

A Split-Merge Markov Chain Sampling Algorithm For Bayesian Mixture Models

Sonia Jain *

Department of Statistics

University of Toronto

Email: sonia@utstat.toronto.edu

Internet: <http://www.utstat.toronto.edu/~sonia/>

Radford M. Neal †

Dept. of Statistics and Dept. of Computer Science

University of Toronto

Email: radford@utstat.toronto.edu

Internet: <http://www.cs.toronto.edu/~radford/>

Abstract. We propose a split-merge Markov chain algorithm to address the problem of inefficient sampling for conjugate Dirichlet process mixture models. Traditional Markov chain Monte Carlo methods for Bayesian mixture models, such as Gibbs sampling, can become trapped in isolated modes corresponding to an inappropriate clustering of data points. This article describes a Metropolis-Hastings procedure that can escape such local modes by splitting or merging mixture components. Our algorithm employs a new technique in which an appropriate proposal for splitting or merging components is obtained by using a restricted Gibbs sampling scan. We demonstrate empirically that our method outperforms the Gibbs sampler in situations where two or more components are similar in structure.

1 Introduction

Mixture models are often applied to density estimation, latent class analysis, and classification problems. The Bayesian approach to mixture models has recently generated interest due to advances in statistical computation, in particular Markov chain Monte Carlo. In this article, we consider Bayesian mixture models in which a Dirichlet process prior on the mixing distribution is used to handle a countably infinite number of mixture components. Computational techniques for Dirichlet process mixture models have been explored previously by Escobar and West (1995), MacEachern (1994), Neal (1992, 2000), and Green and Richardson (2001).

When conjugate priors are used, a Gibbs sampling procedure can easily be constructed for the Dirichlet process mixture model. Here, we fully exploit the conjugacy in the model by analytically integrating away the mixing proportions for the components and the parameters for each component. As a result, the Gibbs sampling procedure only updates the latent indicator variables associating mixture

*The first author acknowledges support from a Natural Sciences and Engineering Research Council of Canada Postgraduate Scholarship.

†The second author's research was supported by the Natural Sciences and Engineering Research Council of Canada and by the Institute for Robotics and Intelligent Systems.

components with data observations. This particular Gibbs sampling method was first discussed by Neal (1992) for models of categorical data and MacEachern (1994) for normal mixture models.

Although the Gibbs sampling approach is straightforward and easily implemented, it can be slow to converge and mix poorly. When two or more mixture components have similar parameters, the Gibbs sampling method may become trapped in a local mode that corresponds to an incorrect clustering of data points. This mixing problem may be attributed to the incremental nature of the Gibbs sampler, which is unable to simultaneously move a group of observations to a new mixture component. Incremental updates are unlikely to move a single observation to a new mixture component because such an intermediate state has low probability. A sampling scheme which allows a group of observations to be updated simultaneously may remedy this problem, since neighbouring observations would support the formation of a new component if appropriate.

This article introduces a new Metropolis-Hastings method that avoids the problems associated with Gibbs sampling and is suitable for high-dimensional data. Typically, Metropolis-Hastings updates involve simple parametric distributions as the proposal distribution. To split mixture components, our method employs a more complex proposal distribution obtained by using a restricted Gibbs sampling scan for the latent class variables. This method is able to quickly traverse the state space and frequently visit high-probability modes because it splits or merges a group of observations in each update, thereby bypassing the incremental updates of the Gibbs sampler. Furthermore, although the proposal distribution used is complex, it does not need to be specially tailored to each model, since the same scheme can be applied to any model with a conjugate prior.

2 Gibbs sampling for this model

In this section, we present the Dirichlet process mixture model (for early references, see Ferguson 1983 and Antoniak 1974) and describe a Gibbs sampling algorithm to sample from the posterior of this model. This procedure completely utilizes the conjugacy in the model to integrate away model parameters and mixing proportions, eliminating them from the algorithm.

2.1 The Dirichlet process mixture model

We consider a hierarchical mixture model in which the observations, y_1, \dots, y_n , are modelled by a mixture of distributions having the form $F(\theta)$. There is no restriction on the dimensionality of the y_i , and the data may be categorical or quantitative. For each observation, y_i , the model parameters, θ_i , are considered to be independent draws from some mixing distribution, G . Rather than requiring G to take some parametric form, a Dirichlet process prior, a distribution over the space of distribution functions, is placed on G . This yields a mixture model of the following form:

$$\begin{aligned} y_i | \theta_i &\sim F(\theta_i) \\ \theta_i | G &\sim G \\ G &\sim DP(G_0, \alpha) \end{aligned} \tag{1}$$

where G_0 defines a baseline distribution for the Dirichlet process prior, and α is a total mass parameter that takes values greater than zero. The usual conditional

independence assumptions for a hierarchical model apply, so that the only dependencies are those that are explicitly shown. Equation (1) represents the most basic Dirichlet process mixture model. Further stages may be added to this hierarchy by placing priors on α and the parameters of G_0 .

This model may be regarded as a countably infinite mixture model (Ferguson 1983), a view that is adopted in the remainder of this article. When G is integrated over its prior distribution in equation (1), the θ_i follow a generalized Polya urn scheme (Blackwell and MacQueen 1973). The prior distribution for the θ_i may be represented in this way by the following conditional distributions:

$$\begin{aligned} \theta_1 &\sim G_0 \\ \theta_i \mid \theta_1, \dots, \theta_{i-1} &\sim \frac{1}{i-1+\alpha} \sum_{j=1}^{i-1} \delta(\theta_j) + \frac{\alpha}{i-1+\alpha} G_0 \end{aligned} \quad (2)$$

where $\delta(\theta_j)$ is the distribution which is a point mass at θ_j . The model of equation (1) has been simplified by integrating away the random distribution, G . We can represent the fact that (2) results in some of the θ_i being identical by setting $\theta_i = \phi_{c_i}$, where c_i represents the ‘‘latent class’’ associated with observation i . The Polya urn scheme for sampling the θ_i is equivalent to the following scheme for sampling the latent variables, c_i , and associated ϕ_c :

$$\begin{aligned} P(c_i = c \mid c_1, \dots, c_{i-1}) &= \frac{n_{i,c}}{i-1+\alpha}, \text{ for } c \in \{c_j\}_{j \neq i} \\ P(c_i \neq c_j \text{ for all } j < i \mid c_1, \dots, c_{i-1}) &= \frac{\alpha}{i-1+\alpha} \end{aligned} \quad (3)$$

where $n_{i,c}$ is the number of c_k for $k < i$ that are equal to c . The labelling of the indicator c_i is irrelevant in the above probabilities; all that matters is which c_i are equal to each other.

