

Bayesian Mixture Labeling by Minimizing Deviance of Classification Probabilities to Reference Labels

Weixin Yao and Longhai Li

Abstract

Solving label switching is crucial for interpreting the results of fitting Bayesian mixture models. The label switching originates from the invariance of posterior distribution to permutation of component labels. As a result, the component labels in Markov chain simulation may switch to another equivalent permutation, and the marginal posterior distribution associated with all labels may be similar and useless for inferring quantities relating to each individual component. In this article, we propose a new simple labeling method by minimizing the deviance of the class probabilities to a fixed reference labels. The reference labels can be chosen before running MCMC using optimization methods, such as EM algorithms, and therefore the new labeling method can be implemented by an online algorithm, which can reduce the storage requirements. Using the Acid data set and Galaxy data set, we demonstrate the success of the proposed labeling method for removing the labeling switching in the raw MCMC samples.

Key words: Bayesian mixtures; Label switching; Markov chain Monte Carlo; Mixture models; Relabeling.

¹Weixin Yao is Assistant Professor, Department of Statistics, Kansas State University, Manhattan, Kansas 66506, U.S.A. Email: wxyao@ksu.edu. Longhai Li is Assistant Professor, Department of Mathematics and Statistics, University of Saskatchewan, Saskatoon, Saskatchewan, S7N5E6, Canada. Email: longhai@math.usask.ca. The research of Longhai Li is supported by fundings from Natural Sciences and Engineering Research Council of Canada, and Canadian Foundation of Innovations.

1 Introduction

Label switching is one of the fundamental issues for Bayesian mixtures if our interests are quantities relating to each individual component. It occurs due to the invariance of the posterior distribution to the permutation of the component labels. Many methods have been proposed to solve the label switching problem. One simple way is to use an explicit parameter identifiability constraint so that only one permutation can satisfy it. See Diebolt and Robert (1994); Dellaportas et al. (1996); Richardson and Green (1997). One problem with the identifiability constraint labeling is that the results are sensitive to the choice of constraint, especially for multivariate problems. Celeux et al. (2000) demonstrated that different order constraints may generate markedly different results; it is difficult to anticipate the overall effect. Moreover, many choices of identifiability constraint do not completely remove the symmetry of the posterior distribution. As a result, label switching problem may remain after imposing an identifiability constraint. See the example by Stephens (2000). Stephens (2000) and Celeux (1998) proposed a relabeling algorithm, which is based on minimizing a Monte Carlo risk. Yao and Lindsay (2009) proposed to label the samples based on the posterior modes and an ascent algorithm (PM(ALG)). PM(ALG) uses each Markov chain Monte Carlo (MCMC) sample as the starting point in an ascending algorithm, and labels the sample based on the mode of the posterior to which it converges. Then PM(ALG) assumes that the samples converged to the same mode have the same labels. Sperrin, Jaki, and Wit (2010) developed several probabilistic relabeling algorithms by extending the probabilistic relabeling of Jasra (2005).

Papastamoulis and Iliopoulos (2010) proposed an artificial allocations based solution to the label switching problem. Yao (2012) proposed to assign the probabilities for each possible labels by fitting a mixture model to the permutation symmetric posterior. Other labeling methods include, for example, Celeux et al. (2000); Fruhwirth (2001); Hurn et al. (2003); Chung et al. (2004); Marin et al. (2005); Geweke (2007); Grun and Leisch (2009); Cron and West (2011). Jasra et al. (2005) provided a good review about the existing methods to solve

the label switching problem in Bayesian mixture modeling.

In this article, we propose a new alternative labeling method by minimizing the deviance of the class probabilities to a fixed reference labels. The reference labels may be chosen before running MCMC using optimization methods, such as EM algorithms, and therefore the new labeling method can be implemented by an online algorithm, i.e., the output of MCMC samples will have been automatically relabeled along with simulating MCMC samples, which can reduce the storage requirements. The reference labels can also be chosen after MCMC sampling by alternating two steps of finding the reference labels and relabeling MCMC samples.

The rest of the paper is organized as follows. Section 2 introduces our new labeling method. In Section 3, we use the Acid data set and Galaxy data set to demonstrate the success of the proposed labeling method. We summarize our proposed labeling method in Section 4.

2 New Method

Generally, the mixture model has the density

$$p(x; \boldsymbol{\theta}) = \pi_1 f(x; \lambda_1) + \pi_2 f(x; \lambda_2) + \cdots + \pi_m f(x; \lambda_m), \quad (2.1)$$

where $\boldsymbol{\theta} = (\pi_1, \dots, \pi_m, \lambda_1, \dots, \lambda_m)$, $f(\cdot)$ is the component density, λ_j is the component specific parameter, which can be scalar or vector and π_j is the proportion of the j th component in the whole population with $\sum_{i=1}^m \pi_i = 1$. If $\mathbf{x} = (x_1, \dots, x_n)$ are independent observations from the m -component mixture model (2.1), the likelihood of $\boldsymbol{\theta}$ given \mathbf{x} is

$$L(\boldsymbol{\theta}; \mathbf{x}) = \prod_{i=1}^n \{\pi_1 f(x_i; \lambda_1) + \pi_2 f(x_i; \lambda_2) + \cdots + \pi_m f(x_i; \lambda_m)\}. \quad (2.2)$$

A permutation $\omega = (\omega(1), \dots, \omega(m))$ of the component labels $\{1, \dots, m\}$ defines a corresponding permutation of the parameter vector θ by

$$\theta^\omega = (\pi^\omega, \lambda^\omega) = (\pi_{\omega(1)}, \dots, \pi_{\omega(m)}, \lambda_{\omega(1)}, \dots, \lambda_{\omega(m)}).$$

A special feature of mixture model is that the likelihood function $L(\theta^\omega; \mathbf{x})$ is exactly the same as $L(\theta; \mathbf{x})$ for any permutation ω .

