

ENTROPY PENALIZED AUTOMATED MODEL SELECTION ON GAUSSIAN MIXTURE*

JINWEN MA

*Department of Information Science, School of Mathematical Sciences & LMAM,
Peking University, Beijing, 100871, P. R. China
jwmamath.pku.edu.cn*

TAIJUN WANG

*Department of Wireless Engineering, Southeast University,
Nanjing, 210018, P. R. China*

Gaussian mixture modeling is a powerful approach for data analysis and the determination of the number of Gaussians, or clusters, is actually the problem of Gaussian mixture model selection which has been investigated from several respects. This paper proposes a new kind of automated model selection algorithm for Gaussian mixture modeling via an entropy penalized maximum-likelihood estimation. It is demonstrated by the experiments that the proposed algorithm can make model selection automatically during the parameter estimation, with the mixing proportions of the extra Gaussians attenuating to zero. As compared with the BYY automated model selection algorithms, it converges more stably and accurately as the number of samples becomes large.

Keywords: Gaussian mixture; model selection; maximum-likelihood estimation; Shannon entropy; penalty.

1. Introduction

Since Gaussian mixture populations are very common in various fields of practical applications, Gaussian mixture modeling is a powerful approach for data analysis. Actually, considerable effort has been made on Gaussian mixture modeling as well as its applications for clustering analysis on a sample data set.¹⁵ Although there have been various statistical or unsupervised competitive learning methods to do such a task, e.g. EM algorithm¹⁷ for maximum-likelihood, k -means algorithm¹⁴ for the Mean Square Error (MSE) clustering and the self-organizing network for hyperellipsoidal clustering (HEC),⁹ it is usually assumed that the number of Gaussians, or clusters, in the data set is pre-known. However, in many instances this key information is not available and then the selection of an appropriate number of Gaussians must be made before or during the estimation of the parameters in the

*This work was supported by the Natural Science Foundation of China for Project 60071004.

mixture. Since the number of Gaussians is just the scale of the mixture model, its determination is actually the problem of model selection for Gaussian mixture modeling. In fact, it is a rather difficult problem.⁸

The earliest possible method for solving this model selection problem may be to choose the optimal number of Gaussians (or clusters) by the Akaike's information criterion² or its extensions.^{4,5,19} Similarly, several indexes of cluster validity have been also proposed to choose the optimal number (e.g. Refs. 3, 16 and 21). But the process of evaluating an information criterion or validity index incurs a large computational cost since we need to repeat the entire parameter estimation process at a number of different values 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 in ISODATA.

From the respect of neural networks, competitive learning (CL) has been developed for clustering analysis and vector quantization.⁶ However, the conventional competitive learning algorithms such as the classical competitive learning algorithm¹⁸ and the frequency sensitive competitive learning algorithm,¹ can be only considered as adaptive versions of k -means algorithm and thus they are unable to solve this model selection problem. In order to do so, Xu, Krzyzak and Oja proposed a new kind of competitive learning algorithm, called rival penalized competitive learning (RPCL) algorithm,²² which has the ability of automatically allocating an appropriate number of weight vectors for a sample data set, with the other extra ones being pushed far away from the sample data. Therefore, the RPCL algorithm can be used to solve the above model selection problem by choosing a number of the weight vectors which is surely larger than the number of actual clusters in the sample data set. In some extended versions of the RPCL algorithm,²⁴ the weight vectors have been generalized to Gaussian densities, which makes it more suitable for the RPCL algorithm to be applied to solving the model selection problem for Gaussian mixture. Recently, Ma, Wang and Xu proposed a cost function by which a generalized RPCL algorithm, called distance sensitive RPCL algorithm, has been proposed with a more powerful performance and easier implementation.¹⁰

Proposed in 1995²³ and systematically developed in the past years,²⁵⁻²⁷ Bayesian Ying-Yang (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 to be implemented automatically during parameter learning via a new class of model selection criteria. Specifically, the BYY harmony learning can be applied to the model selection problem for Gaussian mixture modeling. In fact, Ma, Wang and Xu have already implemented this new mechanism on a BI-architecture and a B-architecture of the BYY system via a gradient, iterative and two annealing algorithms to solve the model selection problem automatically during the parameter learning.^{11,13,20} It has been demonstrated by simulation experiments that the number of Gaussians can be automatically selected for the sample data set, with the mixing proportions of the extra Gaussians attenuating to zero.

