $|J(\Theta_k^{\text{new}}) - J(\Theta_k^{\text{old}})| < 10^{-7}$ .

The first experiment was made on the color image of house which is shown in Figure 6(a). We implemented the adaptive gradient learning algorithm on this color image with the above initial settings and leaded to the segmentation results shown in Figure 6(b). It can be found that two objects (the house and the background) are finally located accurately, while the mixing proportions of the other four Gaussians (objects) are reduced to below 0.04, i.e. these objects are extra and discarded from the figure. That is, the correct number of the actual objects had been detected on the color image with an accurate segmentation. Moreover, the second experiment had been made on the color image of jellies, which is shown in Figure 7(a), with k = 6. As shown in Figure 7(b), two objects (the jellies and the background) are located accurately, while the mixing proportions of the other four extra Gaussians (objects) become less than 0.04. That is, the correct number of the objects could still be detected on this color image. Finally, the adaptive gradient learning algorithm was implemented on the color image of peppers, which is shown in Figure 8(a), with k = 6. As shown in Figure 8(b), the three objects are located accurately, with the mixing proportions of other three extra Gaussians reduced below 0.01.

As a result, the adaptive gradient learning algorithm can detect the number of actual objects in each of these color images. Moreover, the segmentation results of the adaptive gradient learning algorithm are much better than those of the generalized competitive clustering (GCC) algorithm [29] (based on the fuzzy clustering theory) given in the web http://www-rocq.inria.fr/~boujemaa/Partielle2.html.



. The segmentation result on the color image of house. (a) The original color image of house; (b) The segmented image of house via the adaptive gradient learning algorithm.



The segmentation result on the color image of jellies: (a) The original color image of jellies; (b) The segmented image of jellies via the adaptive gradient learning algorithm.



The segmentation result on the color image of peppers: (a) The original color image of peppers; (b) The segmented image of peppers via the adaptive gradient learning algorithm.

By comparison, we can easily find that the adaptive gradient learning algorithms lead to a more accurate segmentation on the contours of the objects in each image.

# 5. Conclusions

After introducing a BI-architecture of the BYY system and the harmony function, we have derived an adaptive gradient leaning algorithm to make parameter learning on finite mixture with automated model selection. It is shown by theoretic analysis that the adaptive gradient learning rule enforces a kind of floating RPCL mechanism such that for each input  $x_t$ , there is a threshold value T(t)such that the components (or posterior probabilities) above T(t) are rewarded to increase, while the components (or posterior probabilities) below T(t) are penalized to decrease, simultaneously with a conscience mechanism on the mixing proportions. Moreover, the adaptive gradient learning algorithm is demonstrated well by the simulation and practical results on the sample data sets from Gaussian mixtures with a certain degree of overlap.

## Acknowledgements

This work was supported by the Natural Science Foundation of China for Projects 60071004, 60471054. The authors thank Prof. Lei Xu for his strong support and helpful discussions, and Prof. Taijun Wang for his simulation supports.

# Appendix

$$\lambda_i(t)$$

(i) According to Equation (9), we have

$$\lambda_{j}(t) = 1 + lnU_{j}(x_{t}) - \sum_{i=1}^{k} p(i|x_{t}) lnU_{i}(x_{t})$$
  
= 1 + lnp(j|x\_{t}) -  $\sum_{i=1}^{k} p(i|x_{t}) lnp(i|x_{t})$   
= 1 + lnp(j|x\_{t}) + H(x\_{t}), (25)

where

$$H(x_t) = H(p(1|x_t), \dots, p(k|x_t)) = -\sum_{i=1}^k p(i|x_t) ln p(i|x_t)$$
(26)

is the Shannon entropy of the posterior probabilities  $p(i|x_t)$  at the sample  $x_t$ . Since  $H(x_t)$  is invariant with j and  $lnp(j|x_t)$  increases with  $p(j|x_t)$ , we certainly have  $\lambda_{j_1}(t) \leq \lambda_{j_2}(t) \leq \cdots \leq \lambda_{j_k}(t)$  if  $p(j_1|x_t) \leq p(j_2|x_t) \leq \cdots \leq p(j_k|x_t)$ . It also follows from Equation (25) that there exists a threshold value  $T(t) = e^{-(1+H(x_t))}$  for  $p(i|x_t)$  such that if  $p(i|x_t) > (\leq)T(t)$ ,  $\lambda_i(t) > (\leq)0$ .