For Bayesian mixtures, if the prior distributions for model parameters are symmetric for all components then the posterior distribution for the parameters will be also symmetric and thus invariant to permutations in the labeling of the component parameters. The marginal posterior distributions for the parameters will be identical for all mixture components. Then the posterior means of each component are the same and are thus poor estimates of these parameters. Similar problem will occur when we try to estimate quantities relating to individual components of the mixture such as predictive component densities, marginal classification probabilities. So in Bayesian analysis, after we get a sequence of simulated values $\theta_1, \dots, \theta_N$ from the posterior distribution of θ given $Y = y$ using MCMC sampling methods, we must first find permutations $\{\omega_1, \dots, \omega_N\}$ such that $\theta_1^{\omega_1}, \dots, \theta_N^{\omega_N}$ have the same label meaning, then we can use the labeled samples to do Bayesian analysis. Many methods (as reviewed in Section 1) have been proposed to find $\{\omega_1, \dots, \omega_N\}$ for relabeling MCMC samples. In this article, we introduce a new method.

Given observations $\mathbf{x} = (x_1, \dots, x_n)$, suppose we have found a set of reference component labels for each observation x_i represented by $\mathbf{Z} = \{Z_{ij}, i = 1, \dots, n, j = 1, \dots, m\}$, where

$$Z_{ij} = \begin{cases} 1, & \text{if the } i^{th} \text{ observation } x_i \text{ is from the } j^{th} \text{ component;} \\ 0, & \text{otherwise.} \end{cases}$$

We will talk about how to find the reference label \mathbf{Z} later. Let $\theta = (\lambda_1, \dots, \lambda_m, \pi_1, \dots, \pi_m)$. Our new method for finding a permutation ω for relabeling a Markov chain sample θ (note

that we drop MCMC index since the relabeling method will be applied to all iterations parallelly) is to minimize the sum of minus log classification probabilities of \mathbf{Z} given by θ^ω with respect to ω :

$$\ell(\omega; \mathbf{Z}, \theta) = - \sum_{i=1}^n \sum_{j=1}^m Z_{ij} \log p_{ij}(\theta^\omega), \quad (2.3)$$

where $p_{ij}(\theta^\omega)$ is the classification probability that the i th observation belongs to j th component based on relabeled parameter θ^ω :

$$p_{ij}(\theta^\omega) = \frac{\pi_\omega(j) f(x_i; \lambda_\omega(j))}{p(x_i; \theta^\omega)} = p_{i, \omega(j)}(\theta). \quad (2.4)$$

The objective function $\ell(\omega; \mathbf{Z}, \mathbf{x})$ in (2.3) can be also considered as the Kullback-Leibler divergence if we consider Z_{ij} as the true classification probability and $p_{ij}(\theta)$ as the estimated classification.

Note that $2\ell(\omega; \mathbf{Z}, \mathbf{x})$ is often called deviance of classification probabilities $p_{ij}(\theta^\omega)$, $i = 1, \dots, n, j = 1, \dots, m$ to the reference labels \mathbf{Z} in the literature of generalized linear models, if \mathbf{Z} is the true response values and $p_{ij}(\theta^\omega)$ is the predictive probabilities based on a generalized linear model. In words, by minimizing $\ell(\omega; \mathbf{Z}, \mathbf{x})$ with respect to ω we will find the optimal permutation ω for a Markov chain sample θ such that the corresponding classification probabilities can best explain the reference label \mathbf{Z} . It is crucial to note that our method uses the differences of the whole probability density functions $f(x; \lambda_j)$ and mixture proportion π_j of all mixture components $j = 1, \dots, m$ in relabeling θ rather than the values of a single or an arbitrarily chosen subset of parameters in θ . Our method therefore works well in the situations where any single parameter in θ cannot clearly distinguish all components but the density functions given the whole set of parameters are clearly different for components.

The proposed objective function (2.3) has another nice interpretation based on complete posterior distribution. Let $\pi(\theta)$ be the prior for θ . Then the posterior for complete data

(\mathbf{x}, \mathbf{Z}) is

$$p_c(\boldsymbol{\theta}; \mathbf{x}, \mathbf{Z}) = \pi(\boldsymbol{\theta})p(\mathbf{x}, \mathbf{Z} \mid \boldsymbol{\theta}) = \pi(\boldsymbol{\theta}) \prod_{i=1}^n \prod_{j=1}^m \{\pi_j f(x_i; \lambda_j)\}^{Z_{ij}}.$$

Note that the above complete posterior is *not* invariant to the component labels and thus can be used to do labeling. Given the reference label \mathbf{Z} , it is natural to do labeling for $\boldsymbol{\theta}$ by maximizing the log complete posterior

$$\log p_c(\boldsymbol{\theta}^\omega; \mathbf{x}, \mathbf{Z}) = \log \pi(\boldsymbol{\theta}^\omega) + \sum_{i=1}^n \sum_{j=1}^m [Z_{ij} \log(\pi_j^\omega f(x_i; \lambda_j^\omega))] \quad (2.5)$$

with respect to (\mathbf{Z}, ω) , where $\pi_j^\omega = \pi_{\omega(j)}$, and $\lambda_j^\omega = \lambda_{\omega(j)}$.

Note that

$$\begin{aligned} & \log p_c(\boldsymbol{\theta}^\omega; \mathbf{x}, \mathbf{Z}) \\ &= \log \pi(\boldsymbol{\theta}^\omega) + \sum_{i=1}^n \sum_{j=1}^m [Z_{ij} \log\{\pi_j^\omega f(x_i; \lambda_j^\omega)/p(x_i; \boldsymbol{\theta}^\omega)\}] + \sum_{i=1}^n \sum_{j=1}^m Z_{ij} \log p(x_i; \boldsymbol{\theta}^\omega) \\ &= \log \pi(\boldsymbol{\theta}^\omega) + \sum_{i=1}^n \sum_{j=1}^m Z_{ij} \log p_{ij}(\boldsymbol{\theta}^\omega) + \sum_{i=1}^n \log p(x_i; \boldsymbol{\theta}^\omega), \end{aligned} \quad (2.6)$$

where $p(x_i; \boldsymbol{\theta}^\omega) = \sum_{j=1}^m \pi_j^\omega f(x_i; \lambda_j^\omega)$. Notice that the first and third terms of (2.6) are invariant to the permutation of ω . Therefore, maximizing (2.6) is equivalent to maximizing the second term of (2.6), which is equivalent to minimizing (2.3).