In the current paper, based on the fact that reducing the Shannon entropy of the mixing proportions in a Gaussian mixture can lead to some of these proportions to degenerate to zero, i.e. to eliminate the corresponding Gaussians from the mixture, we proposed an entropy penalized maximum-likelihood estimation (MLE) iterative algorithm which can also make the parameter estimation with automated model selection. It is further demonstrated by the simulation experiments.

In the sequel, the entropy penalized maximum-likelihood estimation (MLE) iterative algorithm is derived in Sec. 2. In Sec. 3, several simulation experiments are conducted to demonstrate the proposed algorithm, with a comparison between it and those BYY model selection algorithms. A brief conclusion is given in Sec. 4.

Entropy Penalized MLE Iterative Algorithm

We consider the following Gaussian mixture model:

$$(x_j) = \sum_{j=1}^K \pi_j (x_j | \mu_j, \Sigma_j) \quad \pi_j \geq 0 \quad \sum_{j=1}^K \pi_j = 1 \quad (1)$$

where

$$(x_j | \mu_j, \Sigma_j) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma_j|^{\frac{1}{2}}} \exp\left\{-\frac{1}{2}(x-m_j)^T \Sigma_j^{-1} (x-m_j)\right\} \quad (2)$$

and where K is the number of Gaussians, x denotes a sample vector and d is the dimensionality of x . The parameter vector θ consists of the mixing proportions π_j , the mean vectors μ_j , and the covariance matrices $\Sigma_j = (\sigma_{pq}^{(j)})_{d \times d}$ which are assumed symmetric and positive definite.

In this mixture model, π_j is the proportion or probability of Gaussian (component) j in the mixture population. Then, the Shannon entropy of these mixing proportions, i.e. $H(\pi) = -\sum_{j=1}^K \pi_j \log \pi_j$, represents the uncertainty of the distribution of these mixing proportions. Here and hereafter, the base of the logarithm is always the natural number. According to the property of Shannon entropy, as K increases, $H(\pi)$ generally increases. Oppositely, as K decreases, some mixing proportions may be forced to be zero or a small number. That is, the corresponding Gaussians are essentially eliminated from the mixture. Thus, the number of actual Gaussians is decreased. According to this characteristic, we can penalize the MLE by adding certain minus Shannon entropy of the mixing proportions to the log maximum likelihood (ML) function so that the entropy penalized MLE is able to make the parameter estimation with automated model selection. In the following, we will realize this idea by proposing an entropy penalized ML objective function and establishing an iterative algorithm to maximize it.

Given a sample data set $S = \{x_t\}_{t=1}^N$ from a Gaussian mixture with K components and d dimensions, the log likelihood function on the mixture model (x_j) is given by

$$L(\theta) = \sum_{t=1}^N \log (x_t | \theta) \quad (3)$$

The entropy penalized ML objective function can be constructed by

$$L(\theta) = L(\theta) - \lambda \sum_{j=1}^K \theta_j \quad (4)$$

where λ is the penalty factor which is a constant.

