

# NIH Public Access

**Author Manuscript** 

*Comput Stat.* Author manuscript; available in PMC 2011 May 5.

## Published in final edited form as:

Comput Stat. 2010 March ; 25(1): 17-38. doi:10.1007/s00180-009-0159-7.

# Bayesian model-based tight clustering for time course data

#### Yongsung Joo,

Department of Statistics, Dongguk University, Seoul 100-715, Korea, yongsungjoo@dongguk.edu

# G. Casella, and

Department of Statistics, University of Florida, Gainesville, FL 32611, USA

#### J. Hobert

Department of Statistics, University of Florida, Gainesville, FL 32611, USA

# Abstract

Cluster analysis has been widely used to explore thousands of gene expressions from microarray analysis and identify a small number of similar genes (objects) for further detailed biological investigation. However, most clustering algorithms tend to identify loose clusters with too many genes. In this paper, we propose a Bayesian tight clustering method for time course gene expression data, which selects a small number of closely-related genes and constructs tight clusters only with these closely-related genes.

#### Keywords

Bayesian cluster analysis; Tight clustering; Time course gene expression; Microarray

# **1** Introduction

Clustering methods can be categorized into heuristic and model-based frameworks. Methods in heuristic frameworks identify clusters based on non-probabilistic measures. The K-means (Hartigan and Wong 1979) and hierarchical (Johnson and Wichern 2002) algorithms belong to this framework. Methods in model-based frameworks cluster objects based on probabilistic measures. As one of the most popular methods in this framework (Basford and McLachlan 1985; Basford et al. 1997), there is the mixture model

$$f(\mathbf{Y}|\theta) = \prod_{i=1}^{n} \sum_{k=1}^{K} \xi_k f(\mathbf{Y}_i|\theta_k),$$
(1)

where  $\mathbf{Y}_i$  is the response variable of the *i*th object, *K* is the number of components,  $\xi_k$  is the mixing probability of component *k* and  $f(\mathbf{Y}_i|\mathbf{\theta}_k)$  is the probability density function of component *k* with parameter  $\mathbf{\theta}_k$ . This model has been studied for microarray data analyses in numerous papers (Ghosh and Chinnaiyan 2002, McLachlan et al. 2002; Datta and Datta 2003; Ouyang et al. 2004). Recently, the product partition model (Crowley 1997) have been proposed as alternative model-based approaches using the cluster likelihood:

© Springer-Verlag 2009

Correspondence to: Yongsung Joo.

(2)

 $f(\boldsymbol{Y}|\boldsymbol{\theta}, \boldsymbol{\omega}) = \prod_{k=1}^{c(\boldsymbol{\omega})} f(\boldsymbol{Y}_{c_k}|\boldsymbol{\theta}_k)$ 

where  $\omega$  is a fixed unknown partition of *n* objects,  $c(\omega)$  is the number of clusters within  $\omega$ ,  $\mathscr{C}_k$  is a set of object indices for cluster *k*, and  $\mathbf{Y}_{\mathscr{C}_k}$  is the data of objects in cluster *k*. Note that partition  $\omega$  is a free parameter, which should be estimated, and  $\omega$  determines  $c(\omega)$ . Then,

 $\bigcup_{k=1}^{c(\omega)} \mathscr{C}_k = \{1, 2, \dots, n\}$ , and  $\mathscr{C}_i \cap \mathscr{C}_j = \emptyset$  when  $i \neq j$ . These methods assume that the data vectors are partitioned into  $c(\omega)$  clusters according to  $\omega$  and the clusters are independent of each other. While the mixture model (1) is constructed with a known number of clusters, the cluster likelihood (2) contains a partition  $\omega$  as a parameter to be estimated. In other words, pre-specification of the number of clusters is not needed in the cluster likelihood approach.

As microarray technology became more easily available, biologists can measure gene expressions consecutively over time and examine temporal changes of these expressions. Naturally, many new statistical methods have been developed to cluster genes based on these temporal changes (profiles) in both heuristic (Peddada et al. 2003; Hakamada et al. 2006; Lukashin and Fuchs 2001) and model-based (Schliep et al. 2003; Luan and Li 2003; James and Sugar 2003; Tseng and Wong 2005; Ma et al. 2006; Leng and Muller 2006; Ng et al. 2006) frameworks as follows.

- 1. Heuristic framework: Peddada et al. (2003) grouped profiles into, so-called, clusters of inequality profiles. For example, suppose a cluster contains profiles with monotonically increasing temporal trends, which can be characterized with inequalities among means at each time point. Similarly, various types of clusters are pre-specified with inequalities, then profiles are clustered based on bootstrap-based criterion. Hakamada et al. (2006) developed a method that cluster profiles based euclidian distances. Lukashin and Fuchs (2001) applied K-means algorithm to temporal profiles.
- 2. Model-based framework: Leng and Muller (2006) applied functional discriminant analysis considering profiles as independent realizations of a smooth stochastic process. Luan and Li (2003), James and Sugar (2003), Ng et al. (2006), Ma et al. (2006) developed the mixture of mixed effect models to cluster temporal profiles. In their models,  $f(Y_i|\theta_k)$  in (1) is specified with a mixed effect model. Particularly, James and Sugar (2003) emphasized application of their model to sparsely and irregularly measured time course data. The Bayesian objective function approach in Booth et al. (2008) and the hidden Markov model in Schliep et al. (2003) are developed based on the cluster likelihood function (2).
- **3.** Costa et al. (2004), Thalamuthu et al. (2006) and Ma et al. (2006) compared clustering methods for time course data using simulation studies.

As the main example of this paper, we analyze the corneal wound healing data, in which expressions of 646 genes are measured twice at irregular 12 time points, using 24 rats (= 2 replicates  $\times$  12 time points). The goal of this paper is to identify clusters that contain a small number of genes with very similar temporal patterns. We consider that two types of gene sets, closely-related and weakly-related genes, are included in a microarray experiment either intentionally or unintentionally. If a gene has a close relationship with any other gene(s), we will call it a closely-related gene. Otherwise, we will call it a weakly-related gene are weakly-related. These weakly-related genes tend to increase noise in search of the optimal partition (or clusters) without providing significant amounts of information (Tseng and

Wong 2005). Therefore, direct application of conventional methods will provide large and loose clusters that consist of both closely- and weakly-related genes. However, this is not a desirable result, because biologists typically want to conduct further biological research on a small number of closely-related genes after gene expressions are explored with microarray analyses. To overcome this problem, Tseng and Wong (2005) proposed so-called *tight* clustering method for *cross-sectional* data. The main idea of Tseng and Wong (2005) is to construct tight clusters only with closely-related genes, which is a small portion of the whole data set. Their algorithm can be summarized as follows:

Step 1. With the given number of clusters  $\kappa$ , apply K-means algorithm to subsets (e.g. 70% of genes) of the original data from resampling.

Step 2. Construct candidate tight clusters with genes tends to be together in these resampled subsets.

Step 3. Apply Step 1 and 2 with different  $\kappa$  in a certain range.

Step 4. Identify the final tight clusters that tend to be stable even when  $\kappa$  changes.

Performance of this tight clustering method (Tseng and Wong 2005) was demonstrated by Thalamuthu et al. (2006). In this paper, we propose a new tight clustering algorithm for *time course* gene expression data. Our tight clustering algorithm selects closely-related genes that have high relevance probabilities and then identifies clusters of only closely-related genes using a Bayesian objective function approach (Booth et al. 2008).

In Sect. 2, the Bayesian model and the objective function (Booth et al. 2008) are described in detail with an example of the corneal wound experiment. In Sects. 3 and 4, we discuss a stochastic search algorithm that maximizes a Bayesian objective function and provides simulation studies on this search algorithm. In Sect. 5, we explain how to calculate relevance probabilities and propose the tight clustering algorithm for time course data. To distinguish from tight clustering, we will call algorithms that do not employ the idea of tight clustering, *plain* clustering methods. In Sect. 6, our tight clustering method is applied to the corneal wound data.

