

Cluster Analysis, Model Selection, and Prior Distributions on Models

George Casella*
University of Florida

Elías Moreno†
University of Granada

F. Javier Girón ‡
University of Málaga

July 30, 2011

Abstract

Clustering is an important and challenging statistical problem for which there is an extensive literature. Modeling approaches include mixture models and product partition models. Here we develop a product partition model and a model selection procedure based on Bayes factors from intrinsic priors. We also find that the choice of the prior on model space is of utmost importance, almost overshadowing the other parts of the clustering problem, and we examine the behavior of posterior odds based on different model space priors. We find, somewhat surprisingly, that procedures based on the often-used uniform prior (in which all models are given the same prior probability) lead to inconsistent model selection procedures. We examine other priors, and find that a new prior, the *hierarchical uniform prior* leads to consistent model selection procedures and has other desirable properties. Lastly, we examine our procedures, and competitors, on a range of examples.

Keywords and phrases: Hierarchical models, Bayes models, Product partition models, Intrinsic priors, Mixtures, Consistency, Objective priors, Stochastic search.

*Distinguished Professor, Department of Statistics, University of Florida, Gainesville, FL 32611. Supported by National Science Foundation Grants DMS-0631632 and SES-0631588. Email: casella@stat.ufl.edu.

†Professor, Department of Statistics, University of Granada, 18071, Granada, Spain. Supported by Ministerio de Ciencia y Tecnología, Grant SEJ-65200 and Junta de Andalucía Grant SEJ-02814. Email: emoreno@ugr.es.

‡Professor, Department of Statistics, University of Málaga. Email: fj.giron@uma.es.

1 Introduction

Suppose that Y is an observable random variable with sampling distribution in a class of parametric densities $\mathfrak{F} = \{f(y|\theta), \theta \in \Theta\}$, where Θ is in the space \mathbb{R}^k , $k \geq 1$, and we observe a sample of n independent data $\mathbf{y} = (y_1, y_2, \dots, y_n)$ coming from the densities in the class \mathfrak{F} . An interesting problem for its wide range of applications, and theoretical challenges, arises when we look at the sample as being split into clusters, where all of the observations in a cluster come from the same sample density $f(y|\theta)$, and the parameter θ of the density changes across clusters. The clustering problem consists of making inference on the number of clusters in the sample and the location of the sample observations into these clusters.

Assuming that a family of sampling models \mathfrak{F} has been chosen, two difficulties appear in the clustering problem. A first one is the assessment of the prior distribution for the discrete parameters, the number of clusters and the partitions of the sample into these clusters. There is also the assessment of the prior for the (usually) continuous parameters θ of the densities of the partitions. Because of the typical lack of substantive prior information on these parameters, we often do not have enough information for a precise formulation of the priors. This leads us to propose the use of objective priors.

A second difficulty is that of computing the very many posterior model probabilities of the sampling models involved, even when the sample size is moderate. For instance, for a sample size as small as 20, the number of possible models (partitions) is as large as 51724158235372, which makes it infeasible to compute all of the posterior model probabilities, although fortunately only a small number of models will have nonnegligible posterior probabilities. Therefore, there is a need to develop an efficient stochastic search algorithm for computing posterior probabilities for those models having nonnegligible posterior probabilities.

1.1 Background

The literature on clustering is enormous. There are many ways to approach the problem, and here we take a model-based approach that can be divided into approaches based on *mixture models* (see, for instance, Fraley and Raftery 2002, 2007 and references therein), and approaches based on *product partition models* (PPMs). The latter is a model-based approach introduced by Hartigan (1990) and Barry and Hartigan (1992) (see Booth *et al.* 2008 and references therein). In this paper the sampling models for clustering are also product partition models constructed with normal regression models, and objective priors for models and for model parameters will be introduced.

In recent clustering approaches, Dirichlet distributions and PPMs are used as a prior for the clus-

ters in the hierarchical Bayesian framework. Quintana and Iglesias (2003) propose an algorithm for the explicit construction of clusters based on PPM-type priors for partitions of experimental units. They noted that the proposed model is quite flexible because PPMs can be used to express a wide variety of prior distributions on partitions. Lau and Green (2007) propose a general formulation for Bayesian model-based clustering that is optimal with respect to a specified loss function. They compare the new method to some recently discussed methods involving stochastic search or hierarchical clustering under the Dirichlet process by maximizing the posterior probability. In a somewhat different approach, Booth *et al.* (2008) propose a stochastic search algorithm driven by a mixture of two Metropolis-Hastings algorithms, one for small scale changes to individual objects and another for large scale moves involving entire clusters.

Here, we find that the choice of the prior on model space is of utmost importance, almost overshadowing the other parts of the clustering problem. We examine a number of priors, including the typical default uniform prior (where every model has the same prior probability) and find that such a prior leads to inconsistent procedures. We then develop a new prior on model space, the *hierarchical uniform prior*, that does not suffer from such disadvantages. For the continuous parameters we use intrinsic priors, and combine all pieces to evaluate a model using its posterior odds.

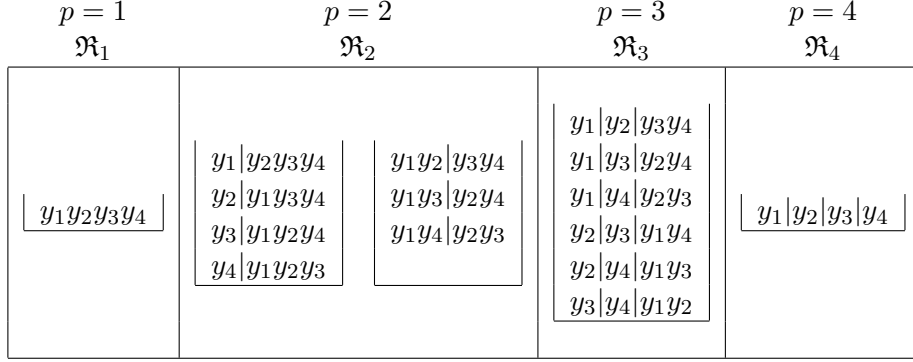
1.2 Summary

The rest of the paper is organized as follows. In Section 2 we describe the structure of the clustering problem. Some background on different approaches to the clustering problem, and a description of the approach taken here, is given in Section 3. In Section 4 we develop priors for the discrete parameters, the partitions and the number of clusters, and Section 5 develops intrinsic priors for the continuous parameters, illustrated with the case of linear models. The model priors are evaluated in Section 6, where we obtain our consistency and inconsistency results. We describe our search algorithm in Section 7, a variation of the biased random walk of Booth *et al.* (2008), and in Section 8 we illustrate the performance of our procedure on both real and simulated data. Finally, Section 9 contains a concluding discussion, and there is a technical appendix with proofs of theorems.

2 Models, Partitions and Configurations

We begin by formally describing the structure of the clustering problem, defining the terms that we will be using in the subsequent analyses. We start with a sample $\mathbf{y} = (y_1, y_2, \dots, y_n)$. For given p , we define a *partition* of the sample into p clusters, denoted by $\mathbf{r}_p = (r_1^{(p)}, \dots, r_n^{(p)})$, where $r_i^{(p)}$, $i = 1, \dots, n$, is an integer between 1 and p denoting the cluster to which y_i is assigned.