According to this objective function $L(\theta)$, we have an entropy penalized maximum-likelihood estimate:

$$\hat{\theta} = \arg \max_{\theta} L(\theta) \quad (5)$$

By Lagrange's method of multipliers for $\hat{\theta}$, we begin to define the \mathcal{L} -function:

$$\mathcal{L}(\theta) = L(\theta) + \lambda \left(1 - \sum_{j=1}^K \theta_j \right) \quad (6)$$

Using the general methods for matrix derivatives,⁷ we are led to the derivatives of $\mathcal{L}(\theta)$ with respect to each θ_j , λ and θ_j as follows:

According to Eqs. (12) and (15), we further have

$$\begin{aligned}
 &= \sum_{j=1}^K \sum_{t=1}^N j(\cdot) + \sum_{j=1}^K (j + j \log j) \\
 &= \dots + (1 \dots) \\
 &= (1 + (1 \dots)) \tag{16}
 \end{aligned}$$

Then, from the Eqs. (12)-(14), we obtain the following iterative algorithm for $\hat{\pi}^+$:

$$\hat{\pi}_j^+ = \frac{1}{(1 + (1 \dots))} \left(\sum_{t=1}^N j(\cdot) + j(1 + \log j) \right); \tag{17}$$

$$\hat{\pi}_j^+ = \frac{1}{\sum_{t=1}^N j(\cdot)} \sum_{t=1}^N j(\cdot) x_t; \tag{18}$$

$$\hat{\pi}_j^+ = \frac{1}{\sum_{t=1}^N j(\cdot)} \sum_{t=1}^N j(\cdot) (x_t - \mu_j)(x_t - \mu_j)^T \tag{19}$$

As compared with the EM algorithm for Gaussian mixture, this iterative algorithm implements a penalty mechanism on the mixing proportions during the iterations, which leads to the automated model selection. The penalty factor can be selected by experience.

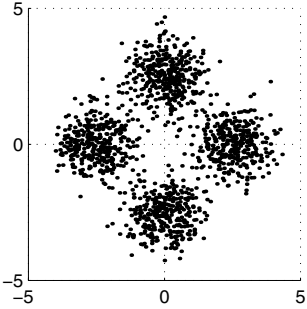
3. Simulation Results

In this section, simulation experiments are carried out to demonstrate the entropy penalized MLE iterative algorithm for both model selection and parameter estimation on a set of sample data from a Gaussian mixture. Moreover, this iterative algorithm is compared with the BYY model selection algorithms.

3.1. The sample data

We begin with a description of the seven 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 densities (i.e. $K = 2$).

As shown in Fig. 1, each data set of samples are generated with a certain degree of overlap among the Gaussians (clusters) in the mixture. They are four typical sets of sample data from Gaussian mixtures. The Gaussians in S_1 are sphere-shaped, with equal number of samples. But those in S_2 are ellipse-shaped, with different numbers of samples. Moreover, S_3 consists of three very flat Gaussians and S_4 has a small number of samples, with similar structure as S_2 . The detailed parameters for these four sets of sample data are given in Table 1 where μ_i , $\Sigma_i = [\Sigma_{jk}]_{2 \times 2}$, π_i and n_i denote the mean vector, covariance matrix, mixing proportion and the number of samples of the i th cluster (Gaussian), respectively.