There are many methods for finding reference labels $\mathbf{Z} = (Z_{ij}, i = 1, \dots, n, j = 1, \dots, m)$. One simple method is to find the posterior mode, say $\hat{\boldsymbol{\theta}}$, and the corresponding classification probabilities, say $p_{ij}(\hat{\boldsymbol{\theta}})$. Then Z_{ij} is estimated by maximizing the classification probabilities over all components, i.e.,

$$Z_{ij} = \begin{cases} 1, & \text{if } p_{ij}(\hat{\boldsymbol{\theta}}) \geq p_{il}(\hat{\boldsymbol{\theta}}) \text{ for all } l = 1, \dots, m; \\ 0, & \text{otherwise.} \end{cases} \quad (2.7)$$

To find the posterior mode, one might simply calculate the posterior for each MCMC

sample and then use the sample that has the largest posterior to approximate the posterior mode. In addition, Yao and Lindsay (2009) also proposed the ECM algorithm for Bayesian mixtures to find the posterior mode. Suppose that we can use Gibbs sampler to get the MCMC samples and there exists a partition of $\boldsymbol{\theta} = \{\boldsymbol{\theta}_{(1)}, \dots, \boldsymbol{\theta}_{(p)}\}$ such that all the conditional complete posterior distributions $\{p(\boldsymbol{\theta}_{(i)} \mid \dots), 1 \leq i \leq p\}$ can be easily found, where $\boldsymbol{\theta}_{(i)}$ can be scalar or vector and $\mid \dots$ denotes conditioning on all other parameters and the latent variable \mathbf{Z} . In the E step, the ECM algorithm calculates the classification probabilities, p_{ij} , for each observation. In the M step, the ECM algorithm maximizes the conditional complete posterior distribution $p(\boldsymbol{\theta}_{(i)} \mid \dots), 1 \leq i \leq p$ sequentially with the latent variable Z_{ij} replaced by the classification probability p_{ij} . The ECM iterates the above E step and M step until convergence. In our examples we will use this ECM algorithm to find the posterior mode and the corresponding reference labels \mathbf{Z} .

Therefore, the above proposed labeling procedure can be summarized as follows.

Algorithm 2.1. *Step 1: Find the posterior mode and the corresponding reference labels $\mathbf{Z} = (Z_{ij}, i = 1, \dots, n, j = 1, \dots, m)$.*

Step 2: For each MCMC sample $\boldsymbol{\theta}_t$, choose $\boldsymbol{\omega}_t$ to minimize $\ell(\boldsymbol{\omega}_t; \mathbf{Z}, \boldsymbol{\theta}_t)$ of (2.3).

One main advantage of the above algorithm is that it can be implemented before MCMC simulation. The reference label \mathbf{Z} will then be used along with simulating MCMC, which saves storage, and therefore the above algorithm is an online algorithm — the output of MCMC samples will have been automatically relabeled.

Following Stephens (2000), we can also find the reference label \mathbf{Z} after simulating MCMC, by simultaneously finding \mathbf{Z} and $\boldsymbol{\Omega} = (\boldsymbol{\omega}_1, \dots, \boldsymbol{\omega}_N)$ that minimizes a Monte Carlo risk:

$$R(\mathbf{Z}, \boldsymbol{\omega}_1, \dots, \boldsymbol{\omega}_N) = \sum_{t=1}^N \ell(\boldsymbol{\omega}_t; \mathbf{Z}, \boldsymbol{\theta}), \quad (2.8)$$

where ℓ is given by (2.3). We propose the following algorithm for minimizing (2.8):

Algorithm 2.2. *Starting with some initial values for $\omega_1, \dots, \omega_N$ (set by order constraint labels for example), iterate the following two steps until a fixed point is reached.*

Step 1: Given \mathbf{Z} , for each t , choose ω_t to minimize $\ell(\omega_t; \mathbf{Z}, \boldsymbol{\theta}_t)$. In other words, relabel all Markov chain iterations with \mathbf{Z} .

Step 2: Estimate \mathbf{Z} by

$$Z_{ij} = \begin{cases} 1, & \text{if } \sum_{t=1}^N \log p_{ij}(\boldsymbol{\theta}_t^{\omega_t}) > \sum_{t=1}^N \log p_{il}(\boldsymbol{\theta}_t^{\omega_t}) \text{ for all } l \neq j; \\ 0, & \text{o.w.} \end{cases},$$

where $i = 1, \dots, n, j = 1, \dots, m$.

Note however, similar to Stephens (2000), the Algorithm 2.2 can only be implemented after saving all MCMC samples and thus is not an online algorithm.

Proposition 2.1. *The Algorithm 2.2 must converge and monotonically decrease the objective function (2.8).*

3 Examples

3.1 Acidity Data

We first consider the acidity data set (Crawford et al., 1992; Crawford, 1994). The observations are the logarithms of an acidity index measured in a sample of 155 lakes in north-central Wisconsin. The data are shown in Figure 1. Crawford et al. (1992), Crawford (1994), and Richardson and Green (1997) have used a mixture of Gaussian distributions to analyze this data set. Here, we fit this data set by a three-component normal mixture based on the result of Richardson and Green (1997). The MCMC samples are generated by Gibbs sampler with the priors given by Phillips and Smith (1996) and Richardson and Green (1997). That is to assume

$$\boldsymbol{\pi} \sim D(\delta, \delta, \delta), \quad \mu_j \sim N(\xi, \kappa^{-1}), \quad \sigma_j^{-2} \sim \Gamma(\alpha, \beta), \quad j = 1, 2, 3,$$

where $D(\cdot)$ is Dirichlet distribution and $\Gamma(\alpha, \beta)$ is gamma distribution with mean α/β and variance α/β^2 , $\delta = 1$, ξ equal the sample mean of the observations, κ equal $1/R^2$, $\alpha = 2$, and $\beta = R^2/200$, where R is the range of the observations. Similar priors are used for the other example. We post processed the 20,000 Gibbs samples by the proposed labeling method.