#### 2 Bayesian modelling and objective function

We look at an example of the corneal wound healing data to explain the Bayesian model and the objective function by Booth et al. (2008). In the experiment, 2 replicates of 646 gene expressions were measured at each of 12 time points (day 0, 1, 2, 3, 4, 5, 6, 7, 14, 21, 42 and 98) with a corneal wound. In our cluster analysis, the time variable is relabeled with orders 1, 2, ..., 12. Justification of relabeling will be discussed in Sect. 6. Because we are interested in the temporal changes of gene expressions rather magnitudes of expressions, gene profiles are centered at the mean of each profile for our analysis.

Averages of two replicates are calculated for each gene at each time point. Then, average expressional profiles of 646 genes are drawn in Fig. 1. Although the patterns of profiles are not easily distinguishable, it seems that there are at least two clusters with increasing and decreasing gene expression patterns in the right end part of the profiles. Also, temporal patterns are not simple enough to use a parametric regression approach and there does not seem to be any periodic pattern. Therefore, within the clustering model, we use the penalized regression spline to explain the mean temporal trend of gene profiles within each cluster.

Given  $\omega$ , let  $\theta = (\theta_k)_{k=1}^{c(\omega)}$  denote a set of cluster-specific parameter vectors  $\theta_k$ . Then, the marginal posterior distribution of  $\omega$  is

$$\pi(\omega|\mathbf{Y}) \propto \int f(\mathbf{Y}|\theta,\omega) \pi(\theta|\omega) \pi(\omega) d\theta$$
  
= 
$$\int \prod_{k=1}^{c(\omega)} \left\{ f(\mathbf{Y}_{c_k}|\theta_k) \pi(\theta_k|\omega) \right\} \pi(\omega) d\theta$$
  
=: Obj (\omega), (3)

where  $f(\mathbf{Y}|\mathbf{\theta}, \omega)$  is the sampling distribution of the whole data set,  $f(\mathbf{Y}_{\mathscr{C}_k}|\mathbf{\theta}_k)$  is the cluster specific sampling distribution, and  $\pi(\cdot)$  denotes a prior distribution. In the analysis of corneal wound data,  $f(\mathbf{Y}_{\mathscr{C}_k}|\mathbf{\theta}_k)$  is set to be the normal probability density function, of which the mean is the cluster-specific penalized regression spline. Booth et al. (2008) proposes obtaining the optimal partition  $\omega^*$  that maximizes the marginal posterior probability  $\pi(\omega|\mathbf{Y})$ . Because the normalizing constant of  $\pi(\omega|\mathbf{Y})$  is difficult to calculate,  $Obj(\omega)$  is used as the actual objective function in optimizing  $\pi(\omega|\mathbf{Y})$ . In Sect. 2.1,  $f(\mathbf{Y}_{\mathscr{C}_k}|\mathbf{\theta}_k)$  is specified in detail with the penalized regression spline. In Sect. 2.2,  $\pi(\mathbf{\theta}_k|\omega)$  and  $\pi(\omega)$  are specified. Finally, in Sect. 2.3, the Bayesian objective function is calculated for our clustering model.

#### 2.1 Penalized regression spline for profiles within a cluster: $f(Y_{\mathscr{C}_k}|\theta_k)$

In our clustering model, all profiles in cluster k are assumed to have a common smooth underlying trend and independent and identically distributed normal errors with a common variance. Detailed explanation of modeling profiles within a cluster is given as follows. Denote the gene expressions by

$$\boldsymbol{Y} = \left(\boldsymbol{Y}_1^T, \ldots, \boldsymbol{Y}_i^T, \ldots, \boldsymbol{Y}_n^T\right)^T,$$

where  $\mathbf{Y}_{i} = (\mathbf{Y}_{i1}^{T}, \dots, \mathbf{Y}_{ij}^{T}, \dots, \mathbf{Y}_{ir}^{T})$ ,  $\mathbf{Y}_{ij}^{T} = (Y_{ij1}, \dots, Y_{ijt}, \dots, Y_{ijp})$ , and  $Y_{ijt}$  is the expression of gene *i* in the *j*th replication at time point *t*. Similarly, define the time variables  $\mathbf{x}_{i} = (\mathbf{x}_{i1}^{T}, \dots, \mathbf{x}_{ij}^{T}, \dots, \mathbf{x}_{ir}^{T})$ ,  $\mathbf{x}_{ij}^{T} = (x_{ij1}, \dots, x_{ijt}, \dots, x_{ijp})$ , and  $x_{ijt} = t$ , and error terms  $\varepsilon_{i=}(\varepsilon_{i1}^{T}, \dots, \varepsilon_{ij}^{T}, \dots, \varepsilon_{ij}^{T})$ ,  $\varepsilon_{ij}^{T} = (\varepsilon_{ij1}, \dots, \varepsilon_{ijt}, \dots, \varepsilon_{ijp})$  and  $\varepsilon_{ijt} \sim N(0, \sigma_{k}^{2})$ . To explain a temporal profile in cluster *k*, we use the penalized regression spline:

$$Y_{ij} = \beta_{0k} + \beta_{1k} \mathbf{x}_{ij} + \dots + \beta_{qk} \mathbf{x}_{ij}^{q} + \sum_{l=1}^{L} u_{lk} (\mathbf{x}_{ij} - \tau_l)_{+}^{q} + \boldsymbol{\varepsilon}_{ij}$$
  
=  $\mathbf{X} \beta_k + \mathbf{Z} \mathbf{U}_k + \boldsymbol{\varepsilon}_{ij},$  (4)

where  $\mathbf{X} = (1, \mathbf{x}_{ij}, \dots, \mathbf{x}_{ij}^q), \beta_k = (\beta_{0k}, \beta_{1k}, \dots, \beta_{qk})^T, \mathbf{Z} = (z_{ij1}, \dots, z_{ijl}, \dots, z_{ijl}), z_{ijl} = (\mathbf{x}_{ij} - \tau_l)_+^2$  with knots  $\tau_l$ 's,  $m_+ = \max(0, m), \mathbf{U}_k = (u_{1k}, \dots, u_{Lk})^T, \mathbf{0}_p$  is the column vector with p zeros,  $\mathbf{I}_p$  is

the  $p \times p$  identity matrix, and  $\varepsilon_{ij} \sim MVN(\mathbf{0}_p, \sigma_k^2 \mathbf{I}_p)$ . If profiles consist of a short time course data, such as 3 time points, the regression spline is not recommended. Then, the simple linear or quadratic regression had better be used instead of the regression spline. The time variable  $\mathbf{x}_{ii}$  is the same for every gene *i* and replicate *j* in the corneal wound data. For

notational simplicity, assume that  $\mathscr{C}_k = \{1, ..., n_k\}$ . Also, let  $Y_{\mathscr{C}_k} = (Y_1^T, ..., Y_{n_k}^T)^T \mathbf{X}_{\mathscr{C}_k} = \mathbf{1}_{n_k r}$  $\otimes \mathbf{X}$  and  $\mathbf{Z}_{\mathscr{C}_k} = \mathbf{1}_{n_k r} \otimes \mathbf{Z}$ , where  $\mathbf{1}_{n_k r}$  is a column vector of ones with  $n_k r$  elements. Then, within cluster *k*, the penalized regression spline method estimates parameters by minimizing

To implement this penalization in the Bayesian framework, two approaches have been widely used. First, the Bayesian Lasso approach of Tibshirani (1996) uses double exponential priors for  $U_k$ , which makes finding the posterior distribution mode equivalent to minimization of (5). Second, Ruppert et al. (2003) suggests using the mixed model, of which BLUP (best linear unbiased prediction) is equivalent to the estimates of the penalized regression spline. We use the second approach in this paper by changing a fixed effect  $U_k$  in

(4) to a cluster-specific random effect  $\mathbf{U}_k \sim N(0, \lambda^2 \sigma_k^2 \mathbf{I}_L)$  where  $\mathbf{I}_L$  is the  $L \times L$  identity matrix.