Figure 1: THE STRUCTURE OF THE CLUSTERING PROBLEM FOR $n = 4$. THERE ARE 15 POSSIBLE PARTITIONS (MODELS), THE BELL NUMBER FOR $n = 4$. IN EACH OF THE CLUSTER CLASSES, $p = 1, 2, 3, 4$ THERE ARE 1, 7, 6, 1 PARTITIONS, THE STIRLING NUMBERS OF THE SECOND KIND. WITHIN THE CLUSTER CLASS FOR $p = 2$ THERE ARE TWO CONFIGURATION CLASSES, CORRESPONDING TO THE CONFIGURATIONS $y|yyy$ AND $yy|yy$. THE NUMBER OF PARTITIONS IN EACH CONFIGURATION CLASS IS GIVEN IN (2), AND THE NUMBER OF CONFIGURATION CLASSES IN EACH CLUSTER CLASS IS $b(n, p)$ OF (3).



As we will see in detail below, the observations within the same configuration class have the same sampling distribution, while those in different classes have different sampling distributions. Thus, under this structure, a given partition \mathbf{r}_p is equivalent to a *model* $M_{\mathbf{r}_p}$. We will provide details in Section 3.2.

We next look at an example that, perhaps, will clarify the definitions, and we will refer to Figure 1 for illustration. Let \mathfrak{R}_p represent the set of partitions of the sample into p clusters, which we call the *cluster class*. The number of partitions in \mathfrak{R}_p is given by the *Stirling number of second kind*, $S(n, p)$. The four cluster classes of Figure 1 each have a Stirling number of partitions. The total of partitions of the sample $\mathfrak{R} = \cup_p \mathfrak{R}_p$ is given by the *Bell number*, $\mathcal{B} = \sum_{p=1}^n S(n, p)$, which in the case of Figure 1 is 15.

Suppose that the sample is split into p clusters, and let n_i be the number of components of the sample located in the i th cluster for $i = 1, 2, \dots, p$. The number of partitions of the sample of size n into p clusters $S(n, p)$ can be written as

$$S(n, p) = \sum_{\substack{n_1 + \dots + n_p = n \\ 1 \leq n_1 \leq \dots \leq n_p}} \binom{n}{n_1 \dots n_p} \frac{1}{R(n_1, \dots, n_p)}, \quad (1)$$

where $\binom{n}{n_1 \dots n_p}$ is the multinomial coefficient and $R(n_1, \dots, n_p) = \prod_{i=1}^n [\sum_{j=1}^p I(n_j = i)]!$ corrects the count by considering the redundant strings corresponding to the vector (n_1, \dots, n_p) . For instance, for the vector (n_1, \dots, n_p) such that $n_1 = \dots = n_{p-4} < n_{p-3} = n_{p-2} < n_{p-1} = n_p$, we have that

$$R(n_1, \dots, n_p) = (p-4)!2!2!$$

Denoting the set of partitions for a fixed vector (n_1, \dots, n_p) by $\mathfrak{R}_{p;n_1, \dots, n_p}$, we can express the class of partitions with p clusters \mathfrak{R}_p as

$$\mathfrak{R}_p = \bigcup_{\substack{n_1 + \dots + n_p = n \\ 1 \leq n_1 \leq \dots \leq n_p}} \mathfrak{R}_{p;n_1, \dots, n_p},$$

where $\mathfrak{R}_{p;n_1, \dots, n_p}$ denotes a *configuration class*, that is, the class of all partitions of the sample that have the same configuration (n_1, \dots, n_p) . The number of partitions in a *configuration class* is given by the corresponding term in (1), that is

$$\text{Number of partitions in } \mathfrak{R}_{p;n_1, \dots, n_p} = \binom{n}{n_1 \cdots n_p} \frac{1}{R(n_1, \dots, n_p)}. \quad (2)$$

As we see in Figure 1, a cluster class \mathfrak{R}_p can have more than one configuration class; for $p = 2$ there are two configuration classes. In general, the number of configuration classes within \mathfrak{R}_p is the number of way the integer n can be partitioned into p parts, which we denote by $b(n, p)$. This number does not seem to have a closed form expression as a function of p and n . However, it can be shown (Jara 2007, personal communication) that $b(n, p)$ satisfies the recursive equation

$$b(n, p) = b(n-1, p-1) + b(n-p, p), \quad 1 \leq p \leq n, \quad (3)$$

with $b(n, 1) = 1$, and $b(n, n) = 1$. The number $b(n, p)$ can be large, even for moderate values of n and p , for instance $b(80, 35) = 89037$. However, it is much smaller than $S(80, 35)$, which has 82 digits.

Now that the structure of the cluster problem is clear, we next want to look at ways in which we can attach sampling models to this structure, and the ways in which prior distributions on the models can be assigned. For example, because the configuration class $\mathfrak{R}_{p;n_1, \dots, n_p}$ contains partitions with the same structure, n_i , we model the data in each configuration class as coming from the same unknown density $\prod_{j=1}^{n_i} f(y_j | \theta_i)$, for $i = 1, \dots, p$. What is, perhaps, less clear is how to assign prior probabilities to the models, and we will see those details in Section 4.

3 Cluster Models

Currently, the most popular approach to cluster analysis is based on mixture models. We briefly describe that approach, and explain why we do not believe it is the best approach. We then describe the product partition model which, we believe, not only respects the structure of the clustering problem, but also can provide meaningful answers to the clustering problem.

3.1 Clustering with Mixture Models

Clustering has been one of the main traditional applications of the mixture model. Under a mixture model the observations y_1, \dots, y_n are assumed to be independently drawn from a model of the form

$$f(y|p, \theta, \omega) = \sum_{j=1}^p \omega_j f(y|\theta_j), \quad \omega_i \geq 0, \quad \sum_i \omega_i = 1,$$

where $p, (\theta_1, \dots, \theta_p)$, and $(\omega_1, \dots, \omega_p)$ are unknown parameters. Not only does this formulation completely ignore the structure outlined in Section 2, it has incurred criticism for a number of other reasons. For example, Stephens (2000), among others, note that there is a “label-switching problem” in that any permutation of the components of the parameter vector $((\omega_1, \theta_1), \dots, (\omega_p, \theta_p))$ has no effect on the density $f(y|p, \theta, \omega)$, and this implies that the individual components (ω_i, θ_i) or ω_i are not identifiable. It is also the case that, due to the symmetry of the model, a new observation y_{n+1} is classified *a posteriori* in any component of the mixture with probability $1/p$, whatever the sample, and thus this model is useless for classifying new observations into groups.

We add that, even if the parameters (p, θ, ω) in $f(y|p, \theta, \omega)$ were all known, we are not able to compute the posterior probability of a given partition of the sample into p clusters by using only this completely specified model. To be able to do so we need to add latent variables to impute the original model from which the mixture model comes.