Figure 2 and 3 are the trace plots and the estimated marginal posterior density plots, respectively, for the original samples and the labeled samples by new method. From Figures 2(a) and 3(a), we can see that the label switching occurred in the raw samples and the marginal density plots display the multi-modality. Based on Figures 2(b) and 3(b), we can see that the new labeling method successfully removed the label switching in the raw output of the Gibbs sampler.

3.2 Galaxy Data

The galaxy data (Roeder, 1990) consists of the velocities (in thousands of kilometers per second) of 82 distant galaxies diverging from our own galaxy. They are sampled from six well-separated conic sections of the corona borealis. A histogram of the 82 data points is shown in Figure 4. This data set has been analyzed by many researchers including, for example, Crawford (1994); Chib (1995); Carlin and Chib (1995); Escobar (1995); Phillips and Smith (1996); Richardson and Green (1997). Stephens (2000) also used this data set to explain the label switching problem. We fit this data by six-component normal mixture. We post processed the 20,000 Gibbs samples by the proposed labeling method.

Figure 5 and 6 are the trace plots and the estimated marginal posterior density plots, respectively, for the original samples and the labeled samples by new method. From Figure 5 and 6, one can see that the new labeling method successfully removed the label switching in the raw output of the Gibbs sampler.

For the above two examples, the methods proposed by Stephens (2000) and Yao and Lindsay (2009) showed similar performance. However, as Stephens (2000) showed that the order constraint labeling *failed* to remove the label switching for Galaxy data.

4 Summary

Label switching has been a long standing problem for Bayesian mixtures. In this paper, we proposed a new alternative labeling method by minimizing deviance of classification probabilities to reference labels. The new labeling method also has a nice interpretation based on the complete posterior likelihood. After finding the reference labels, the new method can be implemented without saving all MCMC samples and classification probabilities, i.e, the output of MCMC samples will have been automatically relabeled along with simulating MCMC samples. Therefore, the new method is an on online algorithm, which can reduce much storage requirements. The two real data applications in Section 3 demonstrates the success of the new method in removing the label switching in the raw MCMC samples.

References

- Bezdek, J. C., Hathaway, R. M., and Huggins, V. J. (1985). Parametric estimation for normal mixtures. *Pattern Recognition*, 3, 79-84.
- Böhning, D. (1999). *Computer-Assisted Analysis of Mixtures and Applications*, Boca Raton, FL: Chapman and Hall/CRC.
- Carlin, B. P. and Chib, S. (1995). *Bayesian model choice via Markov chain Monte Carlo methods*. *Journal of Royal Statistical Society*, B57, 473-484.
- Celeux, G. (1998). Bayesian inference for mixtures: The label switching problem. In *Computat 98-Proc. in Computational Statistics* (eds. R. Payne and P.J. Green), 227-232. Physica, Heidelberg.
- Celeux, G., Hurn, M., and Robert, C. P. (2000). Computational and inferential difficulties with mixture posterior distributions. *J. Am. Statist. Ass.*, 95, 957-970.
- Chib, S. (1995). Marginal likelihood from the Gibbs output. *Journal of American Statistical Association*, 90, 1313-1321.
- Chung, H., Loken, E., and Schafer, J. L. (2004). Difficulties in drawing inferences with finite-mixture models: a simple example with a simple solution. *The American Statistician*, 58, 152-158.
- Crawford, S. L., Degroot, M. H., Kadane, J. B., and Small, M. J. (1992). Modeling lake-chemistry distributions-approximate Bayesian methods for estimating a finite-mixture model. *Technometrics*, 34, 441-453.
- Crawford, S. L. (1994). An application of the Laplace method to finite mixture distributions. *Journal of the American Statistical Association*, 89, 259-267.
- Cron A. J. and West, M. (2011). Efficient classification-based relabeling in mixture models. *The American Statistician*, 65, 16-20.

- Dellaportas, P., Stephens, D. A., Smith, A. F. M., and Guttman, I. (1996). A comparative study of perinatal mortality using a two-component mixture model. In *Bayesian Biostatistics* (eds. D.A. Berry and D.K. Stangl) 601-616, Dekker, New York.
- Dempster, A. P., Laird, N. M., and Rubin, D. B. (1977). Maximum likelihood from incomplete data via the EM algorithm (with discussion). *J. R. Statist. Soc. B*, 39, 1-38.
- Diebolt, J. and Robert, C. P. (1994). Estimation of finite mixture distributions through Bayesian sampling. *J. R. Statist. Soc. B*, 56, 363-375.
- Escobar, M. D. and West, M. (1995), “Bayesian density estimation and inference using mixtures,” *Journal of the American Statistical Association*, 90, 577-588.
- Frühwirth-Schnatter, S. (2001). Markov chain Monte Carlo estimation of classical and dynamic switching and mixture models. *J. Am. Statist. Ass.*, 96, 194-209.
- Frühwirth-Schnatter, S. (2006). *Finite Mixture and Markov Switching Models*, Springer, 2006.
- Geweke, J. (2007). Interpretation and inference in mixture models: Simple MCMC works. *Computational Statistics and Data Analysis*, 51, 3529-3550.
- Grün, B. and Leisch, F. (2009). Dealing with label switching in mixture models under genuine multimodality. *Journal of Multivariate Analysis*, 100, 851-861.
- Hathaway, R. J. (1983). Constrained maximum likelihood estimation for a mixture of univariate normal distributions. *Technical report No. 92* Columbia, South Carolina: University of Carolina.
- Hathaway, R. J. (1985). A constrained formulation of maximum-likelihood estimation for normal mixture distributions. *Ann. Statist.*, 13, 795-800.
- Hathaway, R. J. (1986). A constrained EM algorithm for univariate mixtures. *Journal of Statistical Computation and Simulation*, 23, 211-230.