For data analysis and simulation studies in this paper, we employ a flexible quadratic regression spline by setting q = 2 and L = p - 2 (one knot at each interior time point). Because the cluster memberships are unknown, it is difficult to get a good prior information on the mean temporal trend of profiles in each cluster. Therefore, the clustering model should contain a flexible spline function with a large number of knots so that it can explain any temporal trend. Also, to prevent a unnecessary wiggly fit, we constrained the influence of knots with the penalty function (Ruppert et al. 2003).

Let  $\mathbf{J}_r = \mathbf{1}_r \mathbf{1}_r^T$ . Then, the linear mixed model for cluster k can be expressed as:

where

$$\boldsymbol{\varepsilon}_{\mathscr{C}_{k}} = \left(\boldsymbol{\varepsilon}_{1}^{T}, \dots, \boldsymbol{\varepsilon}_{n_{k}}^{T}\right)^{T} \sim N(0, \sigma_{k}^{2} \mathbf{I}_{n_{k} r p}), \boldsymbol{\varepsilon}_{\mathscr{C}_{k}}^{'} \sim N(0, \boldsymbol{\Sigma}_{\mathscr{C}_{k}}), \boldsymbol{\Sigma}_{\mathscr{C}_{k}} = \sigma_{k}^{2} \boldsymbol{\Omega}_{\mathscr{C}_{k}} \text{ and } \boldsymbol{\Omega}_{\mathscr{C}_{k}} = \mathbf{I}_{n_{k}} \otimes \mathbf{I}_{r} \otimes \mathbf{I}_{p} + \mathbf{J}_{n_{k}} \otimes \mathbf{J}_{r} \otimes \boldsymbol{\lambda}^{2} \mathbf{Z}_{\mathscr{C}_{k}} \mathbf{Z}_{\mathscr{C}_{k}}^{T}$$

Then, within cluster k, the likelihood function is

$$f\left(\mathbf{Y}_{c_{k}}|\beta_{k},\sigma_{k}^{2}\right) = \int f\left(\mathbf{Y}_{c_{k}},\mathbf{U}_{k}|\beta_{k},\sigma_{k}^{2}\right) d\mathbf{U}_{k} = \frac{1}{(2\pi)^{n_{k}pr/2}} \frac{1}{\left|\sum_{c_{k}}\right|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{Y}_{c_{k}}-\mathbf{X}_{c_{k}}\beta_{k})^{T}\sum_{c_{k}}^{-1}(\mathbf{Y}_{c_{k}}-\mathbf{X}_{c_{k}}\beta_{k})\right\}$$

# 2.2 Priors: $\pi(\omega)$ and $\pi$ ( $\beta$ , $\sigma^2|\omega$ )

As for the partition parameter  $\omega$ , we use Crowley's prior (Crowley 1997):

$$\pi(\omega) \propto \varrho^{c(\omega)} \prod_{k=1}^{c(\omega)} (n_k - 1)!, \tag{7}$$

where  $n_k$  is the number of genes in cluster k and  $\rho(>0)$  is the tuning parameter for the size of clusters. Large values of  $\rho$  makes the prior give high probabilities to partitions with a large number of clusters.

For 
$$\beta = (\beta_k)_{k=1}^{c(\omega)}$$
 and  $\sigma^2(\sigma_k^2)_{k=1}^{c(\omega)}$ , we use a non-informative prior,

$$\pi\left(\beta,\sigma^{2}|\omega\right) = \prod_{k=1}^{c(\omega)} \pi\left(\beta_{k},\sigma_{k}^{2}|\omega\right) \propto \prod_{k=1}^{c(\omega)} \left(1/\sigma_{k}^{2}\right)^{\alpha+1}.$$
(8)

#### 2.3 Bayesian objective function: $Obj(\omega)$

The marginal posterior distribution of  $\omega$  is

π

$$\begin{aligned} (\omega|\mathbf{Y}) &\propto \iint f\left(\mathbf{Y}|\boldsymbol{\beta},\sigma^{2},\omega\right)\pi\left(\boldsymbol{\beta},\sigma^{2}|\omega\right)\pi\left(\omega\right)d\boldsymbol{\beta}d\sigma^{2} \\ &=\pi\left(\omega\right)\iint\prod_{k=1}^{c(\omega)}f\left(\mathbf{Y}_{c_{k}}|\boldsymbol{\beta}_{k},\sigma_{k}^{2}\right)\pi\left(\boldsymbol{\beta}_{k},\sigma_{k}^{2}|\omega\right)d\boldsymbol{\beta}d\sigma^{2} \\ &=\pi\left(\omega\right)\prod_{k=1}^{c(\omega)}\iint f\left(\mathbf{Y}_{c_{k}}|\boldsymbol{\beta}_{k},\sigma_{k}^{2}\right)\pi\left(\boldsymbol{\beta}_{k},\sigma_{k}^{2}|\omega\right)d\boldsymbol{\beta}_{k}d\sigma_{k}^{2}. \end{aligned}$$

Also,

$$\begin{split} &\iint f\left(\mathbf{Y}_{c_{k}}|\boldsymbol{\beta}_{k},\sigma_{k}^{2}\right)\pi\left(\boldsymbol{\beta}_{k},\sigma_{k}^{2}|\boldsymbol{\omega}\right)d\boldsymbol{\beta}_{k}d\sigma_{k}^{2} \\ &= \iint \frac{1}{(2\pi)^{n_{k}pr/2}}\frac{1}{|\sum_{c_{k}}|^{1/2}}\exp\left\{-\frac{1}{2}(\mathbf{Y}_{c_{k}}-\mathbf{X}_{c_{k}}\boldsymbol{\beta}_{k})^{T}\sum_{c_{k}}^{-1}(\mathbf{Y}_{c_{k}}-\mathbf{X}_{c_{k}}\boldsymbol{\beta}_{k})\right\}\times\left(1/\sigma_{k}^{2}\right)^{\alpha+1}d\boldsymbol{\beta}_{k}d\sigma_{k}^{2} \\ &= \frac{2^{\alpha}\pi^{(q+1)/2}\Gamma(\boldsymbol{\nu}_{k},2)}{s_{k}^{\nu/2}|\mathbf{\Omega}_{c_{k}}|^{1/2}|\mathbf{X}_{c_{k}}^{T}\mathbf{\Omega}_{c_{k}}^{-1}\mathbf{X}_{c_{k}}|^{1/2}}, \end{split}$$

where

 $\nu_{k} = n_{k} p - q - 1 + 2\alpha, S_{k} = \left( \mathbf{Y}_{\mathscr{C}_{k}} - \mathbf{X}_{\mathscr{C}_{k}} \widehat{\beta}_{k} \right)^{T} \mathbf{\Omega}_{\mathscr{C}_{k}}^{-1} \left( \mathbf{Y}_{\mathscr{C}_{k}} - \mathbf{X}_{\mathscr{C}_{k}} \widehat{\beta}_{k} \right), \text{ and } \widehat{\beta}_{k} = \left( \mathbf{X}_{\mathcal{C}_{k}}^{T} \mathbf{\Omega}_{\mathscr{C}_{k}}^{-1} \mathbf{X}_{\mathscr{C}_{k}} \right)^{-1} \left( \mathbf{X}_{\mathscr{C}_{k}}^{T} \mathbf{\Omega}_{\mathscr{C}_{k}}^{-1} \mathbf{Y}_{\mathscr{C}_{k}} \right)$ . Thus,

$$\pi(\omega|\mathbf{Y}) \propto \varrho^{c(\omega)} \prod_{k=1}^{c(\omega)} \frac{(n_k - 1)! 2^{\alpha} \pi^{(q+1)/2} \Gamma(\nu_k/2)}{S_k^{\nu_k/2} |\mathbf{\Omega}_{C_k}|^{1/2} |\mathbf{X}_{C_k}^T \mathbf{\Omega}_{C_k}^{-1} \mathbf{X}_{C_k}|^{1/2}}$$
  
=: Obj(\omega) (9)

Note that the marginal posterior distribution (9) varies with the scale of **Y**. For example, when the scale of **Y** changes to a**Y** 

$$\pi(\omega|a\mathbf{Y}) \propto \pi(\omega|\mathbf{Y})a^{c(\omega)((q+1)-2\alpha)}.$$

Therefore, to make (9) invariant to the scale of **Y**, we set  $\alpha = (q + 1)/2$ .