3.2 Product Partition Models

For the structure of the cluster problem, a product partition model seems to be a natural way to introduce the sampling model and the prior. Given a partition $\mathbf{r}_p = (r_1^{(p)}, \dots, r_n^{(p)})$, we introduce a vector $\theta_p = (\theta_{r_1^{(p)}}, \dots, \theta_{r_n^{(p)}})$, an unknown parameter vector of dimension n whose components indicate sampling densities in the class \mathfrak{F} . The sampling density of the data \mathbf{y} conditional on a given partition \mathbf{r}_p is

$$f(\mathbf{y}|p, \mathbf{r}_p, \theta_p) = \prod_{i=1}^n f(y_i|\theta_{r_i^{(p)}}) = \prod_{j=1}^p \prod_{i:r_i=j} f(y_i|\theta_j). \quad (4)$$

Thus, the partition of the sample \mathbf{r}_p has associated with it p different values of the parameter θ_p , and the likelihood of the parameters p, \mathbf{r}_p and θ_p is given by (4). We will suppress the superscript (p) in $r_i^{(p)}$ if there is no confusion, and the likelihood will be simply written as $\prod_{i=1}^n f(y_i|\theta_{r_i})$.

The partition $\mathbf{r}_1 = (1, 1, \dots, 1)$ corresponds to the case where there is only one cluster in the sample. Its corresponding likelihood function is given by $f(\mathbf{y}|1, \mathbf{r}_1, \theta) = \prod_{i=1}^n f(y_i|\theta)$, $\theta \in \Theta$.

Since each partition is associated with a sampling density, and recalling that each partition corresponds to a model, it follows that the number of sampling models involved in a clustering

problem of a sample of size n is exactly the Bell number. The class of these models \mathfrak{M} can be seen as the union $\mathfrak{M} = \cup_{1 \leq p \leq n} \mathfrak{M}_p$, where \mathfrak{M}_p is the class of models for p clusters. To complete the specification of the sampling models, we need a prior distribution π for the parameters $p, \mathbf{r}_p, \theta_p$ so the generic Bayesian model is given as

$$M_{\mathbf{r}_p} : \{f(\mathbf{y}|p, \mathbf{r}_p, \theta_p), \pi(p, \mathbf{r}_p, \theta_p|n)\}.$$

A natural decomposition of this prior distribution is $\pi(p, \mathbf{r}_p, \theta_p|n) = \pi(\theta_p|p, \mathbf{r}_p, n) \pi(p, \mathbf{r}_p|n)$, and in the class of models for p clusters, \mathfrak{M}_p , the posterior probability of model $M_{\mathbf{r}_p}$, for $\mathbf{r}_p \in \mathfrak{R}_p$, is

$$\pi(\mathbf{r}_p|\mathbf{y}, p, n) = \frac{\pi(p, \mathbf{r}_p|n)m(\mathbf{y}|\mathbf{r}_p, n)}{\sum_{\mathbf{r}_p \in \mathfrak{R}_p} \pi(p, \mathbf{r}_p|n)m(\mathbf{y}|\mathbf{r}_p, n)}, \quad m(\mathbf{y}|\mathbf{r}_p, n) = \int f(\mathbf{y}|p, \mathbf{r}_p, \theta_p)\pi(\theta_p|p, \mathbf{r}_p, n)d\theta_p,$$

where $m(\mathbf{y}|\mathbf{r}_p, n)$ is the marginal of the data under model $M_{\mathbf{r}_p}$. Since $f(\mathbf{y}|1, \mathbf{r}_1, \theta)$ is nested in $f(\mathbf{y}|p, \mathbf{r}_p, \theta_p)$ it is convenient to write the posterior probability of $M_{\mathbf{r}_p}$ in the class \mathfrak{R}_p as

$$\pi(\mathbf{r}_p|\mathbf{y}, p, n) = \frac{\pi(p, \mathbf{r}_p|n)B_{\mathbf{r}_p, \mathbf{r}_1}(\mathbf{y})}{\sum_{\mathbf{r}_p \in \mathfrak{R}_p} \pi(p, \mathbf{r}_p|n)B_{\mathbf{r}_p, \mathbf{r}_1}(\mathbf{y})}, \quad \text{if } \mathbf{r}_p \in \mathfrak{R}_p,$$

where $B_{\mathbf{r}_p, \mathbf{r}_1}(\mathbf{y}) = \frac{m(\mathbf{y}|\mathbf{r}_p, n)}{m(\mathbf{y}|\mathbf{r}_1, n)}$ represents the Bayes factor for comparing model $M_{\mathbf{r}_p}$ against the model for only one cluster, $M_{\mathbf{r}_1}$.

In the class of all models \mathfrak{M} the posterior probability of model $M_{\mathbf{r}_p}$ is given by

$$\pi(\mathbf{r}_p|\mathbf{y}) = \frac{\pi(p, \mathbf{r}_p|n) B_{\mathbf{r}_p, \mathbf{r}_1}(\mathbf{y})}{\sum_{p=1}^n \sum_{\mathbf{r}_p \in \mathfrak{R}_p} \pi(p, \mathbf{r}_p|n) B_{\mathbf{r}_p, \mathbf{r}_1}(\mathbf{y})}, \quad \text{if } \mathbf{r}_p \in \mathfrak{R} \quad (5)$$

and the posterior probability of p by

$$\pi(p|\mathbf{y}) = \frac{\sum_{\mathbf{r}_p \in \mathfrak{R}_p} \pi(p, \mathbf{r}_p|n)B_{\mathbf{r}_p, \mathbf{r}_1}(\mathbf{y})}{\sum_{p=1}^n \sum_{\mathbf{r}_p \in \mathfrak{R}_p} \pi(p, \mathbf{r}_p|n)B_{\mathbf{r}_p, \mathbf{r}_1}(\mathbf{y})}, \quad 1 \leq p \leq n. \quad (6)$$

We note that $B_{\mathbf{r}_1, \mathbf{r}_1}(\mathbf{y}) = 1$.

Using only these Bayes factors we can also compare the non-nested models contained in the class \mathfrak{M} . Indeed, the Bayes factor for comparing models $M_{\mathbf{r}_p}$ and $M_{\mathbf{r}_q}$, for arbitrary p and q , is given by $B_{\mathbf{r}_p, \mathbf{r}_q}(\mathbf{y}) = \frac{B_{\mathbf{r}_p, \mathbf{r}_1}(\mathbf{y})}{B_{\mathbf{r}_q, \mathbf{r}_1}(\mathbf{y})}$, and, consequently, the posterior odds of models $M_{\mathbf{r}_p}$ and $M_{\mathbf{r}_q}$ can be written as the product of the Bayes factor and the prior odds, that is $\frac{\pi(\mathbf{r}_p|\mathbf{y})}{\pi(\mathbf{r}_q|\mathbf{y})} = B_{\mathbf{r}_p, \mathbf{r}_q}(\mathbf{y}) \frac{\pi(p, \mathbf{r}_p|n)}{\pi(q, \mathbf{r}_q|n)}$.

4 Prior Distributions on Models

Once a parametric family of sampling distributions \mathfrak{F} is specified, the first task for computing the posterior probability in (5) and (6) is that of assessing the prior distribution $\pi(\theta_p|p, \mathbf{r}_p, n)\pi(p, \mathbf{r}_p|n)$. In this section we concentrate on priors for the models, $\pi(p, \mathbf{r}_p|n)$, which we further factor as

$$\pi(p, \mathbf{r}_p|n) = \pi(\mathbf{r}_p|p, n)\pi(p|n).$$

In this representation we note that the factor $\pi(\mathbf{r}_p|p, n)$ is much more important than $\pi(p|n)$. Although both factors depend on n , the former is much more sensitive to n , as the size of the cluster classes grows exponentially with n .