- Hurn, M., Justel, A., and Robert, C. P. (2003). Estimating mixtures of regressions. *Journal of Computational and Graphical Statistics*, 12, 55-79.
- Jasra, A, Holmes, C. C., and Stephens D. A. (2005). Markov chain Monte Carlo methods and the label switching problem in Bayesian mixture modeling. *Statistical Science*, 20, 50-67.
- Lindsay, B. G., (1995). *Mixture Models: Theory, Geometry, and Applications*. NSF-CBMS Regional Conference Series in Probability and Statistics v 5, Hayward, CA: Institute of Mathematical Statistics.
- Marin, J.-M., Mengersen, K. L. and Robert, C. P. (2005). Bayesian modelling and inference on mixtures of distributions. *Handbook of Statistics* 25 (eds. D. Dey and C.R. Rao), North-Holland, Amsterdam.
- McLachlan, G. J. and Peel, D. (2000). *Finite Mixture Models*. New York: Wiley.
- Phillips, D. B. and Smith, A. F. M. (1996). Bayesian model comparison via jump diffusion. *Markov Chain Monte Carlo in Practice*, ch. 13, 215-239, London: Chapman and Hall.
- Papastamoulis, P. and Iliopoulos, G. (2010). An artificial allocations based solution to the label switching problem in Bayesian analysis of mixtures of distributions. *Journal of Computational and Graphical Statistics*, 19, 313-331.
- Richardson, S. and Green, P. J. (1997). On Bayesian analysis of mixtures with an unknown number of components (with discussion). *Journal of Royal Statistical Society*, B59, 731-792.
- Roeder, K. (1990). Density estimation with confidence sets exemplified by superclusters and voids in the galaxies. *Journal of American Statistical Association*, 85, 617-624.
- Sperrin, M., Jaki, T., and Wit, E. (2010). Probabilistic relabeling strategies for the label switching problem in Bayesian mixture models. *Statistics and Computing*, 20, 357-366.

- Stephens, M. (2000). Dealing with label switching in mixture models. *Journal of Royal Statistical Society*, B62, 795-809.
- Yao, W. (2012). Model based labeling for mixture models. *Statistics and Computing*, 22, 337-347.
- Yao, W. and Lindsay, B. G. (2009). Bayesian mixture labeling by highest posterior density. *Journal of American Statistical Association*, 104, 758-767.

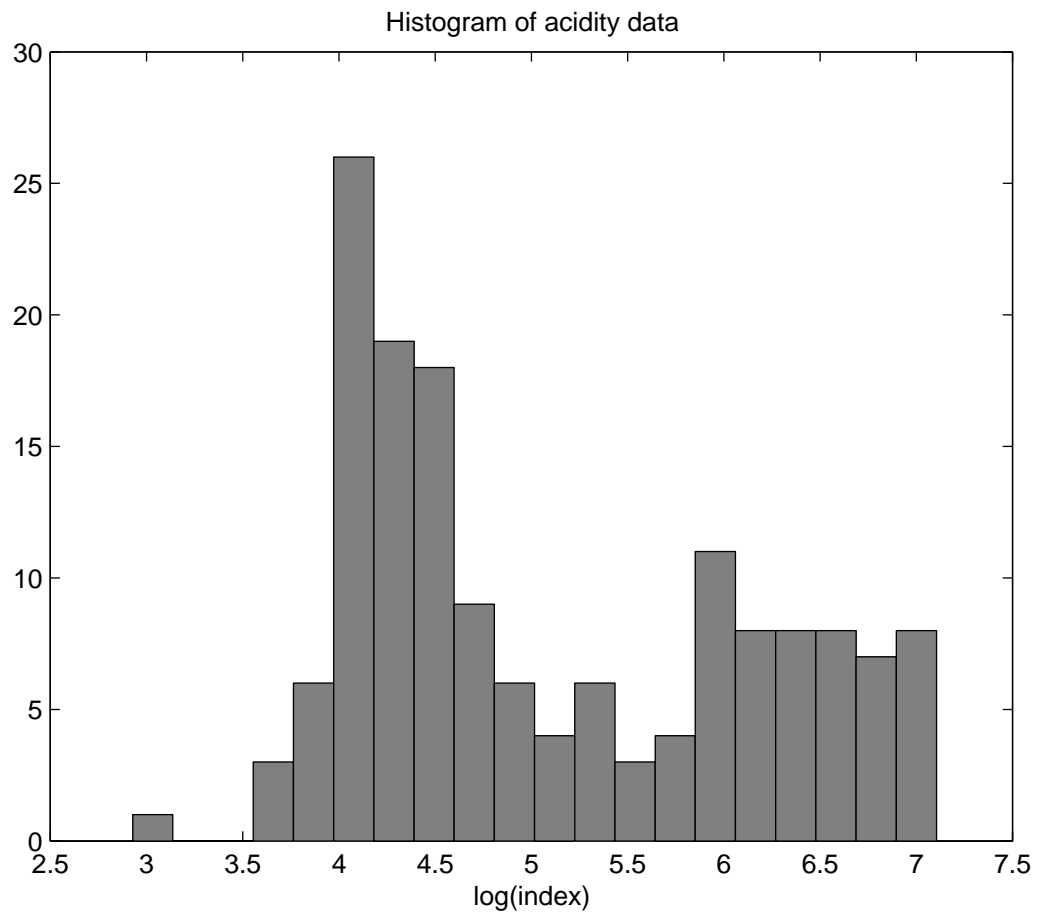
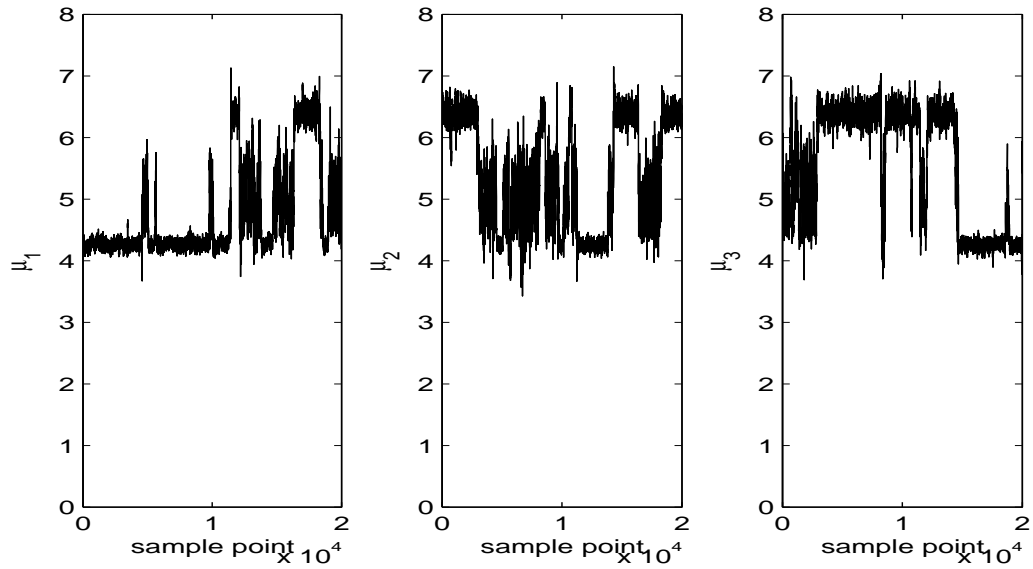
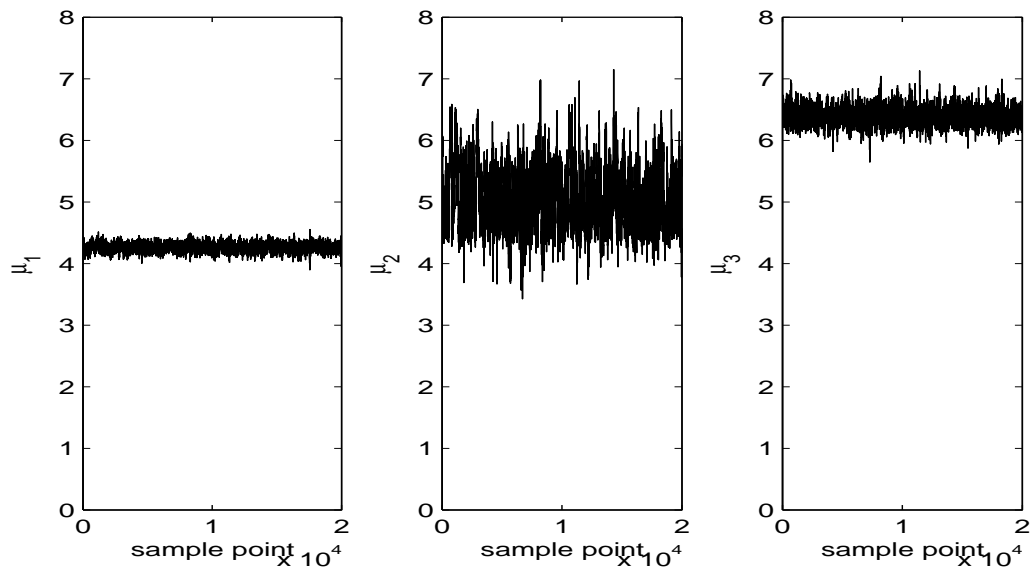