The ϕ_c are drawn independently from the initial distribution G_0 . The probabilities shown in (3) define the Dirichlet process model and are equivalent to the mixture model in equation (1). This notation will be employed in subsequent sections.

2.2 A Gibbs sampling procedure

If G_0 is a conjugate prior for F , it is straightforward to sample from the posterior distribution of the above model using Gibbs sampling. There have been several Gibbs sampling approaches proposed in the Dirichlet process mixture model literature, but we consider the procedure in which conjugacy is fully exploited, which was introduced by Neal (1992) and MacEachern (1994). This procedure integrates away the model parameters, ϕ_{c_i} . Eliminating ϕ_{c_i} simplifies the algorithm considerably, so that the state of the Markov chain for the Gibbs sampler consists only of the class indicators, c_i .

The Markov chain is initialized by setting the c_i to some initial state. The c_i are then updated via Gibbs sampling by repeatedly drawing a new value for each c_i from its conditional distribution given the others, which is proportional to the product of its conditional prior and likelihood. Because the observations are exchangeable, the conditional prior can be derived from equation (3) by considering observation i to be the last of the n observations. This yields the following conditional probabilities:

$$\begin{aligned} P(c_i = c \mid c_{-i}, y_i) &= b \frac{n_{-i,c}}{n-1+\alpha} \int F(y_i, \phi) dH_{-i,c_j}(\phi), \text{ for } c \in \{c_j\}_{j \neq i} \\ P(c_i \neq c_j \text{ for all } j \neq i \mid c_{-i}, y_i) &= b \frac{\alpha}{n-1+\alpha} \int F(y_i, \phi) dG_0(\phi) \end{aligned} \quad (4)$$

where c_{-i} represents the c_j for $j \neq i$, $n_{-i,c}$ is the number of c_j for $j \neq i$ that are equal to c , n is the number of observations, $H_{-i,c}$ is the posterior distribution of ϕ based on the prior G_0 and all observations y_j for which $j \neq i$ and $c_j = c$, $F(y_i, \phi)$ is the likelihood, and b is the appropriate normalizing constant so that the probabilities sum to one. When G_0 is a conjugate prior for F , the integrals $\int F(y_i, \phi) dG_0(\phi)$ and $\int F(y_i, \phi) dH_{-i,c}(\phi)$ can be analytically computed.

3 Split-merge Metropolis-Hastings updates

When two or more mixture components have similar parameters, Gibbs sampling can be inefficient. The Markov chain can become trapped in a local mode, in which two distinct mixture components are merged and assigned parameters which are a compromise between the two separate components. Because the Gibbs sampler incrementally updates each observation, the Markov chain must pass through a low-probability intermediate state in order to split such a component. This leads to slow convergence to the true posterior distribution and slow movement between posterior modes when the data are not sufficient to determine whether one component or two is appropriate. Here, we introduce a nonincremental Markov chain sampling method based on the Metropolis-Hastings algorithm that avoids this problem. Our algorithm is able to split or merge groups of data points, avoiding the need to pass through a low-probability intermediate state in order to make major changes.

3.1 Metropolis-Hastings updates

Our algorithm is a form of the Metropolis-Hastings algorithm (Metropolis *et al* 1953, Hastings 1970). This algorithm samples from a distribution with density $\pi(x)$ by first drawing a candidate state, x^* , according to a proposal density $q(x^*|x)$. This proposed state, x^* , is evaluated by the Metropolis-Hastings acceptance probability which is calculated as follows:

$$a(x^*, x) = \min \left[1, \frac{q(x|x^*)}{q(x^*|x)} \frac{\pi(x^*)}{\pi(x)} \right] \quad (5)$$

The next state will be set to this candidate state with probability $a(x^*, x)$. Otherwise, the new state is the same as the current state, x . These Metropolis-Hastings updates leave the posterior distribution, π , invariant and produce a valid Markov chain Monte Carlo sampling scheme provided the chain is ergodic.

As discussed by Tierney (1994), when constructing Markov chains, it is acceptable to select a transition probability at random from a set of appropriate transition probabilities. In particular, we may randomly choose amongst valid Metropolis-Hastings algorithms by randomly selecting a proposal distribution, $q(x^*|x)$. Note that when calculating the Metropolis-Hastings acceptance probability, the ratio $q(x|x^*)/q(x^*|x)$ may be calculated for the particular proposal distribution that was chosen, rather than by summing over all possible proposal distributions. Both lead to valid Metropolis-Hastings updates, but the latter calculation may be computationally infeasible.

3.2 Random split-merge proposals

First, we introduce the split-merge algorithm when the proposal distribution is based on a simple random split of the subset of observations associated with one

mixture component into two separate components, without reference to the properties of the observed data. This is the simplest version of the split-merge algorithm, which we do not expect to work well, but which illustrates the basic construction. A more elaborate version of this procedure, based on a restricted Gibbs sampling proposal, produces more sensible splits and is discussed in Section 3.3.