4.1 Priors on the Partitions

The Uniform Prior The first prior one may use, in the absence of information about the models, is the *uniform prior*, which gives the same probability to every model, that is,

$$\pi^U(p, \mathbf{r}_p|n) = \frac{1}{\mathcal{B}_n}, \tag{7}$$

where \mathcal{B}_n is the Bell number. As we will see, this seemingly innocuous choice can have unforeseen consequences.

The Ewens-Pitman Prior The Dirichlet random process provides a marginal distribution for (p, \mathbf{r}_p) given by

$$\pi^{\text{EP}}(p, \mathbf{r}_p|\lambda, n) = \frac{\Gamma(\lambda)}{\Gamma(n + \lambda)} \lambda^p \prod_{j=1}^p \Gamma(n_j), \quad p = 1, \dots, n, \quad \mathbf{r}_p \in \mathfrak{R}_p, \tag{8}$$

where λ is an unknown positive hyperparameter which has to be assessed. This prior has been used extensively (Crowley 1997, Quintana and Iglesias 2003, Booth *et al.* 2008, Jensen and Liu 2008, McCullagh and Yang 2008). A detailed scheme to derive the prior (8) is presented in McCullagh and Yang (2008). They also note that the limit, as $p \rightarrow \infty$, is the Ewens process (Ewens 1972, Ishwaran and Zarepour 2002, Pitman 1996), also called the Chinese restaurant process (Aldous 1985, Pitman 1996).¹

¹We have referred to this distribution in a variety of ways, and each way has received criticism from some quarter. We believe that *Ewens-Pitman* allocates the correct degree of recognition.

The Jensen-Liu Prior An alternative to the prior π^{EP} is that given by Jensen and Liu (2008),

$$\pi^{JL}(p, \mathbf{r}_p | \lambda, n) \propto \lambda^{p-1} (\lambda + p) \prod_{i=1}^p (\lambda + i)^{-n_i}, \quad (9)$$

and according to their authors “favors equal allocations of observations, that is, the prior probability that a new observation is placed in any one of the existing clusters is uniform.”

The Hierarchical Uniform Prior (HUP) One of our goals is to develop an objective prior for models, and to do so we consider the structure of the cluster problem as described in Section 2. Thus, we first split the entire set of partitions by conditioning on the number of clusters, and also split these sets into subsets of partitions containing exchangeable configurations.

To carry out our prior specification, we start from the decomposition

$$\pi(p, \mathbf{r}_p | n) = \pi(\mathbf{r}_p | \mathfrak{R}_{p;n_1, \dots, n_p}, n) \pi(\mathfrak{R}_{p;n_1, \dots, n_p} | p, n) \pi(p | n).$$

We note that the difference among the partitions \mathbf{r}_p in $\mathfrak{R}_{p;n_1, \dots, n_p}$ is simply the different ways we assign n_i components of the sample of size n to the density labeled by θ_i , for $i = 1, \dots, p$. This implies that *a priori* these partitions should have the same probability, although the likelihood will vary across the partitions. Then, it seems reasonable to assign the uniform prior to \mathbf{r}_p in $\mathfrak{R}_{p;n_1, \dots, n_p}$ (this ensures exchangeability), that is,

$$\pi(\mathbf{r}_p | \mathfrak{R}_{p;n_1, \dots, n_p}, n) = \left(\binom{n}{n_1 \cdots n_p} \right)^{-1} R(n_1, \dots, n_p), \quad \mathbf{r}_p \in \mathfrak{R}_{p;n_1, \dots, n_p}.$$

Next, we make the reasonable assumption that the sets of partitions $\mathfrak{R}_{p;n_1, \dots, n_p}$ in \mathfrak{R}_p obtained as the vector (n_1, \dots, n_p) varies are *a priori* equally likely. Then, recalling (3), it follows that

$$\pi(\mathfrak{R}_{p;n_1, \dots, n_p} | p, n) = b(n, p)^{-1}, \quad (10)$$

and the hierarchical uniform specification is complete, with the exception of $\pi(p | n)$, which we will discuss in Section 4.2.

Comparisons Both the Ewens-Pitman and Jensen-Liu priors require specification of the hyperparameter λ . For Ewens-Pitman, it is known that the expected prior number of clusters is given by

$$E(p | \lambda, n) = \sum_{p=1}^n \frac{\lambda}{\lambda + p - 1}.$$

Table 1: PRIOR PROBABILITIES FOR EXCHANGEABLE PARTITION SETS IN \mathfrak{R}_3 FOR $n = 10$. THE DESIGNATION $\{1, 3, 6\}$ REFERS TO HAVING THREE CLUSTERS WITH 1, 3 AND 6 OBSERVATIONS; $\{2, 3, 5\}$ HAS 2, 3 AND 5 OBSERVATIONS. (BOTH EWENS-PITMAN AND JENSEN-LIU HAVE $\lambda = 1$.)

Configuration	Prior Probabilities			
	Ewens-Pitman	Hierarchical Uniform	Uniform	Jensen-Liu
$\{1, 3, 6\}$	0.17	0.14	0.09	0.15
$\{2, 3, 5\}$	0.10	0.14	0.27	0.08

Thus, large values of λ will lead to a larger number of clusters. However, the value of the prior probabilities is quite sensitive to the choice of λ , and selection of this hyperparameter is important. (In Quintana and Iglesias (2003) the distribution (8) was ruled out due to its sensitivity to λ ; see their Table 2, page 570.) Rather than presenting a lengthy sensitivity analysis of these priors, we just want to illustrate their performance with a small numerical example; further comparisons are in Section 6. Table 1 shows prior probability specifications for two partitions that are very close, only differing by moving one observation. But we see that only the hierarchical uniform prior gives these partitions the same prior probability.

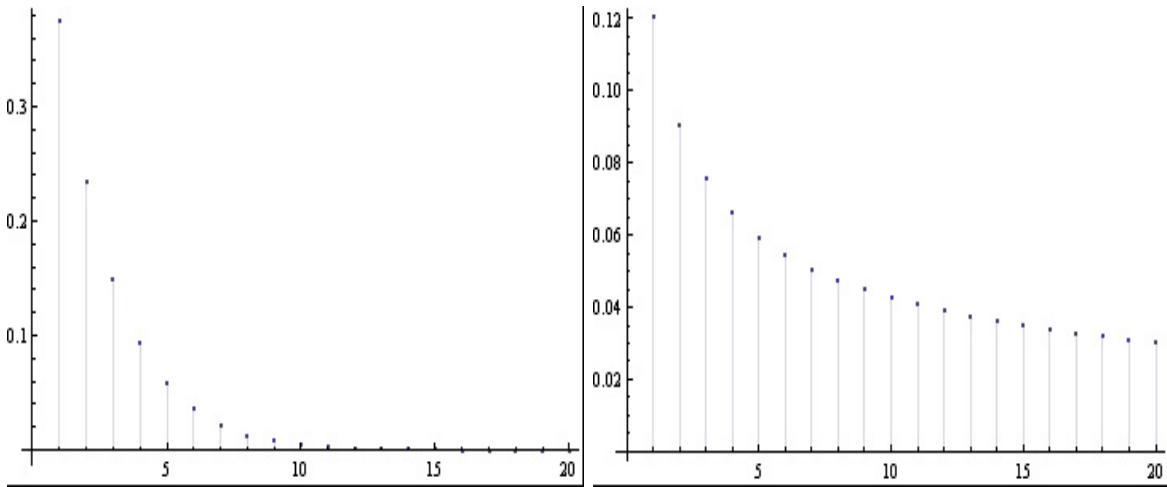
A priori it does not seem that, for $n = 10$, we would have any reason to assign different probabilities to the configurations $\{1, 3, 6\}$ and $\{2, 3, 5\}$, but this is what is done by the other priors. The fact that there are 840 partitions corresponding to the configuration $\{1, 3, 6\}$, and 2525 partitions corresponding to $\{2, 3, 5\}$ not only explains the numbers for the uniform prior in Table 1, but also illustrates its shortcomings. Moreover, the preferences reverse, with Ewens-Pitman and Jensen-Liu favoring $\{1, 3, 6\}$, and the uniform prior favoring $\{2, 3, 5\}$.