Figure 1: Histogram of acidity data. The number of bins used is 20.

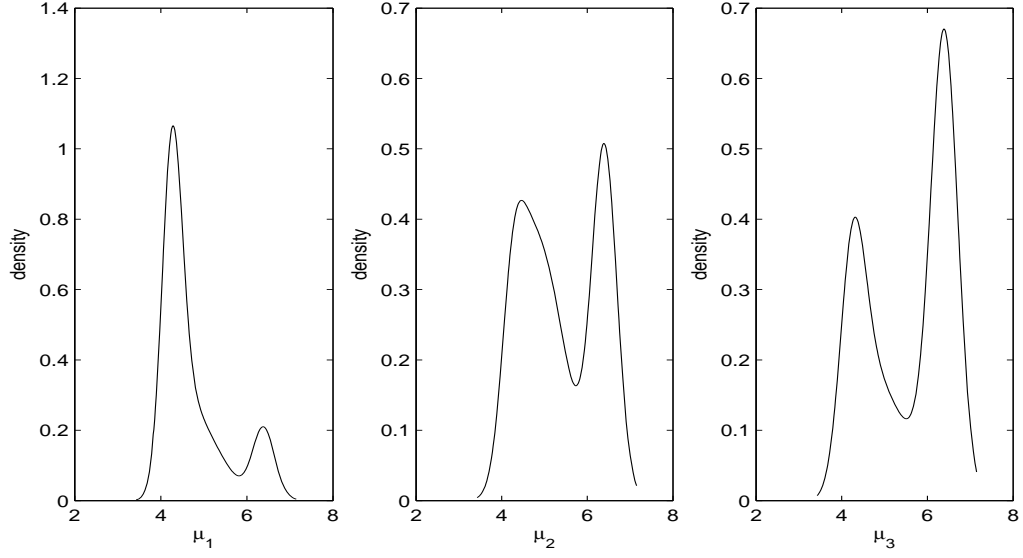


(a)