This split-merge approach will be applied to the conjugate Dirichlet process mixture model, in which the random distribution, G , and the model parameters, ϕ_{c_i} , are integrated away. The state of the Markov chain consists only of the mixture component indicators, c_i . The Markov chain is initialized by assigning each observation to a mixture component. Typical initial states we have used are placing all the data in the same component and assigning each observation to a different component. Below, we outline the steps for the simple random split-merge procedure.

Simple random split-merge procedure

1. Select two distinct observations, i and j , uniformly at random.
2. Let S denote the set of observations, $k \in \{1, \dots, n\}$, for which $k \neq i$ and $k \neq j$, and $c_k = c_i$ or $c_k = c_j$.
3. If items i and j belong to the same mixture component, i.e. $c_i = c_j$, then:
 - (a) Propose a new assignment of data items to mixture components, denoted as \mathbf{c}^{split} , in which component $c_i = c_j$ is split into two separate components, c_i^{split} and c_j^{split} . We define each element of the proposal vector, \mathbf{c}^{split} , as follows:
 - Let c_i^{split} be a new component such that $c_i^{split} \notin \{c_1, \dots, c_n\}$
 - Let $c_j^{split} = c_j$
 - For every observation $k \in S$, let c_k^{split} be randomly set, independently with equal probability, to either component c_i^{split} or c_j^{split}
 - For observations $k \notin S \cup \{i, j\}$, let $c_k^{split} = c_k$
 - (b) Evaluate the proposal in (a) by the Metropolis-Hastings acceptance probability $a(\mathbf{c}^{split}, \mathbf{c})$. If the proposal is accepted, \mathbf{c}^{split} becomes the next state in the Markov chain. If the proposal is rejected, the original vector, \mathbf{c} , remains as the next state.
4. Otherwise, if i and j belong to different mixture components, i.e. $c_i \neq c_j$, then:
 - (a) Propose a new assignment of data items to mixture components, denoted as \mathbf{c}^{merge} , in which components, c_i and c_j , are combined into a single component. Each element of the proposal vector, \mathbf{c}^{merge} , is assigned as follows:
 - Let $c_i^{merge} = c_j$
 - Let $c_j^{merge} = c_j$
 - For every observation $k \in S$, let $c_k^{merge} = c_j$
 - For observations $k \notin S \cup \{i, j\}$, let $c_k^{merge} = c_k$
 - (b) Evaluate the proposal in (a) by the Metropolis-Hastings acceptance probability $a(\mathbf{c}^{merge}, \mathbf{c})$. If the proposal is accepted, \mathbf{c}^{merge} becomes the next state. If the merge proposal is rejected, the original configuration, \mathbf{c} , remains as the next state.

Note that because the numerical values of the c_k are irrelevant, it does not matter in steps 3(a) and 4(a) which item, i or j , remains fixed at its original mixture component. The labels are significant only in that they distinguish which items are grouped in the same mixture component. Also note that the vectors \mathbf{c}^{split} , \mathbf{c}^{merge} , and \mathbf{c} designate which mixture component is assigned to each observation

in the data — not just to observations that are involved in the split or merge steps. However, items not associated with i , j , or set S remain unchanged and unaffected during the Metropolis-Hastings update.

The Metropolis-Hastings acceptance probability (equation 5) used in steps 3 and 4 of this procedure takes the following form:

$$a(\mathbf{c}^*, \mathbf{c}) = \min \left[1, \frac{q(\mathbf{c}|\mathbf{c}^*)}{q(\mathbf{c}^*|\mathbf{c})} \frac{P(\mathbf{c}^*)}{P(\mathbf{c})} \frac{L(\mathbf{c}^*|\mathbf{y})}{L(\mathbf{c}|\mathbf{y})} \right] \quad (6)$$

where \mathbf{c}^* is either the vector \mathbf{c}^{split} or \mathbf{c}^{merge} depending on the type of proposal. The posterior distribution, $\pi(\mathbf{c})$, in equation (5) has been expanded into a product of its factors: the prior, $P(\mathbf{c})$, and the likelihood, $L(\mathbf{c}|\mathbf{y})$, where $\mathbf{y} = (y_1, \dots, y_n)$ is the vector of observations. Note that factors not involving \mathbf{c} may be ignored.

The prior distribution, $P(\mathbf{c})$, for the entire vector \mathbf{c} will be a product over distinct $c \in \{c_1, \dots, c_n\}$ of the factors presented in equation (3). This product yields the following prior distribution:

$$P(\mathbf{c}) = \alpha^D \frac{\prod_{c \in \{c_1, \dots, c_n\}} (n_c - 1)!}{\prod_{k=1}^n (\alpha + k - 1)} \quad (7)$$

where D is the number of distinct mixture components contained in vector \mathbf{c} and n_c is the count of items belonging to mixture component c in \mathbf{c} .

Notice that the ratio of the prior distributions in equation (6) simplifies considerably because the denominator in equation (7) will cancel, as well as factors in equation (7) associated with components that are not directly involved in the Metropolis-Hastings update. For the split proposal, the prior distribution ratio reduces to the following:

$$\frac{P(\mathbf{c}^{split})}{P(\mathbf{c})} = \alpha \frac{(n_{c_i^{split}} - 1)! (n_{c_j^{split}} - 1)!}{(n_{c_i} - 1)!} \quad (8)$$

where \mathbf{c} represents the original state in which i and j belong to the same mixture component. Here, $n_{c_i^{split}}$ and $n_{c_j^{split}}$ represent the number of observations that belong to the two split mixture components. The factor of α in the ratio arises from D being one greater in \mathbf{c}^{split} than in \mathbf{c} .

Similarly, for the merge proposal, the prior ratio simplifies to:

$$\frac{P(\mathbf{c}^{merge})}{P(\mathbf{c})} = \frac{1}{\alpha} \frac{(n_{c_i^{merge}} - 1)!}{(n_{c_i} - 1)! (n_{c_j} - 1)!} \quad (9)$$

where \mathbf{c} represents the original state in which items i and j belong to separate components.