# 3 Stochastic search for the optimal partition

Because  $\omega$  is not a continuous variable, many popular numeric search algorithms, i.e. the Newton–Rapson method, may not be applied. Also, the space of  $\omega$  is very large, because the Bell number ( $B_n$ , the number of all possible partitions) grows rapidly with *n* (super-exponentially). For example, when *n* is 10, the Bell number is approximately 10<sup>5</sup>. When *n* is 75, the Bell number becomes approximately 10<sup>78</sup>, which is a rough estimate for the number

of atoms in the observable universe. The corneal wound data has n = 646 genes, which, in a practical sense, have an infinite number of possible partitions. Therefore, searching for the optimal partition  $\omega^*$  is a challenging problem.

To search for the optimal partition, Booth et al. (2008) proposed to generate random partitions from the posterior function  $\pi(\omega|\mathbf{Y})$ , which is proportional to our objective function Obj( $\omega$ ), and select a partition that has the highest Obj( $\omega$ ). This algorithm tends to generate many partitions that are close to the mode of the posterior distribution, which makes the search algorithm efficient. This is so-called Markov Chain Monte Carlo (MCMC) optimization (Jerrum and Sinclair 1996). If the number of objects is not large or only a small number of partitions have high posterior probability  $\pi(\omega|\mathbf{Y})$ , then the optimal partition can be easily found. Otherwise, the algorithm may work slowly or may find only a suboptimal partition within a reasonable time. However, according to our experience, these suboptimal partitions are close enough to the optimal partition in practical applications.

Here is a detailed description on the MCMC optimization algorithm (Booth et al. 2008). Within the optimization algorithm, a biased random walk (Metropolis-Hastings algorithm) generates the targeting Markov chain of random partitions from  $\pi(\omega|\mathbf{Y})$ . Suppose that the biased random walk is iterated *R* times. Let  $\omega_i$  be the partition in the *i*<sup>th</sup> iteration,  $c(\omega_i)$  be the number of clusters in  $\omega_i$ ,  $\tilde{\omega}_i^*$  be the partition with the highest value of the objective function

during the first *i* iterations, and  $\tilde{\omega}^* \left(=\tilde{\omega}_R^*\right)$  be the best partition that is found by an application of the optimization algorithm. Then, consider the following algorithm:

#### Algorithm 1

Step 1. Choose an initial partition  $\omega_1$  and set  $\tilde{\omega}_1^* = \omega_1$  and i = 1.

Step 2. Generate a candidate partition  $\omega'$ :

- If  $c(\omega_i)=1$ , select one uniform random object out of *n* and make it as a new singleton cluster (one profile in one cluster). This makes two clusters with 1 and n-1 objects.
- If  $c(\omega_i) \ge 2$ , select one uniform random object.
- If the selected object is a singleton cluster, move it to one of  $c(\omega_i) 1$  clusters with probability  $1/(c(\omega_i) 1)$ .
- Otherwise, move it to one of other clusters with probability  $1/c(\omega_i)$  or make its own singleton cluster with probability  $1/c(\omega_i)$ .

Step 3. Accept  $\omega'$  with probability  $min(1, \pi(\omega'_i)/\pi(\omega))$ . If accepted,  $\omega_{i+1} = \omega'$ . Otherwise,  $\omega_{i+1} = \omega_i$ .

Step 4. If  $Obj(\omega_{i+1}) > Obj(\tilde{\omega}_i^*)$ ,  $\tilde{\omega}_{i+1}^* = \omega_{i+1}$ . Otherwise,  $\tilde{\omega}_{i+1}^* = \tilde{\omega}_i^*$ . Set i = i + 1 and repeat Steps 2, 3 and 4.

Note that, if *R* is not large enough,  $\omega^*$  may not be the same as the optimal partition  $\omega^*$ . In this paper, Algorithm 1 is applied three times for each data set of interest. If all chains provide the same  $\omega^*$ , then  $\omega^*$  is considered as  $\omega^*$ .

# 4 Simulation studies on stochastic search algorithm

Simulation studies are conducted with the corneal wound data, to check the convergence-indistribution of the biased random walk in Algorithm 1 and to examine the needed number of iterations before Algorithm 1 finds the optimal partition.

#### 4.1 Convergence of the biased random walk

To examine the convergence or performance of the biased random walk algorithm, it is necessary to start chains with fair and non-informative initial partitions. As one of the most reasonable initial partitions, we may consider a uniform random partition that is drawn randomly with a probability of  $1/B_n$ . However, the generation of a uniform random partition is a challenging problem when the partition space is very large. In this subsection, we describe an algorithm to generate uniform random partitions with a large *n*, and demonstrate the convergence of chains using the corneal wound data.

**Generation of uniform random partitions with a large n**—Let  $\mathbb{P}_n$  be the set of all possible partitions of  $\mathbb{N}_n \coloneqq \{1, ..., n\}$ . Also, let  $\Pi$  denote a uniform random partition such that  $P(\Pi = \pi) = 1/B_n$  for all  $\pi \in \mathbb{P}_n$ , so  $\Pi$  has the uniform distribution on  $\mathbb{P}_n$ . Here we describe a method of simulating a random uniform partition  $\Pi$ .

Let *M* be a random variable on the set  $\{1, 2, ..., n\}$  with probabilities given by

$$P(M=m) = \frac{m^n}{B_n} \left[ \frac{1}{m!} \sum_{s=0}^{n-m} \frac{(-1)^s}{s!} \right]$$
(10)

for m = 1, 2, ..., n. Pitman (1997) gave an algorithm for drawing a value of  $\Pi$ , which goes as follows. First, draw a value of M = m from (10), then randomly distribute *n* balls with labels  $\mathbb{N}_n$  into the *m* different urns. Note that some of the urns may end up empty. After excluding empty urn(s), the resulting partition has the uniform distribution on  $\mathbb{P}_n$ .

Unfortunately, this method does not work well with large n because computation of the distribution (10) requires the evaluation of large factorials, which results in numerical difficulties. However, it is possible to circumvent this problem by approximating (10) to a high degree of accuracy and drawing M from this approximation. Based on this this idea, we propose a new algorithm as follows (See the Appendix for detailed calculation of this algorithm).

#### Algorithm 2:

Step 1. Choose  $\epsilon > 0,$  which controls the degree of approximation. Here we use  $\epsilon = 10^{-30}$ 

Step 2. Calculate  $l_m = n \log m - \log \Gamma(m + 1)$  for  $m \in \{1, 2, ..., n\}$  and set

$$l^* = \max \{l_1, l_2, \ldots, l_n\}.$$

Step 3. Find

$$N = \inf \left\{ m > n: l^* - l_m > \log \left( 2^{n+1} n! / \varepsilon \right) \right\}.$$

Step 4. Draw  $M^*$  with probabilities given by

$$P(M^*=m) = \frac{\exp\{l_m\}}{\sum_{j=1}^{N} \exp\{l_j\}} = \frac{\exp\{l_m - l^*\}}{\sum_{j=1}^{N} \exp\{l_j - l^*\}}$$
(11)

for  $m \in \{1, 2, ..., N\}$ .