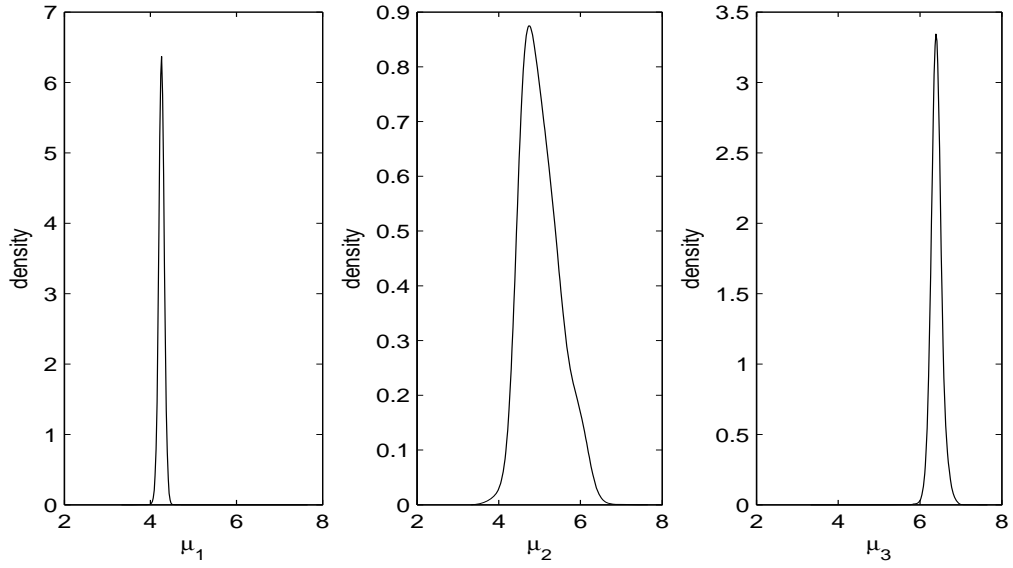


(b)

Figure 2: Trace plots of the Gibbs samples of component means for acidity data: (a) original Gibbs samples; (b) labeled samples by the new method.



(a)



(b)

Figure 3: Plots of estimated marginal posterior densities of component means for acidity data based on: (a) original Gibbs samples; (b) labeled samples by the new method.

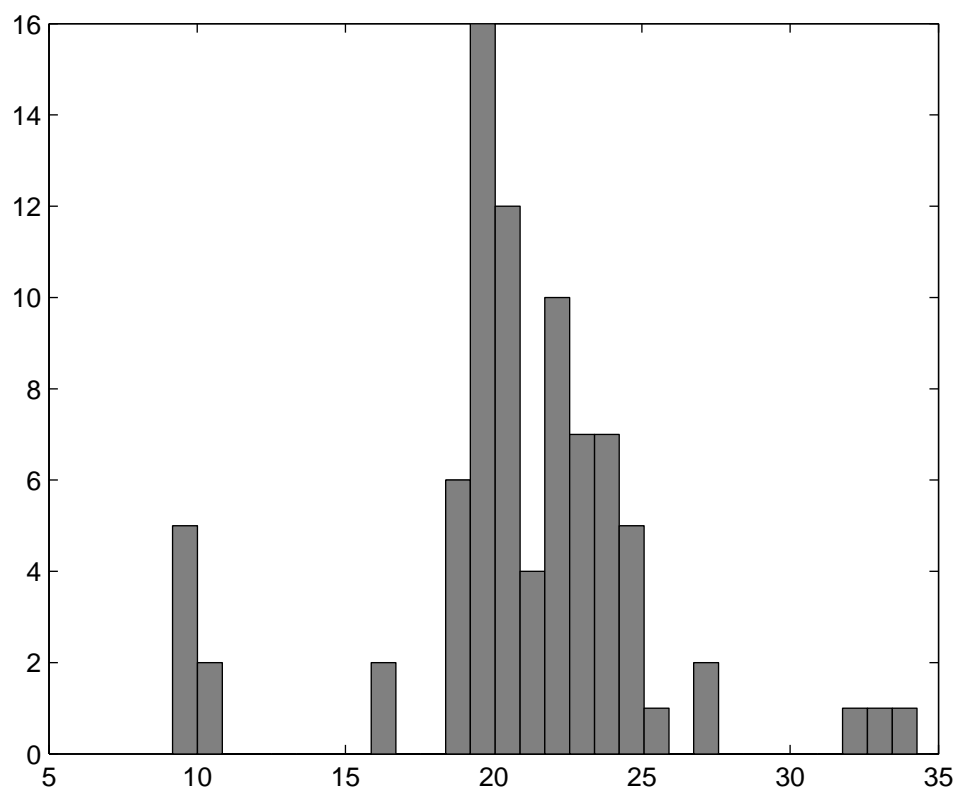
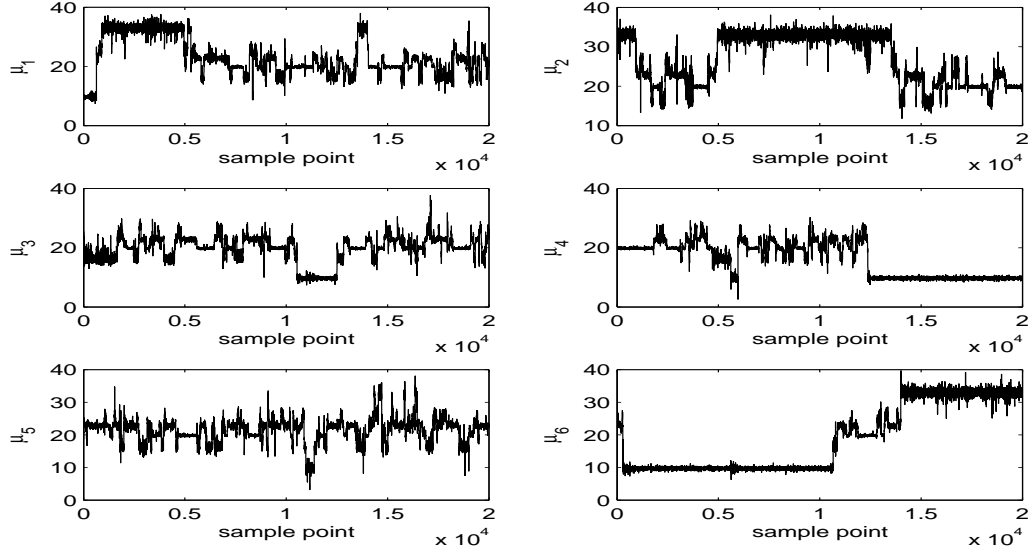
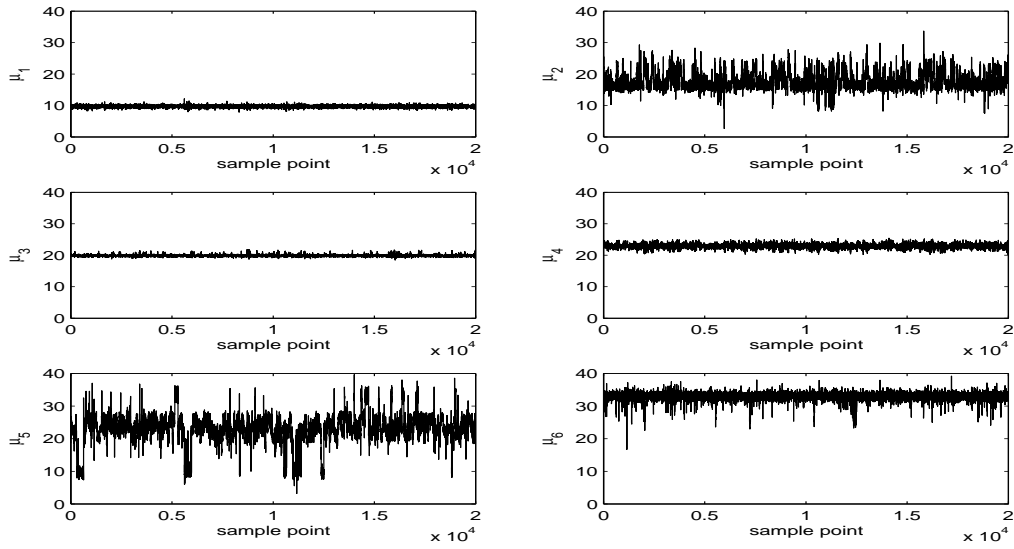