The likelihood for the vector of component indicators will be a product over the n observations:

$$L(\mathbf{c}|\mathbf{y}) = \prod_{k=1}^n \int F(y_k, \phi) dH_{k,c_k}(\phi) \quad (10)$$

where H_{k,c_k} is the posterior distribution of ϕ based on the prior G_0 and all observations y_g for which $g < k$ and $c_g = c_k$. We assume that the integral $\int F(y_k, \phi) dH_{k,c_k}(\phi)$ is analytically tractable, which is the case if G_0 is a conjugate prior. Note that when k is the first item observed from a particular component, then $H_{k,c}$ will be the prior

distribution, G_0 , since no data from that mixture component precedes item k . Alternatively, $L(\mathbf{c}|\mathbf{y})$ may be expressed as a double product over components, c , and items, $k \in \{1, \dots, n\}$, associated with each component:

$$L(\mathbf{c}|\mathbf{y}) = \prod_{c=1}^D \prod_{k: c_k=c} \int F(y_k, \phi) dH_{k,c}(\phi) \quad (11)$$

where D is the number of distinct components.

Since factors involving items associated with components not directly involved in the split proposal will cancel, the ratio of likelihoods in equation (6) reduces to the following:

$$\frac{L(\mathbf{c}^{split}|\mathbf{y})}{L(\mathbf{c}|\mathbf{y})} = \frac{\prod_{k: c_k^{split}=c_i^{split}} \int F(y_k, \phi) dH_{k,c_i^{split}}(\phi) \prod_{k: c_k^{split}=c_j^{split}} \int F(y_k, \phi) dH_{k,c_j^{split}}(\phi)}{\prod_{k: c_k=c_i} \int F(y_k, \phi) dH_{k,c_i}(\phi)} \quad (12)$$

Similarly, for the merge proposal, the ratio of likelihoods is:

$$\frac{L(\mathbf{c}^{merge}|\mathbf{y})}{L(\mathbf{c}|\mathbf{y})} = \frac{\prod_{k: c_k^{merge}=c_i^{merge}} \int F(y_k, \phi) dH_{k,c_i^{merge}}(\phi)}{\prod_{k: c_k=c_i} \int F(y_k, \phi) dH_{k,c_i}(\phi) \prod_{k: c_k=c_j} \int F(y_k, \phi) dH_{k,c_j}(\phi)} \quad (13)$$

In the first step of this procedure, the selection of observations, i and j , decides which Metropolis-Hastings proposal will be used. As a result, when calculating the acceptance probability, i and j are fixed. The probability of proposing a particular split of the items in set S from the merged state is:

$$q(\mathbf{c}^{split}|\mathbf{c}) = \left(\frac{1}{2}\right)^{n_{c_i^{split}} + n_{c_j^{split}} - 2} \quad (14)$$

Notice that $q(\mathbf{c}^{split}|\mathbf{c})$ is equivalent to $q(\mathbf{c}|\mathbf{c}^{merge})$.

The probability of proposing a merge move for the items in S from a split state is:

$$q(\mathbf{c}^{merge}|\mathbf{c}) = 1 \quad (15)$$

since there is only one way to assign all items in S to the same component. Note that $q(\mathbf{c}^{merge}|\mathbf{c})$ is equivalent to $q(\mathbf{c}|\mathbf{c}^{split})$.

It follows from equations (14) and (15) that the appropriate ratio of transition probabilities for the split proposal is:

$$\frac{q(\mathbf{c}|\mathbf{c}^{split})}{q(\mathbf{c}^{split}|\mathbf{c})} = \frac{1}{\left(\frac{1}{2}\right)^{n_{c_i^{split}} + n_{c_j^{split}} - 2}} \quad (16)$$

Similarly, the ratio of transition probabilities for the merge proposal is:

$$\frac{q(\mathbf{c}|\mathbf{c}^{merge})}{q(\mathbf{c}^{merge}|\mathbf{c})} = \left(\frac{1}{2}\right)^{n_{c_i} + n_{c_j} - 2} \quad (17)$$

Therefore, the resulting acceptance probability (6) for a split proposal is based on the product of equations (8), (12), and (16). Likewise, the acceptance probability

for a merge proposal is based on the product of equations (9), (13), and (17). By employing the Hastings (1970) version of the Metropolis (1953) algorithm when calculating the acceptance probability, we correct for the fact that the probability of proposing a particular split is smaller than the probability of proposing to merge the two resulting components.

This basic form of our procedure illustrates how we may nonincrementally update groups of observations. If a proposed split is appropriate for the data, it will likely be accepted, since neighbouring observations will lend support for the creation of a new component, bypassing the problem of being trapped in a local mode. Unfortunately, as stated earlier, we do not expect the simple random split version of this algorithm to perform well. Since components are split without reference to the observed data, the split proposals are unlikely to be appropriate, and hence are unlikely to be accepted.

3.3 Restricted Gibbs sampling split-merge proposals

Next, we describe a proposal distribution in which properties of the observed data are used to decide how to split mixture components via a restricted Gibbs sampling scan. This yields a method in which reasonable splits of components are more frequently proposed. First, as a way to validate the main algorithm, we introduce the Gibbs sampling split-merge proposal when the state immediately prior to the Gibbs sampling scan is fixed. This pre-Gibbs state will be referred to as the *launch* state. We then present a generalized version in which the launch state is itself randomly selected; in particular, we can select the launch state by conducting several “intermediate” restricted Gibbs sampling scans.