(ii) If  $j_c = \operatorname{argmax} p(j|x_t)$ , i.e.,  $p(j|x_t) \leq p(j_c|x_t)$  for any j, and thus  $lnp(j|x_t) \leq lnp(j_c|x_t)$  for any j. Therefore, we have

$$\sum_{j=1}^{k} p(j|x_t) ln p(j|x_t) \leqslant \sum_{j=1}^{k} p(j|x_t) ln p(j_c|x_t) = ln p(j_c|x_t),$$
(27)

from which we further have

$$\lambda_{j_c}(x_t) \ge 1. \tag{28}$$

(iii) If  $p(j|x_t) > \frac{1}{e}$ , i.e.,  $ep(j|x_t) > 1$ , then  $ln[ep(j|x_t)] = 1 + lnp(j|x_t) > 0$ . Since  $H(x_t) \ge 0$ , we have

$$\lambda_{j}(t) = 1 + lnp(j|x_{t}) + H(x_{t}) > 0.$$
<sup>(29)</sup>

If  $p(j|x_t) \leq \frac{1}{ke}$ , i.e.,  $ep(j|x_t) \leq \frac{1}{k}$ , then  $ln[ep(j|x_t)] = 1 + lnp(j|x_t) \leq -lnk$ . By the property of Shannon entropy,  $H(x_t) = H(p(1|x_t), \dots, p(k|x_t)) \leq lnk$  with the equality only holding when each  $p(i|x_t) = \frac{1}{k}$ , we have

$$\lambda_{j}(x_{t}) = 1 + lnp(j|x_{t}) + H(x_{t})$$

#### BYY HARMONY LEARNING ON FINITE MIXTURE

- 7. Bozdogan, H: Model selection and Akaike's information criterion: the general theory and its analytical extensions, **52** (1987), 345–370.
- 8. Ball, G. H. and Hall, D. J.: ISODATA: a novel method of data analysis and pattern classification, Technique Report No. 699616, 1965, Stanford Research International.
- 9. Makhoul, J., Rpucos, S. and Gish, H.: Vector quantization in speech coding, 73 (1985), 1551–1558.
- Rumelhart, D. E. and Zipser, D.: Feature discovery by competitive learning, 9 (1985), 75–112.
- Ahalt, S. C., Krishnamurty, A. K., Chen, P. and Melton, D. E.: Competitive learning algorithm for vector quantization, 3 (1990), 277–291.
- Xu, L., Krzyzak, A. and Oja, E.: Rival penalized competitive learning for clustering analysis, RBF net, and curve detection, 4 (1993), 636–648.
- Xu, L.: Rival penalized competitive learning, finite mixture, and multisets clustering, In: 3 (1998),
  - 251-2530.
- 14. Xu, L.: Ying-Yang machine: a Bayesian-Kullback scheme for unified learnings and new results on vector quantization, In:

**2** (1995), 977–988.

- Xu, L.: Best harmony, unified RPCL and automated model selection for unsupervised and supervised learning on Gaussian mixtures, three-layer nets and ME-RBF-SVM models, 11 (2001), 43–69.
- 16. Xu, L.: Ying-Yang learning, In: M. A. Arbib (2nd ed.),
  - pp. 1231–1237. The MIT Press, Cambridge, MA, 2002.
- 18. Xu, L.: Bayesian-Kullback YING-YANG learning scheme: reviews and new results, In:

- 19. Hu, X. and Xu, L.: A comparative study of several cluster number selection criteria, In: vol. 2690 pp. 195–202, 2003.
- 20. Liu, Z. and Xu, L.: Smoothed local PCA by BYY data smoothing learning, In:

, 621–622 (2001).

21. Ma, J., Wang, T. and Xu, L.: An annealing approach to BYY learning on Gaussian mixture with automated model selection, In:

( , , & , \_ 1 (2003), 23–28.

- 22. Ma, J., Wang, T., and Xu, L.: A gradient BYY harmony learning rule on Gaussian mixture with automated model selection, **56** (2004), 481–487.
- 23. Ma, J., Gao, B., Wang, Y., and Cheng Q.: Conjugate and natural gradient rules for BYY harmony learning on Gaussian mixture with automated model selection,

*fi* **19** (2005), 701–713.

- 24. Robbins, H. and Monro, S.: A stochastic approximation method, 22 (1951), 400–407.
- 25. Robert, S. J., Everson, R., and Rezek, I.: Maximum certainty data partitioning, **33** (2000), 833–839.
- 26. Ueda, N. and Ghahramani, Z.: Bayesian model search for mixture models based on optimizing variational bounds, 15 (2002), 1223–1241.

- 27. Blake, C. L. and Merz, C. J.: University of California, Irvine, Department of Information and Computer Science, 1998, http://www.ics.uci.edu/~mlearn/MLRepository.html.
- Loo, C. K., and Rao, M. V. C.: Accurate and reliable diagnosis and classification using probablistic ensemble simplified fuzzy ARTMAP, 17 (2005), 1589–1593.
- 29. Boujeman, N.: Generalized competitive clustering for image segmentation, In:

, 133–137 (2000).