Figure 4: Histogram plot of galaxy data. The number of bins used is 30.

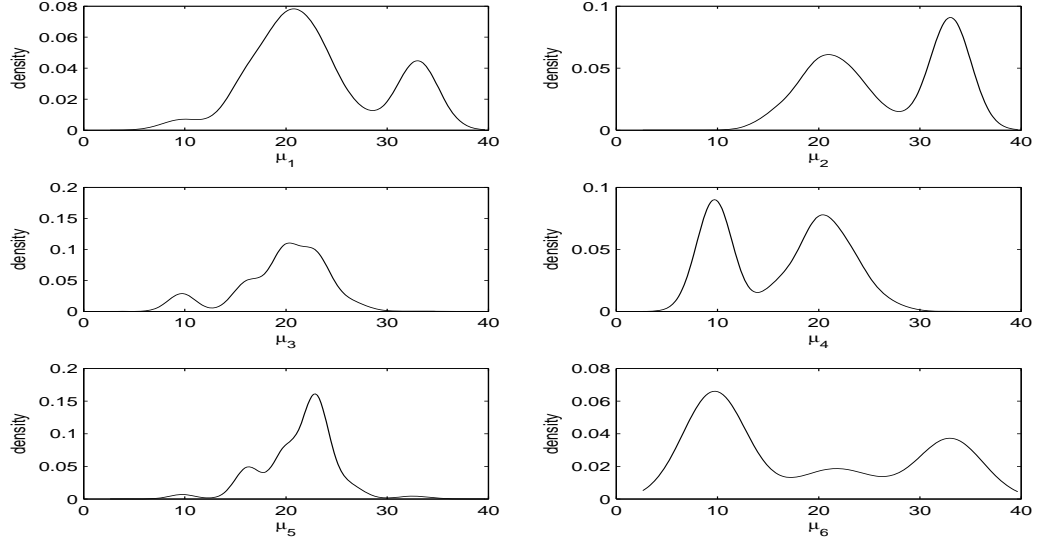


(a)

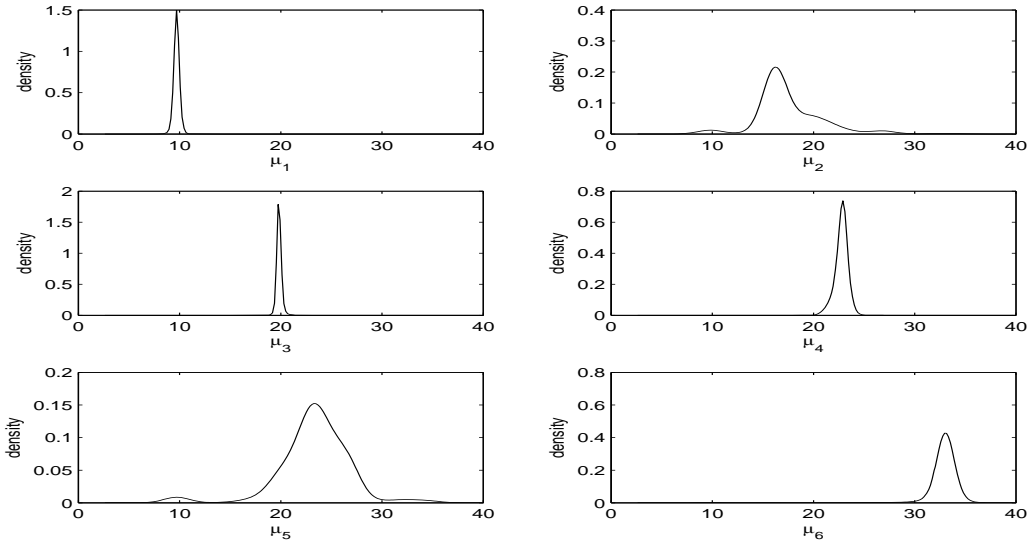


(b)

Figure 5: Trace plots of the Gibbs samples of component means for galaxy data: (a) original Gibbs samples; (b) labeled samples by the new method.



(a)



(b)

Figure 6: Plots of estimated marginal posterior densities of component means for galaxy data based on: (a) original Gibbs samples; (b) labeled samples by the new method.