4.2 Prior Distributions for the Number of Clusters

To complete the specification of the prior $\pi(p, \mathbf{r}_p | n)$ we need to construct $\pi(p | n)$. For doing that we observe that when analyzing a cluster problem of a sample of size n we desire a relatively small number of clusters in the sample, and therefore the extreme case of having n clusters should be given *a priori* a very small probability. Extending this argument for any n , it seems reasonable that the prior distribution of the number of clusters be, in a smooth way, a decreasing function of the number of clusters p .

A candidate for $\pi(p | n)$ can be obtained by assuming that the number of clusters p follows a truncated Poisson distribution $\mathcal{P}(p | \lambda)$, where λ is an unknown parameter. Assuming the default improper Jeffreys distribution for λ , $\pi^J(\lambda) \propto \lambda^{-1/2}$, the marginal distribution for p is given by $\int \mathcal{P}(p | \lambda) \pi^J(\lambda) d\lambda$ which we now truncate to the set $\{1, \dots, n\}$. This prior is rather close to the

Figure 2: PRIOR PROBABILITIES FROM THE POISSON-INTRINSIC PRIOR OF (11) (LEFT PANEL), AND THE POISSON-JEFFREYS PRIOR (RIGHT PANEL), $n = 20$.



uniform prior in the sense that it has a very flat tail, and consequently for large n it will dilute the prior probabilities in $\{1, \dots, n\}$.

A way to derive a prior for p with thinner tails than the above one is obtained by replacing the Jeffreys prior $\pi^J(\lambda)$ with the intrinsic prior $\pi^I(\lambda|\lambda_0 = 1)$ constructed by testing the Poisson null hypothesis $H_0 : \lambda = \lambda_0$ versus $H_1 : \lambda \in R^+$. This prior, given by (Moreno *et al.* 2005), is

$$\pi^I(\lambda|\lambda_0 = 1) = \frac{\lambda^{-1/2}}{\Gamma(1/2)} e^{-(\lambda+1)} {}_0F_1(1/2, \lambda),$$

where ${}_0F_1(\frac{1}{2}, \lambda)$ denotes the confluent hypergeometric function. The reason for taking $\lambda_0 = 1$ is that the one cluster model is the reference model throughout the analysis. The resulting marginal intrinsic distribution for p is

$$\pi^I(p|n) = \frac{m^I(p)}{\sum_{p=1}^n m^I(p)}, \quad p = 1, \dots, n, \quad m^I(p) = \int_0^\infty \frac{\lambda^p e^{-\lambda}}{p!} \pi^I(\lambda|\lambda_0 = 1) d\lambda. \quad (11)$$

Figure 2 shows a plot of the Poisson-Intrinsic prior $m^I(p|n = 20)$ along with the Poisson-Jeffreys prior. There we see that the drop in prior probabilities, as we move from one cluster, is very flat for the Poisson-Jeffreys prior, but steeper for the Poisson-Intrinsic prior.

Then, using (11) together with the development above, we obtain the prior distribution for (p, \mathbf{r}_p) ,

for $\mathbf{r}_p \in \mathfrak{R}_{p;n_1, \dots, n_p}$ and $p = 1, \dots, n$,

$$\pi(p, \mathbf{r}_p | n) = \binom{n}{n_1 \cdots n_p}^{-1} R(n_1, \dots, n_p) b(n, p)^{-1} \frac{m^I(p)}{\sum_{p=1}^n m^I(p)}.$$

5 Intrinsic Priors for the Continuous Parameters θ_p

As we noted in Section 4, for computing the Bayes factors $B_{\mathbf{r}_p, \mathbf{r}_1}(\mathbf{y})$, $\mathbf{r}_p \in \mathfrak{R}$, prior distributions for the continuous parameters θ and θ_p are needed. The usual objective choices are the reference priors $\pi^N(\theta)$ and $\pi^N(\theta_p)$ associated with the sampling models $f(\mathbf{y}|1, \mathbf{r}_1, \theta)$ and $f(\mathbf{y}|p, \mathbf{r}_p, \theta_p)$, respectively (Berger *et al.* 2009). However, these priors are typically improper, and while this is not an inconvenience for estimating θ and θ_p , it is a serious problem for model comparison, as they leave the Bayes factor $B_{\mathbf{r}_p, \mathbf{r}_1}(\mathbf{y})$ defined only up to an arbitrary multiplicative constant.

Fortunately, the sampling model $f(\mathbf{y}|1, \mathbf{r}_1, \theta)$ is nested in $f(\mathbf{y}|p, \mathbf{r}_p, \theta_p)$ and then the reference priors can be converted into the so-called intrinsic priors (Berger and Pericchi 1996a, Moreno 1997, Moreno *et al.* 1998), for which not only is the Bayes factor well-defined, but also the intrinsic prior for θ_p concentrates probability mass around θ , a desirable condition for model comparison (known as the *Savage continuity condition*). Furthermore, the intrinsic priors have been proved to behave extremely well in a wide variety of problems (Berger and Pericchi 1996b; Berger and Mortera 1999; Kim and Sun 2000; Casella and Moreno 2005; 2009; Girón *et al.* 2006; Moreno 2005; Moreno *et al.* 2010; Casella *et al.* 2009, among others).