3.3.1 Restricted Gibbs sampling proposals from a fixed launch state

Here, we replace the simple random split step in our earlier procedure by a restricted Gibbs sampling scan on the component indicators, c_k , starting from a predetermined fixed launch state. The fixed state version of this algorithm is not expected to be of any particular use, except as a method to prove the validity of the subsequent random launch state algorithm. The restricted Gibbs sampling proposal distribution is more elaborate than typical Metropolis-Hastings proposals, but the proposal probabilities can still be explicitly computed. Each fixed launch state for the Gibbs sampling scan defines a particular Metropolis-Hastings algorithm, all of which are valid, since they satisfy the usual requirements, such as independence of proposals from past states.

For the split proposal, once observations i and j have been assigned to different components, other observations (i.e. those in S) that belong to the merged component will first be assigned to one of these two split components in some predetermined manner. Once this launch state, e^{launch} , is determined, one restricted Gibbs sampling scan is conducted to decide how the items in S will be allocated between the two split components. The Gibbs sampling scan is restricted in that it is only performed on a subset of the data (set S) and can assign these items to only two of the mixture components. To update a c_k in S via restricted Gibbs sampling, a new value of c_k is drawn from its (restricted) conditional distribution as follows:

$$P(c_k | c_{-k}, y_k) = \frac{n_{-k, c_k} \int F(y_k, \phi) dH_{-k, c_k}(\phi)}{n_{-k, c_i} \int F(y_k, \phi) dH_{-k, c_i}(\phi) + n_{-k, c_j} \int F(y_k, \phi) dH_{-k, c_j}(\phi)} \quad (18)$$

To simplify notation, we refer to component indicators for the launch state as c_k in equation (18). However, throughout the Gibbs sampling scan, these values (as well

as the values for the other terms) are continually modified as the c_k are incrementally updated and used for the next computation leading to c^{split} . Here, c_{-k} represents the c_g for $g \neq k$ in $S \cup \{i, j\}$, $n_{-k,c}$ is the number of c_g for $g \neq k$ in $S \cup \{i, j\}$ that are equal to c , $F(y_k, \phi)$ is the likelihood, and $H_{-k,c}$ is the posterior distribution of ϕ based on the prior G_0 and data observations y_g such that $c_g = c$ where $g \in S \cup \{i, j\}$, for which $g \neq k$. Again, when G_0 is a conjugate prior for F , the above integrals may be analytically computed.

In general, the transition probability for a full sequential Gibbs sampling scan is a product of the conditional probabilities of each individual update. The probability that c^{split} will be produced by a restricted Gibbs sampling scan starting from c^{launch} is the product of the probabilities of assigning each observation $k \in S$ to a particular split mixture component via Gibbs sampling from the fixed launch state as given by equation (18). In our algorithm, this product is the Metropolis-Hastings proposal probability, $q(c^{split}|\mathbf{c})$.

For the merge proposal, as in the simple random split-merge procedure, there is still only one way to merge items in two components to one component, so $q(c^{merge}|\mathbf{c}) = 1$. However, to obtain the corresponding probability, $q(\mathbf{c}^{merge})$, we need to calculate the probability of generating the original split state from the fixed launch state in one Gibbs sampling scan. This is done in the same way as for the split proposal, except that no actual sampling is performed since the ‘‘split’’ state is already known.

As in the simple random split case, to obtain the Metropolis-Hastings acceptance probability, the appropriate split or merge proposal distribution ratio (now based on restricted Gibbs sampling) is substituted into equation (6). The prior and likelihood ratios in equation (6) remain as shown in Section 3.2.

Since only one scan of Gibbs sampling is conducted, we do not expect that the allocation of items between the two components has reached equilibrium. Because the Metropolis-Hastings proposal distribution can take any form and still produce a valid algorithm, lack of convergence will not invalidate this algorithm. However, it is quite likely that the proposed splits using a single iteration of Gibbs sampling may not be that sensible. We would like to improve the proposals further so that the splits proposed are appropriate for the data.

3.3.2 Restricted Gibbs sampling proposals from a random launch state

Every fixed launch state for the algorithm of the previous section produces a valid Metropolis-Hastings update. As was discussed in Section 3.1, we may select a Markov chain transition at random from the set of valid transitions (Tierney 1994). Therefore, a launch state for the restricted Gibbs sampling scan may be chosen at random from the set of all fixed states. We could, for example, choose the launch state uniformly at random. However, if only a single scan of Gibbs sampling is performed from a random launch state, it may still lead to an unreasonable assignment of observations to the two mixture components.

To achieve more reasonable splits, several intermediate restricted Gibbs sampling scans are performed before the final scan. When calculating the split proposal probability, the result of the last intermediate Gibbs sampling scan is considered the random launch state, from which the restricted Gibbs sampling transition probability is explicitly calculated. We would prefer to incorporate all of the intermediate Gibbs sampling scans in the proposal distribution, but summing probabilities over all of these intermediate states is computationally infeasible. Although equilibrium will probably not be reached after only a few restricted Gibbs sampling scans, the

clustering of observations between the two mixture components will be a better reflection of the actual attributes of the data than is produced by a single scan of Gibbs sampling.

Split proposal probabilities are calculated in the same way as for the fixed launch state Gibbs sampling proposals (equation 18). For the merge proposal, to obtain $q(\mathbf{c}|\mathbf{c}^{merge})$, the same intermediate Gibbs sampling operations that are performed when proposing a split must be conducted here to arrive at a launch state, even though no actual split is performed. The Gibbs sampling transition probability is calculated from the launch state (which is the last intermediate Gibbs sampling state) to the original split state. These operations are necessary in order to produce the correct proposal ratios.

We could modify the algorithm by replacing the intermediate restricted Gibbs sampling scans by Markov chain updates of some other type. However, replacing the final Gibbs sampling scan (from the launch state) with some other update would be possible only if the transition probability for this update could be calculated. Below, the procedure for the restricted Gibbs sampling split-merge update with a random launch state is summarized.