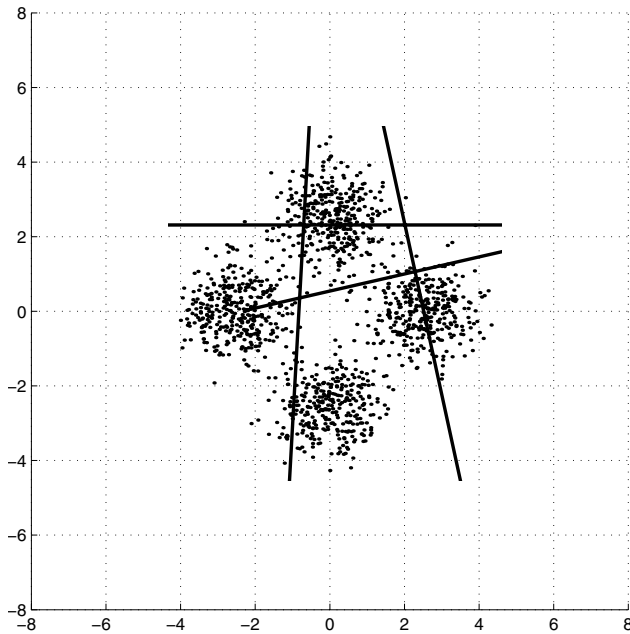


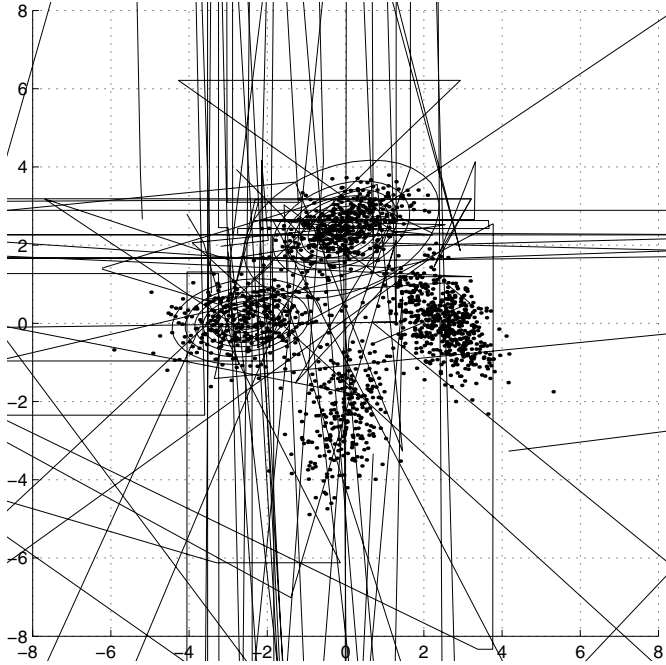
2. Automated model selection

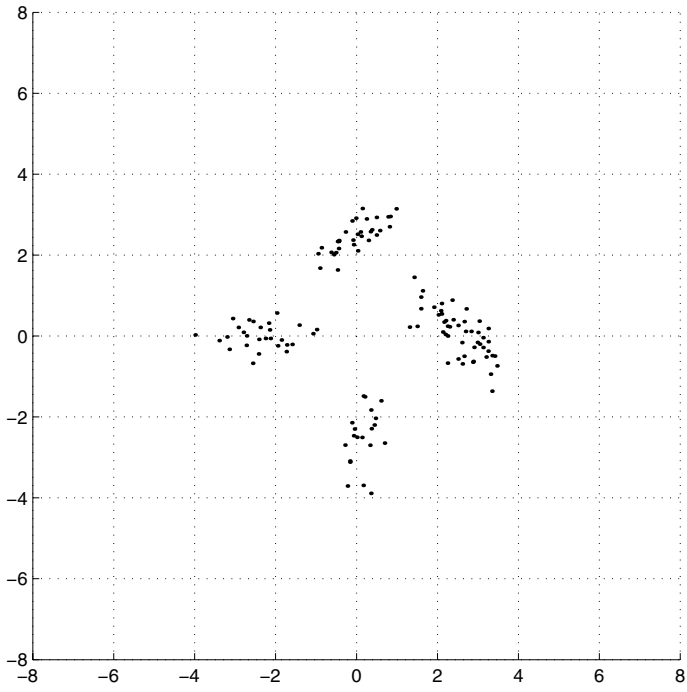
We implement the entropy penalized MLE iterative algorithm on those four sample data sets always with k^* , respectively. The parameters are initialized randomly within certain intervals, satisfying the required constraints. The penalty factor λ is selected to be 0.15. In all the experiments, the algorithm is stopped when $\|x_k^{new} - x_k^{old}\|_j < 10^{-7}$

The experimental results on S_1 and S_2 are given in Figs. 2 and 3, respectively, with case $k = 8$ and $k^* = 4$. We observe that four Gaussians are naturally located accurately, while the mixing proportions of the other four Gaussians were reduced to 0.002 or below, i.e. these Gaussians are extra and can be discarded. That is, the correct number of the Gaussians has been automatically determined on these data sets.