The intrinsic prior for the parameter θ_p , conditional on an arbitrary but fixed point θ , is given by

$$\pi^I(\theta_p | \theta) = \pi^N(\theta_p | \mathbf{r}_p) E_{\mathbf{y}(\ell_p) | \theta_p} \frac{f(\mathbf{y}|1, \mathbf{r}_1, \theta)}{\int f(\mathbf{y}(\ell_p) | p, \mathbf{r}_p, \theta_p) \pi^N(\theta_p) d\theta_p},$$

where the expectation is taken with respect to the sampling distribution $f(\mathbf{y}(\ell_p) | p, \mathbf{r}_p, \theta_p)$ with $\mathbf{y}(\ell_p)$ a vector of dimension $\ell_p = kp + 1$. Here ℓ_p denotes the minimal training sample size needed for estimating θ_p with the prior $\pi^N(\theta_p)$; that is, the minimal sample size for which

$$0 < \int f(y_1, \dots, y_k | \theta_p, \mathbf{r}_p) \pi^N(\theta_p) d\theta_p < \infty.$$

It can be easily checked that $\pi^I(\theta_p | \theta)$ is a probability distribution. The unconditional intrinsic prior for θ_p is given by

$$\pi^I(\theta_p) = \int \pi^I(\theta_p | \theta) \pi^N(\theta) d\theta,$$

and the pair $(\pi^N(\theta), \pi^I(\theta_p))$ is the intrinsic prior for comparing models $M_{\mathbf{r}_1}$ and $M_{\mathbf{r}_p}$. We note that they are improper priors whose moments typically do not exist, which seems to be a reasonable

property for objective priors, although they are well-calibrated priors in the sense that both depend on a unique arbitrary multiplicative constant, the arbitrary constant inherited from $\pi^N(\theta)$, which cancels out in the ratio. Therefore, the Bayes factor for intrinsic priors

$$B_{\mathbf{r}_p, \mathbf{r}_1}(\mathbf{y}) = \frac{\int f(\mathbf{y}|p, \mathbf{r}_p, \theta_p) \pi^I(\theta_p) d\theta_p}{\int f(\mathbf{y}|1, \mathbf{r}_1, \theta) \pi^N(\theta) d\theta}, \quad \mathbf{r}_p \in \mathfrak{R},$$

is free of arbitrary constants, needs neither subjective input nor actual (data-dependent) training samples, and is completely automatic.

5.1 The Case of Linear Models

We now consider the case where the class of parametric sample densities, \mathfrak{F} , is the class of linear models with k regressors. For example, suppose that the sample (y_1, \dots, y_n) follows a normal linear model

$$\mathbf{y} = \mathbf{X}\beta + \varepsilon, \quad \varepsilon \sim N_n(\varepsilon|0, \tau^2 \mathbf{I}_n),$$

where \mathbf{X} is an $n \times k$ design matrix of full rank, β is a vector of regression coefficients with dimension k , and τ^2 is the common variance of the error terms. This is the sampling model for one cluster in the sample and, in the notation of the preceding section, we have

$$f(\mathbf{y}|1, \mathbf{r}_1, \beta, \tau) = N_n(\mathbf{y}|\mathbf{X}\beta, \tau^2 \mathbf{I}_n).$$

The reference prior for the parameters of this model is $\pi^N(\beta, \tau) = c/\tau$, where c is an arbitrary positive constant.

Without loss of generality, suppose we split the sample into p clusters, where one cluster is formed with the first n_1 components of the sample, a second cluster is formed with the second n_2 components, and so on. These clusters correspond to the partition of the sample as $\mathbf{y}' = (\mathbf{y}'_1, \dots, \mathbf{y}'_p)$ and the corresponding partition of the design matrix as

$$\mathbf{X} = \begin{pmatrix} \mathbf{X}_1 \\ \dots \\ \mathbf{X}_p \end{pmatrix},$$

where \mathbf{X}_i has dimensions $n_i \times k$. Then, the sampling model for p clusters is given by

$$f(\mathbf{y}|p, \mathbf{r}_p, \beta_1, \dots, \beta_p, \sigma_p) = \prod_{i=1}^p N_n(\mathbf{y}_i | \mathbf{X}_i \beta_i, \sigma_p^2 \mathbf{I}_{n_i}),$$

assuming that σ_p^2 is the common variance of the p clusters model.

We note that the model for p clusters can be written as the linear model

$$\mathbf{y} = \mathbf{V}\gamma + \eta,$$

where \mathbf{V} is the following upper triangular $n \times kp$ design matrix

$$\mathbf{V} = \begin{pmatrix} \mathbf{X}_1 & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{X}_2 & \mathbf{X}_2 & \cdots & \mathbf{0} \\ \vdots & \vdots & \cdots & \vdots \\ \mathbf{X}_p & \mathbf{X}_p & \cdots & \mathbf{X}_p \end{pmatrix},$$

γ being the $k \times p$ vector $\gamma' = (\gamma'_1, \dots, \gamma'_p)$, where $\gamma_1 = \beta_1, \gamma_2 = \beta_2 - \beta_1, \dots, \gamma_p = \beta_p - \beta_{p-1}$ are k dimensional vectors, and the random vector η is distributed as $N_n(\eta | \mathbf{0}, \sigma_p^2 \mathbf{I}_n)$. Thus, the sampling model for p clusters is now given by

$$f(\mathbf{y}|p, \mathbf{r}_p, \gamma_1, \dots, \gamma_p, \sigma_p) = N_n(\mathbf{y} | \mathbf{V}\gamma, \sigma_p^2 \mathbf{I}_n).$$

Since $f(\mathbf{y}|1, \mathbf{r}_1, \beta, \tau)$ is nested in $f(\mathbf{y}|p, \mathbf{r}_p, \gamma_1, \dots, \gamma_p, \sigma_p)$, by simply making $\beta_1 = \beta_2 = \dots = \beta$ or, equivalently, $\gamma_1 = \beta, \gamma_2 = \dots = \gamma_p = 0$, direct application of the standard intrinsic methodology gives the intrinsic prior for the parameters $(\gamma_1, \dots, \gamma_p, \sigma_p)$ conditional on a fixed point (β, τ) , as

$$\pi^I(\gamma_1, \dots, \gamma_p, \sigma_p | \beta, \tau) = \frac{2}{\pi\tau(1 + \sigma_p^2/\tau^2)} N_{pk}(\gamma | (\mathbf{X}\beta, \mathbf{0}, \dots, \mathbf{0})', (\sigma_p^2 + \tau^2)\mathbf{W}^{-1}),$$

where $\mathbf{W}^{-1} = n/(pk + 1)(\mathbf{V}'\mathbf{V})^{-1}$. Note that the conditional intrinsic prior is centered at the null model (the one cluster model), and its covariance structure depends on the covariance matrix of the model of p clusters only. The unconditional intrinsic priors are given by the pair $(\pi^N(\beta, \tau), \pi^I(\gamma_1, \dots, \gamma_p, \sigma_p))$, where

$$\pi^I(\gamma_1, \dots, \gamma_p, \sigma_p) = \int \pi^I(\gamma_1, \dots, \gamma_p, \sigma_p | \beta, \tau) \pi^N(\beta, \tau) d\gamma d\tau. \quad (12)$$

5.2 Bayes Factors for Intrinsic Priors

The objective intrinsic Bayesian model for one cluster is

$$M_{\mathbf{r}_1} : \{N_n(\mathbf{y}|\mathbf{X}\beta, \tau^2\mathbf{I}_n), \pi^N(\beta, \tau)\},$$

and for p clusters

$$M_{\mathbf{r}_p} : \{N_n(\mathbf{y}|\mathbf{V}\gamma, \sigma_p^2\mathbf{I}_n), \pi^I(\gamma_1, \dots, \gamma_p, \sigma_p)\},$$

where $\pi^I(\gamma_1, \dots, \gamma_p, \sigma_p)$ is given in (12).