Restricted Gibbs sampling split-merge procedure

1. Select two distinct observations, i and j , uniformly at random.
2. Let S denote the set of observations, $k \in \{1, \dots, n\}$, for which $k \neq i$ and $k \neq j$, and $c_k = c_i$ or $c_k = c_j$.
3. Define the launch state, \mathbf{c}^{launch} , that will be used to compute Gibbs sampling probabilities. If $c_i = c_j$, then let c_i^{launch} be set to a new component such that $c_i^{launch} \notin \{c_1, \dots, c_n\}$ and let $c_j^{launch} = c_j$. Otherwise, if $c_i \neq c_j$, then let $c_i^{launch} = c_i$ and $c_j^{launch} = c_j$. For every $k \in S$, set c_k^{launch} to either of the distinct components, c_i^{launch} or c_j^{launch} , as follows:
 - Select an initial state by randomly setting, independently with equal probability, c_k^{launch} to either c_i^{launch} or c_j^{launch} .
 - Modify \mathbf{c}^{launch} by performing t intermediate restricted Gibbs sampling scans or some other type of Markov chain update.
4. If items i and j are in the same mixture component, i.e. $c_i = c_j$, then:
 - (a) Propose a new assignment of data items to mixture components, denoted as \mathbf{c}^{split} , in which component $c_i = c_j$ is split into two separate components, c_i^{split} and c_j^{split} . Define each element of the proposal vector, \mathbf{c}^{split} , as follows:
 - Let $c_i^{split} = c_i^{launch}$ (note that $c_i^{launch} \notin \{c_1, \dots, c_n\}$)
 - Let $c_j^{split} = c_j^{launch}$ (which is the same as c_j)
 - For every observation $k \in S$, let c_k^{split} be set to either component c_i^{split} or c_j^{split} by conducting **one** final Gibbs sampling scan from the launch state, \mathbf{c}^{launch}
 - For observations $k \notin S \cup \{i, j\}$, let $c_k^{split} = c_k$
 - (b) Calculate the proposal probability, $q(\mathbf{c}^{split}|\mathbf{c})$, by computing the Gibbs sampling transition probability from the launch state, \mathbf{c}^{launch} , to the final proposed state, \mathbf{c}^{split} . The Gibbs sampling transition probability is the product, over $k \in S$, of the probabilities of setting each c_k^{split} to its final value in the final Gibbs sampling scan.
 - (c) Evaluate the proposal by the Metropolis-Hastings acceptance probability $a(\mathbf{c}^{split}, \mathbf{c})$. If the proposal is accepted, \mathbf{c}^{split} becomes the next state in the Markov chain. If the proposal is rejected, the original vector, \mathbf{c} , remains as the next state.

5. Otherwise, if i and j are in different mixture components, i.e. $c_i \neq c_j$, then:
 - (a) Propose a new assignment of data items to mixture components, denoted as \mathbf{c}^{merge} , in which components, c_i and c_j , are combined into a single component. Assign each element of the proposal vector, \mathbf{c}^{merge} , as follows:
 - Let $c_i^{merge} = c_j$
 - Let $c_j^{merge} = c_j$
 - For every observation $k \in S$, let $c_k^{merge} = c_j$
 - For observations $k \notin S \cup \{i, j\}$, let $c_k^{merge} = c_k$
 - (b) Calculate the proposal probability, $q(\mathbf{c}|\mathbf{c}^{merge})$, by computing the Gibbs sampling transition probability from the launch state, \mathbf{c}^{launch} , to the original split configuration, \mathbf{c} . The Gibbs sampling transition probability is the product, over $k \in S$, of the probabilities of setting each c_k in the original split state to its original value in a (hypothetical) Gibbs sampling scan from the launch state.
 - (c) Evaluate the proposal by the Metropolis-Hastings acceptance probability $a(\mathbf{c}^{merge}, \mathbf{c})$. If the proposal is accepted, \mathbf{c}^{merge} becomes the next state. If the merge proposal is rejected, the original configuration, \mathbf{c} , remains as the next state.

3.4 Cycling Metropolis-Hastings and Gibbs sampling updates

The split-merge Metropolis-Hastings algorithm produces a Markov chain that leaves the posterior distribution invariant. The Markov chain is also irreducible, since for any statistical model in which the data have non-zero prior probability, there is a non-zero probability that the chain started from any initial state will assign every observation to a separate mixture component as a result of a series of split moves. Further, except for some degenerate models, the Markov chain is aperiodic, since there is a non-zero probability that the chain will remain in its current state (i.e. at least some split or merge proposals have a non-zero probability of being rejected). The split-merge algorithm is therefore ergodic.

Even though the above procedure is ergodic and produces nonincremental splits or merges of components, further improvements in convergence may be obtained by combining this algorithm with traditional Gibbs sampling. This procedure addresses the problem of making major changes in the allocation of items by moving observations as a cluster during a single iteration. However, it may take longer to move a single observation between components. In this situation, a “fine tuning” approach is required, which the regular Gibbs sampling scan can provide. Consequently, we propose combining these two algorithms by alternately performing a Metropolis-Hastings update and a full scan of Gibbs sampling. By doing this, we exploit the nonincremental (major) changes from the Metropolis-Hastings step, and the incremental (minor) refinement from the Gibbs sampling step.

Tierney (Section 2.4, 1994) notes that if Markov chain transition kernels are applied in cycles, and one of the kernels is ergodic, then the cycle kernel is not guaranteed to be ergodic. However, in our case, since both the Metropolis-Hastings and Gibbs sampling steps have a non-zero probability of leaving the state unchanged, applying each transition in turn will produce an ergodic Markov chain.

We can tune this algorithm by modifying the number of Metropolis-Hastings updates and the number of final Gibbs sampling scans in each full iteration. As expected, by increasing the values for both of these tuning parameters, convergence (measured in full iterations) is improved, but at the cost of computation time per iteration.