Uniform random partitions are generated with n = 646 as in the corneal wound data. Then, the number of clusters in each partition is counted to draw the histogram in Fig. 2. Most uniform random partitions have from 120 to 143 clusters with mean 131.3 and variance  $4.6^2$ . The gray line in Fig. 2 shows the distribution of the number of bins, *m* in (11), which has mean 132.3 and variance  $4.7^2$ .

**Convergence-in-distribution with the corneal wound data**—When multiple chains are simulated with a good MCMC algorithm, they get mixed well regardless of the initial partitions. Using the biased random walk, five chains are started from four non-informative partitions (two uniform random partitions, one partition with n = 646 singleton clusters and one partition that has only one cluster for all objects) and one informative partition (with 200 clusters found using the K-means algorithm).

The Bayesian objective function (9) has two parameters,  $\lambda$  and  $\varrho$ , that should be set before Algorithm 1 starts. To estimate  $\lambda$ , which is the smoothing parameter or the ratio of two variances in the linear mixed model (6), we clustered profiles into an arbitrary number of groups, K = 40, with K-means algorithm, and fitted the linear mixed model to the grouped data assuming a homogeneous error variance across clusters. From this procedure, we got an estimate  $\hat{\lambda} = 1.31$ . Even though  $\hat{\lambda}$  depends on K,  $\hat{\lambda}$  is robust to K as long as K is not close to 1 or n. For example,  $\hat{\lambda} = 1.27$  when K = 20, and  $\hat{\lambda} = 1.33$  when K = 200. Also, the resulting optimal partition is usually not sensitive to a small variation of  $\hat{\lambda}$ , such as  $\pm 0.2$ . For the tuning parameter in Crowley's prior,  $\log(\varrho) = 8$  is chosen. Detailed discussion about the selection of  $\log(\varrho)$  is in Sect. 6.

Simulation histories of the five chains are shown with the number of clusters on Fig. 3a, and with the objective function  $Obj(\omega)$  on Fig. 3b. Because most uniform random partitions have a similar number of clusters, around 131, and similar values of  $Obj(\omega_1) \approx -2,300$ , two initial partitions are overlapped on Fig. 3(a, b). Note that uniform random partitions have very low initial values of the objective function, which indicates that they are very non-informative partitions. On the contrary, the initial partition from K-means algorithm has the highest objective function value,  $Obj(\omega_1) = 4,597$ . Even when five chains have very different initial partitions, Fig. 3 shows that the biased random walk results in excellent agreement at convergence, and good mixing over the stationary distribution after  $10^6$  iterations. Therefore, we consider that  $2 \times 10^6$  iterations are enough to be a burn-in period in analyzing the corneal wound data.

The convergence-in-distribution does not guarantee that Algorithm 1 can find the optimal partition  $\omega^*$  within a reasonable time. For example, in Table 1, the algorithm ( $R = 3 \times 10^6$ ) is applied to the corneal wound data twice with  $\log(\varrho) = 0, 3, 5, 8, 10, 13, 15$  and 20. Regardless of the value of  $\log(\varrho)$ , two chains provide different  $\tilde{\omega}^*$ 's, which have different Obj ( $\tilde{\omega}^*$ )'s and different numbers of clusters. However, two Obj ( $\tilde{\omega}^*$ ) are very close to each other, considering that this search is conducted in almost infinite space of a discrete parameter. More details of Table 1 will be discussed in Sect. 6.

#### 4.2 Number of iterations before the optimal partition is found

To examine the relationship between the number (n) of gene profiles to be clustered and the number  $(R^*)$  of iterations before the optimal partition  $\omega^*$  is found by Algorithm 1, data sets are simulated by randomly selecting *n* gene profiles from the corneal wound data without replacements. Then, three long chains with  $5 \times 10^7$  iterations are simulated for each data set using the biased random walk. If these three chains have the same  $\tilde{\omega}^*$ , we consider that the optimal partition  $\omega^* = \tilde{\omega}^*$  is found for a simulated data set. The average of three  $R^*$ 's,

 $Avg(R^*)$ , is reported in Table 2. For example, for the first simulated data set with randomly selected n = 20 profiles, the same  $\omega^*$  is found by three chains after 6,220 iterations on average. In total, 5 data sets are simulated with n = 20. The other four data sets have  $Avg(R^*) = 1,637, 663, 4,439$  and 934. All  $Avg(R^*)$ 's are less than  $10^4$  when n = 20. When n increases to 50,  $Avg(R^*)$ 's have magnitude 10<sup>6</sup> or 10<sup>7</sup> with the first four simulated data sets. With the fifth simulated data set, the same optimal partition is not found by three chains, which indicates that more than  $5 \times 10^7$  iterations are needed. When n = 100, the optimal partition is not found with any data set within  $5 \times 10^7$  iterations. It seems that  $Avg(R^*)$ increases super exponentially as the Bell number does  $(B_{n=100} \approx 10^{115})$ . Our C program spent about 8 h to iterate  $5 \times 10^7$  times for each data set with n = 100. Therefore, it is practically impossible to find the optimal partition of 646 gene profiles within a reasonable time. There can be two major reasons that make the optimization algorithm slow. First, the discrete space of partitions is too large with  $n \ge 100$ . Second, there seems to be many suboptimal partitions, of which  $Obj(\omega)$ 's are close to  $Obj(\omega^*)$ . This may happen when many weakly-related genes are included in the data set and cause loose clusters. Also, weaklyrelated genes have a tendency to change cluster memberships easily during MCMC iterations, because they don't make significant contributions to the objective function. In the next section, we propose a tight clustering algorithm that selects a small number of closelyrelated genes and applies Algorithm 1 to only these genes.

# 5 Tight clustering algorithm

Recall that, during the biased random walk of Algorithm 1, a candidate partition  $\omega'$  is generated by moving one object (gene) in the current partition  $\omega_i$ . If a gene is closely-related with other genes in the current cluster and build a tight cluster, then this gene have low chances to move during iterations. We will consider a gene to be stable, when a gene doesn't move at all or change only a small number of times over MCMC iterations of Algorithm 1.

Let the relevance probability (RP) of a pair of genes be the probability of an event that two genes in a pair belong to one cluster together in a random partition. In the Bayesian objective function approach framework, RP can be estimated easily by counting how often two genes are together in a chain of the biased random walk. Because closely-related genes tend to stay together within a cluster over the course of MCMC iterations, they have high RP's. Our definition of relevance probability is slightly different from the original concept by Hartigan (1991) where the relevance probability is the probability having a set Sof genes (objects) as a cluster in a random partition. It is inappropriate to use Hartigan's definition of RP for finding closely-related genes because there can be too many possible clusters, S For example, when Hartigan's definition of RP is used to find closely-related genes in corneal wound data, RPs of  $2^{646} - 2 \approx 2.9 \times 10^{194}$  (the number of all possible non-empty subsets) possible clusters must be calculated for a simulation chain. When our definition of RP is

used, RPs of only  $\binom{646}{2}$  =208, 335 pairs must be calculated.

Here, we propose the following tight clustering algorithm:

#### Algorithm 3

Step 1. Select  $log(\varrho^*)$  that makes a small number of genes stable in a chain of the biased random walk.

Step 2. Apply Algorithm 1 with  $log(\rho^*)$  and estimate RP's of all possible pairs.

Step 3. Apply Algorithm 1 with only closely-related genes (RP $\geq \eta$ ) to construct tight clusters.

In Step 1, we suggest selecting a value of log ( $\varrho^*$ ) that makes a small number of genes stable because this makes a small number of genes have high RP's in Step 2. However, an analytic optimization for log ( $\varrho^*$ ) is very difficult. Therefore, we suggest running Algorithm 1 with different log( $\varrho$ ) values to calculate the number of stable genes. The cutoff value  $\eta$  of RP in Step 3 can be considered as a tuning parameter that determines the tightness of clusters. If  $\eta$ is close to 1, Algorithm 3 will select a small number of genes that construct very tight clusters. If  $\eta$  is close to 0, Algorithm 3 will select most genes without constructing tight clusters. Based on our experience with real data and simulation studies, we suggest setting  $\eta$ = 0.80 for reasonably tight clusters.