To compute the Bayes factor of model $M_{\mathbf{r}_p}$ versus model $M_{\mathbf{r}_1}$, we note that the residual sum of squares of a partition, with cluster sizes (n_1, \dots, n_p) , is equal to

$$RSS_{n_1, \dots, n_p} = \mathbf{y}'(\mathbf{I} - \mathbf{H}_\mathbf{V})\mathbf{y},$$

where $\mathbf{H}_\mathbf{V} = \mathbf{V}(\mathbf{V}'\mathbf{V})^{-1}\mathbf{V}'$. Simple but cumbersome algebra shows that, due to the upper diagonal structure of the matrix \mathbf{V} , the residual sum of squares for that partition, with cluster sizes (n_1, \dots, n_p) , can be written as

$$RSS_{n_1, \dots, n_p} = \sum_{i=1}^p RSS_{n_i},$$

where RSS_{n_i} is the residual sum of squares from the regression in the i^{th} cluster. Some lengthy calculations render a quite simple form for the Bayes factor for intrinsic priors. The following theorem summarizes the result.

Theorem 1 *The Bayes factor for the model $M_{\mathbf{r}_p}$ versus model $M_{\mathbf{r}_1}$ is given by*

$$B_{\mathbf{r}_p, \mathbf{r}_1}(\mathbf{y}) = \frac{2}{\pi} (pk + 1)^{(p-1)k/2} \int_0^{\pi/2} \frac{\sin^{(p-1)k} \varphi (n + (pk + 1) \sin^2 \varphi)^{(n-pk)/2}}{(n\mathcal{R}_{\mathbf{r}_p} + (pk + 1) \sin^2 \varphi)^{(n-k)/2}} d\varphi \quad (13)$$

where the statistic $\mathcal{R}_{\mathbf{r}_p}$ is

$$\mathcal{R}_{\mathbf{r}_p} = \frac{RSS_{n_1} + \dots + RSS_{n_p}}{RSS_n},$$

with $RSS_{n_i} = \mathbf{y}'_i(\mathbf{I} - \mathbf{H}_i)\mathbf{y}_i$, $i = 1, \dots, p$, $RSS_n = \mathbf{y}'(\mathbf{I} - \mathbf{H})\mathbf{y}$, and \mathbf{H}_i and \mathbf{H} the hat matrices associated with \mathbf{X}_i and \mathbf{X} , respectively.

Proof. The marginal under model \mathbf{r}_p is given by

$$m_{\mathbf{r}_p}(\mathbf{y}) = \int \left\{ \int \left\{ \int N_n(\mathbf{y}|\mathbf{X}\beta, \tau^2\mathbf{I}_n) \pi^I(\gamma_1, \dots, \gamma_p, \sigma_p | \beta, \tau) d\gamma \right\} \pi^N(\beta, \tau) d\beta \right\} d\sigma_p d\tau.$$

The integral with respect to γ and β can be done in closed form, and a change of variables from (σ_p, τ) to polar coordinates gives the above expression for the Bayes factor. ■

Substituting the Bayes factor for the intrinsic priors (13) in expressions (5) and (6), the intrinsic posterior probability of model $M_{\mathbf{r}_p}$ and the posterior probability of p clusters in the sample, are obtained. We want to note that with a large number of observations, the factor $\mathcal{R}_{\mathbf{r}_p}$ in (13) can get very close to zero, causing numerical problems in computation (the integral is returned as infinite). However, the transformation $t = \frac{pk+1}{\mathcal{R}_{\mathbf{r}_p}} \sin^2(\varphi)$ results in a representation that is numerically very stable, and allows for doing all computations on the log scale.

Lastly, there is one technicality to note. If a partition contains a value of n_i with $n_i < k$ then, of course, the regression cannot be fit in that cluster. We proceed by fitting the largest model feasible in that particular partition, with the limiting case being a cluster of size 1, to which we assign $\text{Var}(Y)$ as the residual sum of squares. (One might consider assigning a residual sum of squares of zero to such a cluster, but this unduly rewards clusters of size 1.)

6 The Effect of the Prior on the Limiting Bayes Procedure

It is well-known that for regular sampling models, the Bayesian model selection procedure is consistent when the dimension of the sampling model is fixed and comparisons are pairwise. In that case consistency of the Bayesian model selection procedure is inherited from the consistency of the Bayes factor, because of the model prior does not play any role in the consistency of the procedure. However, when the dimension of the model grows with the sample size, the model prior plays an important role for obtaining consistency, and this is the case in clustering. As we will see, the choice of the prior on model space is of major importance in determining the asymptotic behavior of a clustering procedure. Surprisingly, the actual choice of Bayes factor is of almost no consequence in determining consistency, as many Bayes factors have the same asymptotic representation.

In this section we look at the asymptotic behavior of Bayesian clustering procedures when using the four model priors of Section 4. Perhaps the most surprising result is that the uniform prior, which gives the same prior probability to every model, leads to an inconsistent procedure. Furthermore, this is the case when sampling from the simpler model, and the number of clusters is finite, a situation in which consistency is typically obtained for a Bayesian testing procedure. This behavior is explained by observing that, in clustering, the model prior depends on the sample size and, as the sample size tends to infinity, the speed of convergence of the prior to its limit compared with that of the Bayes factor is now crucial.

In what follows we will assume that the number of clusters p is bounded, that is, $p \leq T < \infty$.

We could be more general and put a growth rate on the number of clusters (Moreno *et al.* 2010) but, in practice, assuming that p is bounded is certainly a realistic constraint.

For large samples, approximations of Bayes factors for intrinsic priors depend on the dimensions of the competing models and a *pseudodistance* between them. If M_i and M_j are arbitrary general normal linear models, the pseudodistance from M_i to M_j is defined as

$$\delta_{ij} = \frac{1}{\sigma_i^2} \beta_i' \frac{\mathbf{X}_i' (\mathbf{I}_n - \mathbf{H}_j) \mathbf{X}_i}{n} \beta_i.$$

Note that this pseudodistance is defined for every pair of models, not only nested models, and it is not symmetric. Some useful properties of δ_{pi} are the following: (a) The distance from any model M_i to itself is always 0, (b) If M_i is nested in M_j , then $\delta_{ij} = 0$, and (c) if model M_i is nested in M_j , then $\delta_{ki} \geq \delta_{kj}$ for any model M_k .

We start with the following lemma, where we recall that the singular class $\mathfrak{A}_{1,n}$ contains only the one cluster model $M_{\mathbf{r}_1}$. (We use the notation $[M]$ to denote the model that generated the sample.)

Lemma 1 *Suppose that model M_i , of dimension i , is nested in model M_j , of dimension j , and M_t is the true model. Under the sampling model M_t , as $n \rightarrow \infty$, the Bayes factor can be approximated by*

$$B_{r_j r_i}(\mathbf{y}) \approx \exp \left\{ \left(\frac{i-j}{2} \right) \log \left(\frac{n}{j+1} \right) \right\} \left(\frac{1 + \delta_{ti}}{1 + \delta_{tj}} \right)^{n/2} [M_t],$$

where δ_{ti} and δ_{tj} are the pseudodistances from the true model to models M_i and M_j , respectively.