4 Example: Bernoulli data with a Beta prior

In this section, we empirically compare the split-merge procedure (and its variants) to Gibbs sampling. We consider a categorical mixture model, in which the data, $\mathbf{y} = (y_1, \dots, y_n)$, are independent and identically distributed, such that each observation, y_i , has m Bernoulli attributes, (y_{i1}, \dots, y_{im}) . Given the class, c_i , that observation i belongs to, the item's attributes are independent of each other. Neal (1992) considered a similar model when examining the performance of the Gibbs sampling procedure discussed in Section 2. For simplicity of exposition, we consider only dichotomous attributes, but the model and algorithms easily generalize to categorical attributes with more than two values.

4.1 The statistical model

The observations, $y_i = (y_{i1}, \dots, y_{im})$, are multivariate Bernoulli, giving the following likelihood:

$$F(y_i, \theta_i) = \prod_{h=1}^m \theta_{ih}^{y_{ih}} (1 - \theta_{ih})^{1-y_{ih}} \quad (19)$$

The parameters, θ_i , of component i give the probabilities of each attribute having the value one. Each such probability is given a Beta distribution prior with parameters (β_1, β_0) . Under G_0 , which is the prior over the vector $\theta = (\theta_1, \dots, \theta_m)$, the θ_h are independent. (Note that here we use subscripts to denote different attributes rather than different observations.) The density for G_0 is:

$$P(\theta) = \prod_{h=1}^m \left(\frac{\Gamma(\beta_{1,h} + \beta_{0,h})}{\Gamma(\beta_{1,h}) \Gamma(\beta_{0,h})} \theta_h^{\beta_{1,h}-1} (1 - \theta_h)^{\beta_{0,h}-1} \right) \quad (20)$$

where $\beta_{0,h}$ and $\beta_{1,h}$ are greater than zero.

Because this is a conjugate prior for $F(y_i, \theta_i)$, the model parameters may be integrated away. To update c_i via Gibbs sampling, a new value of c_i is drawn from its conditional distribution (equation 4), which for this model is the following, when the Beta prior and Bernoulli likelihood are substituted:

$$\begin{aligned} \text{for } c \in \{c_j\}_{j \neq i}, P(c_i = c \mid c_{-i}, y_i) &= b \frac{n_{-i,c}}{n-1+\alpha} \prod_{h=1}^m \frac{\sum_{k \neq i} \delta(c_k, c) \delta(y_{kh}, y_{ih}) + \beta_{y_{ih},h}}{n_{-i,c} + \beta_{0,h} + \beta_{1,h}} \\ P(c_i \neq c_j \text{ for all } j \neq i \mid c_{-i}, y_i) &= b \frac{\alpha}{n-1+\alpha} \prod_{h=1}^m \frac{\beta_{y_{ih},h}}{\beta_{0,h} + \beta_{1,h}} \end{aligned} \quad (21)$$

The delta function $\delta(x, y)$ is equal to one if $x = y$ and zero otherwise. The term $\sum_{k \neq i} \delta(c_k, c) \delta(y_{kh}, y_{ih})$ counts the number of observations associated with component c that match y_i with respect to attribute h . The second formula gives the probability for setting c_i to a new mixture component that is currently not assigned to any other observation. In both equations, b is the factor that normalizes the distribution to sum to one.

For the Metropolis-Hastings acceptance probability in equation (6), the prior is calculated as shown in equation (7). The appropriate ratio of the transition probabilities based on restricted Gibbs sampling is obtained by using the first formula

in equation (21). The likelihood (equation 11) based on the Bernoulli-Beta model is as follows:

$$L(\mathbf{c}|\mathbf{y}) = \prod_{c=1}^D \prod_{k:c_k=c} \prod_{h=1}^m \frac{\sum_{i<k} \delta(c_i, c) \delta(y_{ih}, y_{kh}) + \beta_{y_{kh},h}}{n_{k,c} + \beta_{0,h} + \beta_{1,h}} \quad (22)$$

4.2 The synthetic data sets

Our primary goal is to partition observations into appropriate latent classes using the Bernoulli-Beta Dirichlet process mixture model. Computationally, this classification problem becomes more difficult as the dimensionality increases and as the sets of attributes that distinguish the various components become more similar in structure. We examine this difficulty by considering two simulated data sets, in which the number of attributes is increased so that the different components appear more alike as the dimensionality increases.

The data are composed of five equally-probable mixture components, in which each component produces a distribution over m dichotomous attributes. To maintain uniformity between the examples, $n = 100$ observations were produced for each example, and 20 observations were generated from each of the five mixture components.

Data for the two examples were randomly generated from mixture distributions, in which the mixture components are distinguished by the first four attributes, which for consistency, have been kept constant in both examples. The true mixture distribution for the first example with six attributes is shown in Table 1. The second example is similar to the first except that it has eighteen attributes. Dimensionality is increased by simply replicating the distribution for the last attribute, which makes the components more similar, and thereby, more difficult to distinguish. Note the intentional asymmetry in the construction of the mixture components, in which the first three components are more similar than the last two components. This is intended to test whether the split-merge algorithms can handle “three-way” splits.

Table 1: True mixture distribution for Example 1.

c	$P(c_i = c)$	$P(y_{ih} = 1 c_i = c), h = 1, \dots, 6$
1	0.2	.95 .95 .95 .95 .95 .95
2	0.2	.05 .05 .05 .05 .95 .95
3	0.2	.95 .05 .05 .95 .95 .95
4	0.2	.05 .05 .05 .05 .05 .05
5	0.2	.95 .95 .95 .95 .05 .05

For the following demonstrations of the algorithms, the Dirichlet process parameter, α , is set to one. A small value of α implies that the number of mixture components present in the data is likely to be small. The $\beta_{0,h}$ and $\beta_{1,h}$ parameters for the Beta prior distribution have also been set to one. These priors may not be realistic, but for consistency, these values are fixed at one during the simulations. In actual problems, α and the β 's would be set by prior knowledge or given higher-level priors.