Moreover, the experiment has been made on S_3 with case $k = 8$ $k^* = 3$. As shown in Fig. 4, the actual Gaussians are very flat. However, three Gaussians are still located accurately, while the mixing proportions of the other five extra Gaussians are no more than 0.0012. That is, the correct number of the Gaussians can still be determined on such a special data set. Furthermore, the iterative algorithm is also implemented on S_4 with $k = 8$ $k^* = 4$. As shown in Fig. 5, each cluster has a small number of samples, the correct number of Gaussians can still be detected, with the mixing proportions of other four extra Gaussians reduced below 0.008.







3.3. *Parameter estimation*

References

- S. C. Ahalt, A. K. Krishnamurthy, P. Chen and D. E. Melton, Competitive learning algorithm for vector quantization, **Neural Networks** **3**(3) (1990) 277–291.
- H. Akaike, A new look at the statistical model identification, **IEEE Trans. Autom. Contr.** **AC-19** (1974) 716–723.
- J. C. Bezdek and N. R. Pal, Some new indexes of cluster validity, **Contr.**

21. X. L. Xie and G. A. Beni, A validity method for fuzzy clustering, **IEEE Trans. Patt. Anal. Mach. Intell.** **13**(8) (1991) 843–847.
22. L. Xu, A. Krzyzak and E. Oja, Rival penalized competitive learning for clustering analysis, RBF net, and curve detection, **IEEE Trans. Neural Networks** **4** (1993) 636–648.
23. L. Xu **et al.**, A unified learning scheme: Bayesian–Kullback YING–YANG machine, in **Advances in Neural Information Processing Systems**, eds. D. S. Touretzky **et al.** (MIT Press, Cambridge, MA, 1996), Vol. 8, pp. 444–450; A part of its preliminary version on **Proc. 1995 Int. Conf. Neural Information Processing**, pp. 977–988.
24. L. Xu, Rival penalized competitive learning, finite mixture, and multisets clustering, in **Proc. 1998 IEEE Int. Joint Conf. Neural Networks**, Vol. 3, pp. 251–253.
25. L. Xu, Best harmony, unified RPCL and automated model selection for unsupervised and supervised learning on Gaussian mixtures, three-layer nets and ME-RBF-SVM models, **Int. J. Neural Syst.** **11**(1) (2002) 43–69.
26. L. Xu, BYY harmony learning, structural RPCL, and topological self-organizing on mixture modes, **Neural Networks** **15** (2002) 1231–1237.
27. L. Xu, Ying–Yang learning, in **The Handbook of Brain Theory and Neural Networks**, ed. M. A. Arbib, 2nd edn. (The MIT Press, 2002), pp. 1231–1237.



Jinwen Ma received the M.S. degree in applied mathematics from Xi'an Jiaotong University in 1988 and the Ph.D. in probability theory and statistics from Nankai University in 1992. From July 1992 to November 1999, he

was a Lecturer and Associate Professor at the Department of Mathematics, Shantou University. He has been a full Professor at the Institute of Mathematics, Shantou University since December, 1999. In September 2001, he was transferred to the Department of Information Science at the School of Mathematical Sciences, Peking University. During 1995 and 2003, he also made several visits to the Department of Computer Science and Engineering, the Chinese University of Hong Kong as a Research Associate or Fellow.

He has published over 40 academic papers on neural networks, pattern recognition, artificial intelligence and information theory.



Taijun Wang received the M.S. degree in signal and information processing from Southeast University, China, in 1982. He has been an Associate Professor at the Department of Radio Engineering of Southeast University

since 1988.

His research interests are related to digital signal processing, pattern recognition, artificial neural networks, statistical learning theory and scientific visualization.