To check performance of the tight clustering algorithm (Algorithm 3), we conducted simulation studies with the following six true clusters:

 $C_{1}:y_{i}=\sin(x_{i}/2) \times \log(x_{i})^{2} + \varepsilon_{1i}$   $C_{2}:y_{i}=-1+\exp(x_{i}/5) \times \cos((x_{i}-2)/2)^{2} + \varepsilon_{2i}$   $C_{3}:y_{i}=2-x_{i}+\exp(x_{i}/5) \times \cos(x_{i}/2)^{2} + \varepsilon_{3i}$   $C_{4}:y_{i}=1+1.5 \times \sin(x_{i}/2) + \varepsilon_{4i}$   $C_{5}:y_{i}=\exp(x_{i}/5) + \varepsilon_{5i}$   $C_{6}:y_{i}=3-x_{i}+\cos(x_{i}/2)^{2} + \varepsilon_{6i},$ 

where  $\varepsilon_{1i}$ ,  $\varepsilon_{2i}$ ,  $\varepsilon_{3i}$ ,  $\varepsilon_{4i}$  and  $\varepsilon_{5i}$  have independent and identical normal distributions with mean 0 and variance  $\sigma_4^2$ ,  $\varepsilon_{6i}$  has the independent and identical normal distribution with mean 0 and variance  $\sigma_{B}^{2}$ , and  $x_{i} = 1, 2, ..., 6$ . Four sets of simulations are considered with different  $\sigma_{A}$  and  $\sigma_B$ . In each set, we simulated 60 profiles (10 profiles from each cluster) 100 times. Then, our tight clustering and plain Bayesian objective function approaches (Booth et al. 2008) with various tuning parameter  $log(\rho)$ 's are compared in Table 3. In Step 1 of the tight clustering algorithm, we chose  $log(\varrho^*)$  for each simulated data set by examining the number of stable genes when  $\log(\rho) = -10, -8, -6, -4, -2$  and 0. In general,  $\log(\rho)$  is a very important parameter in the clustering algorithm, controlling the number of clusters. However, when two values of  $\log(\rho)$  differed by less than 2 in our simulation experience, the same or very similar clusters were found. In this paper, we will consider a pair of profiles to be correctly clustered if the pair is in the same cluster of the simulation model and is grouped together by the clustering algorithm or if two profiles in the pair are in different clusters of the simulation model and are not grouped together. Otherwise, we will consider the pair to be incorrectly clustered. There are  $N_{\text{pairs}} = n(n-1)/2$  possible pairs, where *n* is the number of profiles to be clustered. In the plain clustering, n = 60 and  $N_{\text{pairs}} = 1,770$  for every simulated data set. However, in the tight clustering, n and  $N_{\text{pairs}}$  varies with simulated data because only closely-related profiles are interested for clustering and the number of closely-related profiles varies with data. We define that the percentage classification error rates of pairs is:

$$ER=100 \times \frac{\sum_{i}^{N_{\text{pairs}}} I_i}{N_{\text{pairs}}},$$

where  $i = 1, ..., N_{\text{pairs}}$  and

 $I_i = \begin{cases} 1, & \text{if ith pair is incorrectly clustered;} \\ 0, & \text{otherwise.} \end{cases}$ 

See Table 3. For the first set of simulations, we set  $\sigma_A = \sigma_B = 0.5$ . These are relatively small standard deviations that make clustering algorithms find correct clusters easily. Tight clustering method has the lowest error rate of 6.0% with a standard error of 0.8%. Among plain clusterings, the lowest error rate, 8.8%, is gained when  $log(\rho) = -10$ . Recall that a higher  $\rho$  makes the algorithm choose a large number of small clusters. When  $\log(\rho) = 5$ , the algorithm provides the same optimal partition of n singleton clusters for all 100 simulated data sets. A simulated data set has  $270 (= 45 \times 6)$  pairs that have two profiles from the same true cluster. Therefore, when all profiles are separated as singleton clusters, the error rate becomes  $15.3\% (= 100 \times 270/1,770)$ . When we increase the standard deviations to  $\sigma_A = \sigma_B =$ 0.7, the tight clustering algorithm achieves the lowest error rate again, 13.7%. When the standard deviations are as large as  $\sigma_A = \sigma_B = 1.0$ , it is practically impossible to get good clustering results because most parts of the true clusters overlap each other. Therefore, the error rates are high for both tight and plain clustering methods. In this case, the plain clustering method with  $log(\rho) = 5$  works the best, only because it separates all profiles as singletons and gains an error rate of 15.3%. However, this is not a desirable result because it does not provide any non-singleton cluster. As the last set of simulations, we considered  $\sigma_A$ = 0.7 and  $\sigma_B$  = 3.0, which makes cluster  $\mathscr{C}_6$  relatively diffused. The tight clustering method worked the best in this case. Compared to simulations with  $\sigma_A = \sigma_B = 0.7$ , it seems clustering was easier because  $\mathscr{C}_6$  with  $\sigma_B = 3.0$  is not similar to other clusters with  $\sigma_A = 0.7$ .

Overall, we found that the tight clustering method has lower error rates compared to plain clustering when true clusters are reasonably separated. This result can be easily expected because tight clustering uses only closely-related profiles (genes) and excludes weakly-related profiles that may cause erroneous clusters.

#### 6 Analysis of the corneal wound data

In the corneal wound experiment, gene expressions are measured at day = 0, 1, 2, 3, 4, 5, 6, 7, 14, 21, 42 and 98 because expressions are expected to change more intensely in the first week and then get stabilized in the later part of experiment. In other words, a less smooth pattern is expected in the first week. When spline knots are located with equal interval on the real time scale (days), this change of smoothness will not be reflected. In this paper, we suggest relabeling time points with 1,...,12 before the cluster analysis. This may not be the optimal solution, but an easy and reasonable solution to the unequal smoothness problem. Here, this argument is illustrated with an example of a gene profile in our data. Using the same smoothing parameter ( $\lambda = 1.3$ ) and knots at all interior time points, the gene profile is fitted on real scale in Fig. 4a and on relabeled equal interval scale in Fig. 4c. Also, to make comparison of these fits easier, Fig. 4a is redrawn with equal interval time scale in Fig. 4(b, c) is redrawn with real time scale in Fig. 4d. When Fig. 4(a, d) (or, equivalently, Fig. 4b, c) is compared, the difference between two spline fits are negligible from 1st (day=0) to 7th time points (day=6). However, the later part of the profile (i.e. from 7th to 12th time points) is overfitted in Fig. 4a, while it is fitted with a proper smooth line in Fig. 4d. Similar patterns are found with other gene profiles. If analyst wants to use the real time scale, this overfitting problem may be fixed by carefully choosing the number and locations of knots.

Now on, we will provide a detailed discussion about application of the tight clustering algorithm to the corneal wound data.

#### Step 1

Simulation studies are conducted with  $log(\varrho) = 0,3,5,8,10,13,15$  and 20. For consistency of comparisons, the same initial partition with 200 clusters from the K-means algorithm is used for every simulation. Algorithm 1 is applied twice with each  $log(\varrho)$  and the movements of genes are monitored in the last  $10^6$  iterations out of  $3 \times 10^6$  total iterations. Table 1 describes

Obj  $(\tilde{\omega}^*)$ ,  $c(\tilde{\omega}^*)$ , and the number of stable genes, which do not move at all during the last  $10^6$  iterations. For example, from the first run with  $\log(\varrho) = 0$ , we got Obj  $(\tilde{\omega}^*) = 5,280$ ,  $c(\tilde{\omega}^*) = 27$  and 58 stable genes out of 646. When  $\log(\varrho)$  gets larger, the number of clusters  $c(\tilde{\omega}^*)$  increases because a larger value of  $\log(\varrho)$  makes the objective function support a larger number of clusters, and hence smaller numbers of genes per cluster. In Fig. 5, the number of stable genes is plotted with  $\log(\varrho)$ . It shows that a relatively small number of genes are stable (do not switch clusters in the last  $10^6$  iterations) when  $\log(\varrho)$  is between 3 and 13. When  $\log(\varrho)$  is as large as 20, most genes have a strong tendency to stay as singleton clusters during the iterations. This makes many genes appear stable. When  $\log(\varrho)$  is as small as 0, even weakly-related genes may stay together easily within large clusters. This also causes many genes to appear stable. Figure 5 shows that the number of stable genes is the minimum around  $\log(\rho) = 8$ . Therefore,  $\log(\rho^*) = 8$  is selected for the corneal wound data.