4.3 Performance of the algorithms

For each example, the Gibbs sampling algorithm was compared to five versions of the split-merge algorithm: Simple Random Split, Split-Merge (0,1,0), Split-Merge (0,1,1), Split-Merge (5,1,0), and Split-Merge (5,1,1). The first number in parentheses is the number of intermediate Gibbs sampling scans to reach the launch state, the second is the number of Metropolis-Hastings updates in a single iteration, and the third is the number of complete Gibbs sampling scans after the final Metropolis-Hastings update. For each algorithm, all observations were assigned to the same mixture component for the initial state, and each algorithm was run for 2000 iterations. The performance of each algorithm was evaluated by examining trace plots and the computation time per iteration.

The first example is the simplest. It is relatively low-dimensional (six attributes) and has five well-separated mixture components. From the trace plots (not shown), it appears that all of the algorithms except for the Simple Random Split have appropriately separated the data into the five mixture components. As discussed earlier, the Simple Random Split is not expected to converge rapidly, even in simple situations, because the split proposals are usually nonsensical. Gibbs sampling, Split-Merge (0,1,1) and (5,1,1) seem to have short burn-in times and mix equally well. Similarity in performance is also confirmed by approximately equal autocorrelation times. In this simple problem, Gibbs sampling is successful in correctly splitting the items amongst the five components, so the split-merge algorithms are not necessary.

Example 2 is a high dimensional problem (eighteen attributes), in which the posterior distribution, given the priors assigned, gives substantial probability to configurations with (mainly) four or five components. The trace plots (not shown) indicate that Gibbs sampling remains in an incorrect split that is not typical of the true posterior distribution for the entire 2000 iteration run. This happens because the mixture components in this example are quite similar, so incremental creation of a new component via Gibbs sampling is quite rare. If each item is initially assigned to a different mixture component (plot not shown), Gibbs sampling splits the data into five components immediately, but takes roughly 1000 iterations to move to the four-component configuration, showing that it mixes poorly between the four and five component configurations. Split-Merge (5,1,1) separates the observations into the proper configuration immediately and mixes well between the four and five components. Split-Merge (5,1,0) also mixes between four and five components, but the minor adjustments are slow. The two split-merge algorithms without any intermediate Gibbs sampling scans find the four and five component configurations, but are stuck in the four-component split for a long time. This is a result of non-optimal Metropolis-Hastings split proposals.

5 Discussion

The split-merge Metropolis-Hastings procedure has been shown to be an improvement over traditional Gibbs sampling in high-dimensional problems in which mixture components are similar. The cycled split-merge version that includes both intermediate Gibbs sampling scans and a full overall Gibbs sampling scan is the most successful split-merge variation. Computation time per iteration is greater than for Gibbs sampling, but in situations where Gibbs sampling is unable to arrive at the correct stationary distribution in any reasonable length of time, this burden is clearly acceptable. The quality of the proposals can be controlled by varying

the number of intermediate Gibbs sampling scans. Implementing this method is relatively simple and does not become more difficult in higher dimensions. It is straightforward to apply this method to any conjugate model, including normal mixture models for real-valued data with the conjugate normal-inverse gamma priors for the mean and variance. We have implemented the split-merge algorithm for conjugate normal mixture models when the variance is known and observed similar improvements as described here.

Currently, we are extending the algorithm to handle non-conjugate models, in which the model parameters cannot be analytically integrated away. We believe that this is possible, but how well the resulting algorithm will work remains to be seen.

The full version of this article is available at:

<http://www.utstat.toronto.edu/~sonia/paper1.ps>

References

- Antoniak, C. E. (1974) "Mixtures of Dirichlet processes with applications to Bayesian nonparametric problems", *Annals of Statistics*, vol. 2, pp. 1152-1174.
- Blackwell, D. and MacQueen, J. B. (1973) "Ferguson distributions via Pólya urn schemes", *Annals of Statistics*, vol. 1, pp. 353-355.
- Escobar, M. D. and West, M. (1995) "Bayesian density estimation and inference using mixtures", *Journal of the American Statistical Association*, vol. 90, pp.577-588.
- Ferguson, T. S. (1983) "Bayesian density estimation by mixtures of normal distributions", in H. Rizvi and J. Rustagi (editors) *Recent Advances in Statistics*, pp. 287-303, New York: Academic Press.
- Green, P. J. and Richardson, S. (2001) "Modelling heterogeneity with and without the Dirichlet process", *Scandinavian Journal of Statistics*, vol. 28, pp. 355-375.
- Hastings, W. K. (1970) "Monte Carlo sampling methods using Markov chains and their applications", *Biometrika*, vol. 57, pp. 97-109.
- MacEachern, S. N. (1994) "Estimating normal means with a conjugate style Dirichlet process prior", *Communications in Statistics: Simulation and Computation*, vol. 23, pp. 727-741.
- Metropolis, N., Rosenbluth, A. W., Rosenbluth, M. N., Teller, A. H., and Teller, E. (1953) "Equation of state calculations by fast computing machines", *Journal of Chemical Physics*, vol. 21, pp. 1087-1092.
- Neal, R. M. (1992) "Bayesian mixture modeling", in C. R. Smith, G. J. Erickson, and P. O. Neudorfer (editors) *Maximum Entropy and Bayesian Methods: Proceedings of the 11th International Workshop on Maximum Entropy and Bayesian Methods of Statistical Analysis*, Seattle, 1991, pp. 197-211, Dordrecht: Kluwer Academic Publishers.
- Neal, R. M. (2000) "Markov chain sampling methods for Dirichlet process mixture models", *Journal of Computational and Graphical Statistics*, vol. 9, pp. 249-265.
- Tierney, L. (1994) "Markov chains for exploring posterior distributions" (with discussion), *Annals of Statistics*, vol. 22, pp. 1701-1762.