#### Step 2 and Step 3

After applying Algorithm 1 to the whole data with  $2 \times 10^7$  iterations, the 1,169 pairs are found to have RP≥ 0.80. Then, Algorithm 1 is applied 3 times to only the 139 closelyrelated genes that compose these 1,169 pairs. The same  $\omega^*$  is found easily in all three applications with less than  $10^6$  iterations. Recall that, when a data set of 100 profiles is simulated in Sect. 4.2, the optimal partition  $\omega^*$  is not found even after  $5 \times 10^7$  iterations. The convergence is fast with closely-related genes because only a small number of genes are considered in tight clustering and close relationships among selected genes make Obj( $\omega^*$ ) much higher than any Obj( $\omega$ ), where  $\omega \neq \omega^*$ . See Fig. 6 for the final result of our analysis. Each cluster has a distinctive pattern and gene profiles build tight clusters. Cluster 3 and 9 seem similar, but scales distinguish them.

#### 7 Concluding remark

A typical purpose of cluster analysis with microarray data is to identify a small number of closely-related genes that biologists can study further in future studies. Also, weakly-related genes make the clustering algorithm work slowly and build large and loose clusters. Therefore, we propose to select closely-related genes using relevance probabilities and get tight clusters only with closely-related genes.

In our tight clustering, the stochastic search algorithm (Algorithm 1) is used three times for different purposes. In the first two applications with all gene profiles, the stochastic algorithm is used for the selection of  $log(\varrho^*)$  and the calculation of relevance probabilities of all pairs, rather than for the search of the optimal partition. Because the stochastic algorithm is implemented by MCMC simulation, the calculation of relevance probabilities could be done easily. Instead of the Bayesian objective function approach, if the mixture model is used in tight clustering, it will be difficult to calculate relevance probabilities because the mixture model requires prior knowledge on a fixed number of components, *K*. Even though *K* can be determined with a model selection criterion (i.e. BIC), it is difficult to measure uncertainties in determining *K* and take account of it when calculating relevance probabilities.

In addition to closely-related genes, if a small number of weakly-related genes are interesting for biological reasons, then they can also be included in Step 3 of tight clustering algorithm. We expect that inclusion of these weakly-related genes will not make significant impact on the optimal partition of tight clusters, because closely-related genes build a stable structure of clusters.

In this paper, we clustered expressions of 646 genes that are preselected by biologists. However, in general, most microarray experiments generate 10,000–40,000 gene

expressions at once. When the number of genes are such high, computational time can be considered as a limitation of our approaches. For example, if genes are measured twice at 12 time points as in corneal wound data set, it takes about a week with 3.0Hz PC for our tight clustering method to handle 1,000 or less genes (or objects). Therefore, we recommend preselecting gene profiles that vary largely over time, using the one-way ANOVA model with time as the covariate. For example, we may preselect genes that have the 1,000 highest *F*-values. A similar approach was also proposed in Peddada et al. (2003).

# Appendix

It turns out that the Pitman's algorithm (Pitman 1997) still works if we take *M* to be the random variable on the set  $\mathbb{N} \coloneqq \{1, 2, 3, ...\}$  with probabilities given by

$$P(M=m) = \frac{m^n e^{-1}}{B_n m!}$$

for m = 1, 2, 3,... We now describe a method for sampling from (an arbitrarily good approximation of) M without having to calculate  $B_n$ . Specifically, given  $\varepsilon > 0$ , we find a positive integer  $N = N(\varepsilon)$  such that  $P(M \ge N + 1) < \varepsilon$  and we approximate M with a random variable  $M^*$  with probabilities given by  $P(M^* = m) \propto m^n/m!$  for  $m \in \{1, 2, ..., N\}$ . If we choose  $\varepsilon$  small enough, there will be no discernible difference between draws from M and draws from  $M^*$ .

Fix N > 1 and note that

$$\sum_{j=N+1}^{\infty} P(M=j) = \sum_{j=1}^{\infty} P(M=N+j),$$

and

$$\begin{array}{ll} \frac{P\left(M=N+j\right)}{P\left(M=N\right)} &= \frac{(N+j)^n}{(N+j)!} \frac{N!}{N^n} \\ &= (1+j/N)^n \frac{N!}{(N+j)!} \\ &= (1+j/N)^n \frac{1}{(N+j)(N+j-1)\cdots(N+1)} \\ &\leq (1+j)^n \frac{1}{N^j} \\ &\leq \frac{(2j)^n}{N^j}. \end{array}$$

Therefore,

$$P(M=N+j) \le \frac{(2j)^n}{N^j}P(M=N),$$

so

Thus, if we can find *N* such that  $P(M = N) < \varepsilon/(2^{n+1}n!)$ , then

$$P(M \ge N+1) = \sum_{j=N+1}^{\infty} P(M=j) < \varepsilon.$$

It's easy to show that P(M = m) is decreasing for m > n. Define  $l_m = n \log m - \log \Gamma(m + 1)$  so that  $\log P(M = m) = l_m - c$ , where  $c = 1 + \log B_n$ . Now define  $l^* = \max\{l_1, l_2, ..., l_n\}$  and note that, for any  $m \in \mathbb{N}$ ,

$$P(M=m) = \frac{\exp\{l_m\}}{\sum_{j=1}^{\infty} \exp\{l_j\}} \le \frac{\exp\{l_m\}}{\exp\{l^*\}} = \exp\{l_m - l^*\}.$$

Therefore, if  $l^* - l_m > \log(2^{n+1} n! / \varepsilon)$ , we have

$$P(M=m) \le \exp\{l_m - l^*\} \le \varepsilon/(2^{n+1}n!).$$

Define

$$N = \inf \left\{ m > n: l^* - l_m > \log \left( 2^{n+1} n! / \varepsilon \right) \right\}$$

We then calculate the probabilities of  $M^*$  as

$$P(M^*=m) = \frac{\exp\{l_m\}}{\sum_{j=1}^{N} \exp\{l_j\}} = \frac{\exp\{l_m - l^*\}}{\sum_{j=1}^{N} \exp\{l_j - l^*\}}$$

for  $m \in \{1, 2, ..., N\}$ .

# References

Basford KE, McLachlan GJ. Likelihood estimation with normal mixture models. Appl Stat. 1985; 34:282–289.

Basford KE, Greenway DR, McLachlan GJ, Peel D. Standard errors of fitted means under normal mixture models. Comput Stat. 1997; 12:1–17.

- Booth J, Casella G, Hobert J. Clustering using objective functions and stochastic search. J R Stat Soc B. 2008; 70(1):119–140.
- Costa IG, Carvalho FAT, Souto MCP. Comparative analysis of clustering methods for gene expression time course data. Genet Mol Biol. 2004; 27:623–631.
- Crowley EM. Product partition models for normal means. J Am Stat Assoc. 1997; 92:192–198.
- Datta S, Datta S. Comparisons and validation of statistical clustering techniques for microarray gene expression data. Bioinformatics. 2003; 19:459–466. [PubMed: 12611800]
- Ghosh D, Chinnaiyan AM. Mixture modelling of gene expression data from microarray experiments. Bioinformatics. 2002; 18:275–286. [PubMed: 11847075]
- Hakamada K, Okamoto M, Hanai T. Novel technique for preprocessing high dimensional time-course data from DNA microarray: mathematical model-based clustering. Bioinformatics. 2006; 22:843– 848. [PubMed: 16434440]
- Hartigan JA. Partition models. Commun Stat Theory Methods. 1991; 19:2745-2756.
- Hartigan JA, Wong MA. Algorithm AS 136: a K-means clustering algorithm. Appl Stat. 1979; 28:100–108.
- James GM, Sugar CA. Clustering for sparsely sampled functional data. J Am Stat Assoc. 2003; 98:397–408.
- Jerrum, M.; Sinclair, A. Approximation algorithms for NP-hard problems. Boston: PWS Publishing; 1996. The Markov Chain Monte Carlo method: an approach to approximate counting and integration.
- Johnson, RA.; Wichern, DW. Applied multivariate statistical analysis. 5th edn.. Prentice Hall, Upper Saddle River: 2002.
- Leng X, Muller H. Classification using functional data analysis for temporal gene expression data. Bioinformatics. 2006; 22:68–76. [PubMed: 16257986]
- Luan Y, Li H. Clustering of time-course gene expression data using a mixed-effects model with Bsplines. Bioinformatics. 2003; 19:474–482. [PubMed: 12611802]
- Lukashin AV, Fuchs R. Analysis of temporal gene expression profiles: clustering by simulated annealing and determining the optimal number of clusters. Bioinformatics. 2001; 17:405–414. [PubMed: 11331234]
- Ma P, Castillo-Davis CI, Zhong W, Liu JS. A data-driven clustering method for time course gene expression data. Nucleic Acids Res. 2006; 34:1261–1269. [PubMed: 16510852]
- McLachlan, GJ.; Baford, KE. Mixture models: inference and applications to clustering. New York: Marcel Dekker, Inc; 1988.
- McLachlan GJ, Bean RW, Peel D. A mixture model-based approach to the clustering of microarray expression data. Bioinformatics. 2002; 18:413–422. [PubMed: 11934740]
- Ng SK, McLachlan GJ, Wang K, Jones LB, Ng SW. A mixture model with random-effects components for clustering correlated gene-expression profiles. Bioinformatics. 2006; 22:1745– 1752. [PubMed: 16675467]
- Ouyang M, Welsh WJ, Georgopoulos P. Gaussian mixture clustering and imputation of microarray data. Bioinformatics. 2004; 20:917–923. [PubMed: 14751970]
- Park T, Yi S, Lee S, Lee SY, Yoo D, Ahn J, Lee Y. Statistical tests for identifying differentially expressed genes in time-course microarray experiments. Bioinformatics. 2003; 19:694–703. [PubMed: 12691981]
- Peddada SD, Lobenhofer EK, Li L, Afshari CA, Weinberg CR, Umbach DM. Gene selection and clustering for time-course and dose response microarray experiments using order-restricted inference. Bioinformatics. 2003; 19:834–841. [PubMed: 12724293]
- Pitman J. Some probabilistic aspects of set partitions. Am Math Mon. 1997; 104:201-209.
- Ruppert, D.; Wand, MP.; Caroll, RJ. Semiparametric regression. New York: Cambridge University Press; 2003.
- Schliep A, Schonhuth A, Steinhoff C. Using hidden Markov models to analyze gene expression time course data. Bioinformatics. 2003; 19:i255–i263. [PubMed: 12855468]
- Thalamuthu A, Mukhopadhyay I, Zheng X, Tseng GC. Evaluation and comparison of gene clustering methods in microarray analysis. Bioinformatics. 2006; 22:2405–2412. [PubMed: 16882653]

Tibshirani R. Regression shrinkage and selection via the lasso. J R Stat Soc B. 1996; 58:267-288.

Tseng GC, Wong WH. Tight clustering: a resampling-based approach for identifying stable and tight patterns in data. Biometrics. 2005; 61:10–16. [PubMed: 15737073]



**Fig. 1.** Average gene-expressional profiles

Joo et al.







#### Fig. 3.

Convergence of chains with different initial partitions. **a**, **b** demonstrates the convergence with the number of clusters and the objective function  $Obj(\omega)$ . Simulation history in the last  $2 \times 10^6$  iterations is magnified in an insert of each figure



#### Fig. 4.

Comparisons of spline fits of a randomly selected gene profile when the real time scale or the relabeled equal interval time scale is used: (a) gene profile is fitted on the real time scale (days), (b) plot  $\mathbf{a}$  is redrawn with the relabeled time scale, (c) the profile is fitted on the equal interval scale, and (d) plot  $\mathbf{c}$  is redrawn with the relabeled time scale

Joo et al.







The effect of  $log(\varrho)$  in the Crowley prior on the number of stable genes. A quadratic regression fit is provided as a *reference line* 

Joo et al.



#### Fig. 6.

Clustered gene profiles of closely related genes. The number of genes in each cluster is described inside parentheses. The light-colored thick lines are from the BLUP of the mixed model (6)

# Table 1

Convergence of algorithm with different  $\log(\varrho)$ 's and stable genes in the last  $10^6$  iterations

log( <b>g</b> )	0		3		S		æ	
Run No.	1st	2nd	1st	2nd	1st	2nd	1st	2nd
$Obj(\tilde{\omega}^*)$	5,280	5,273	5,338	5,344	5,411	5,404	5,561	5,577
$c(\tilde{\omega}^*)$	27	30	42	44	64	56	135	121
No. of stable genes	58	75	40	52	43	41	50	51
log(ϱ)	10		13		15		20	
Run No.	1st	2nd	1st	2nd	1st	2nd	1st	2nd
Obj(@*)	5,849	5,846	6,940	6,974	8,351	8,342	11,631	11,630
$c(\tilde{\omega}^*)$	224	225	500	510	636	635	645	645
No. of stable genes	36	38	52	43	96	83	162	178

#### Table 2

Average number of iterations before the optimal partition is found  $(Avg(R^*))$ 

Simulated data set	<i>n</i> = 20	<i>n</i> = 50	<i>n</i> = 100
1	6,220	4,075,804	$>5  imes 10^7$
2	1,637	7,927,721	$>5 imes10^7$
3	663	26,765,931	$>5  imes 10^7$
4	4,440	12,796,622	$>5 imes10^7$
5	934	$>5 imes10^7$	$>5 imes10^7$

#### Table 3

Comparisons of tight and plain Bayes clustering algorithm based on ER: the standard deviation of ER is denoted inside parentheses

σΑ/σΒ	0.5/0.5	0.7/0.7	1.0/1.0	0.7/3.0
Tight Clustering	6.0 (0.8)	13.7 (1.3)	26.8 (0.9)	8.8 (0.8)
$\log(\varrho) = -50$	10.3 (0.3)	25.9 (0.9)	40.7 (0.6)	22.4 (0.9)
$\log(\varrho) = -30$	11.1 (0.4)	26.4 (0.9)	40.6 (0.6)	22.7 (0.9)
$\log(\varrho) = -20$	9.8 (0.3)	26.1 (0.9)	40.4 (0.6)	22.0 (0.8)
$\log(\varrho) = -15$	9.8 (0.3)	25.3 (0.9)	40.5 (0.6)	21.5 (0.8)
$\log(\varrho) = -10$	8.8 (0.2)	21.7 (0.8)	39.0 (0.6)	17.4 (0.7)
$\log(\varrho) = -5$	9.8 (0.2)	17.6 (0.5)	29.8 (0.7)	14.9 (0.4)
$\log(\varrho) = 0$	14.6 (0.1)	15.2 (0.1)	16.5 (0.2)	15.1 (0.1)
$\log(\varrho) = 5$	15.3 (0.0)	15.3 (0.0)	15.3 (0.0)	15.3 (0.0)