In particular, when sampling from model $M_{\mathbf{r}_1}$, for large n , the Bayes factor $B_{\mathbf{r}_p \mathbf{r}_1}(\mathbf{y})$ can be approximated by

$$B_{\mathbf{r}_p \mathbf{r}_1}(\mathbf{y}) \approx \left(\frac{pk+1}{n} \right)^{k(p-1)/2} [M_{\mathbf{r}_1}].$$

Proof. The first expression immediately follows from Lemma 3 in Girón *et al.* (2010), and the second expression follows from noting that $\delta_{r_j r_1}$ is zero for all models $M_{\mathbf{r}_j}$. ■

Lemma 1 shows that for large n , when sampling from $M_{\mathbf{r}_1}$, the Bayes factor $B_{\mathbf{r}_p \mathbf{r}_1}(\mathbf{y})$ is constant across partitions \mathbf{r}_p in the class \mathfrak{A}_p . Moreover, this asymptotic result is not limited to intrinsic Bayes factors; for example we know that BIC is asymptotically equivalent to the intrinsic Bayes factor. Moreover, Casella *et al.* (2009) show that the approximation holds for a wide class of priors.

In the following two subsections we analyze the limiting behavior of the four priors of Section 4.1, and then examine their effect on the consistency of the resulting Bayes procedures.

6.1 Limiting Behavior of the Model Priors

For a fixed value of p , we denote by \mathbf{Z}_n a p -dimensional random vector which takes values (n_1, \dots, n_p) on the set of integers satisfying the conditions $n_1 \leq \dots \leq n_p$ and $n_1 + \dots + n_p = 1$. We write

$$\Pr(\mathbf{Z}_n = (n_1, \dots, n_p) | \mathfrak{R}_p) = \Pr(M_{\left(\frac{n_1}{n}, \dots, \frac{n_p}{n}\right)}),$$

which represents the prior probability of any model arising from the priors of Section 4.1 in the set of models $\mathfrak{R}_{p;n_1, \dots, n_p}$. The study of the limiting behavior in this way cannot be done as the sample space where the random vector \mathbf{Z}_n takes values changes with the sample size n . If we consider instead the common space of the simplex \mathcal{S}_p , all models can be regarded as points of this simplex. Thus, we define the random vector $\mathbf{X}_n = \frac{1}{n} \mathbf{Z}_n$, and study the limiting behavior of \mathbf{X}_n for the different priors. Note that the sample space of \mathbf{X}_n is \mathcal{S}_p for all n , and $\Pr(\mathbf{X}_n = \left(\frac{n_1}{n}, \dots, \frac{n_p}{n}\right) | \mathfrak{R}_p) = \Pr(\mathbf{Z}_n = (n_1, \dots, n_p) | \mathfrak{R}_p)$.

Theorem 2 (a) *For fixed p and the hierarchical uniform prior, $\mathbf{X}_n = \frac{1}{n} \mathbf{Z}_n$ converges in distribution to a uniform prior on the simplex \mathcal{S}_p .*

(b) *For fixed p and λ and the Ewens-Pitman prior, $\mathbf{X}_n = \frac{1}{n} \mathbf{Z}_n$ converges in distribution to an improper prior distribution proportional to*

$$\frac{1}{x_1 \cdot x_2 \cdots (1 - x_1 - \cdots - x_{p-1})}, \quad (14)$$

restricted to the simplex \mathcal{S}_p . Note that this function does not depend on λ .

(c) *For fixed p and the uniform prior, $\mathbf{X}_n = \frac{1}{n} \mathbf{Z}_n$ converges in distribution, and in probability, to a degenerate distribution concentrated at the vertex $(\frac{1}{p}, \dots, \frac{1}{p})$ of the simplex \mathcal{S}_p .*

(d) *For fixed p and λ and the Jensen-Liu prior, $\mathbf{X}_n = \frac{1}{n} \mathbf{Z}_n$ converges in distribution, and in probability, to a degenerate distribution concentrated at the interior point $(\frac{\lambda_T^H}{p(\lambda+T)}, \dots, \frac{\lambda_T^H}{p(\lambda+1)})$ of the simplex \mathcal{S}_p , where λ_p^H denotes the harmonic mean of $\lambda + 1, \dots, \lambda + p$.*

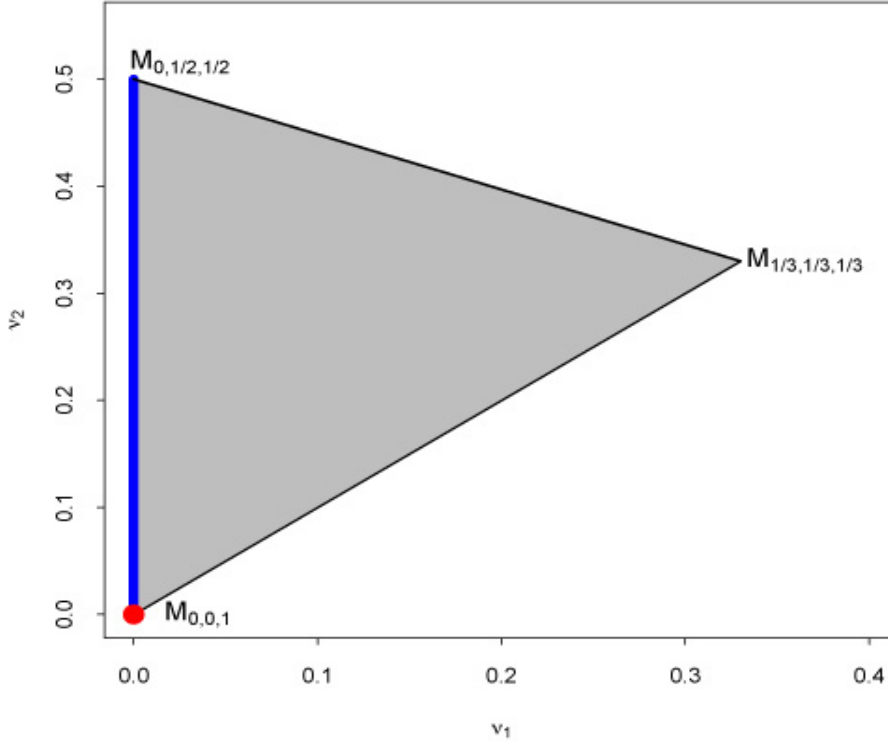
Proof. The proof is given in Appendix A.2 ■

Thus, of all four priors, only the HUP converges to a proper distribution. What is, perhaps, most distressing is the limiting behavior of the uniform and Jensen-Liu priors, which degenerate to a point and, thus, could have undo influence in the clustering algorithm.

6.2 Consistency of Bayes Procedures

To analyze consistency in clustering we realize that, as the sample size tends to infinity, the way we allocate the components of the sample in the clusters is not an issue since the sequence is not

Figure 3: THE SIMPLEX FOR $p = 3$ IS IN GREY. THE POINT IN RED REPRESENTS THE CLUSTER CLASS \mathfrak{R}_1 , THAT IS, THE ONE CLUSTER MODEL, WHILE THE BLUE LINE REPRESENTS \mathfrak{R}_2 . THE REST OF THE SIMPLEX REPRESENTS THE CLUSTER CLASS \mathfrak{R}_3 . NOTE THAT CLUSTER CLASS \mathfrak{R}_1 IS A VERTEX (EXTREME POINT) OF THE SIMPLEX, THE CLUSTER CLASS \mathfrak{R}_2 IS AN EDGE, AND THE MODEL WHICH ASSIGNS AN EQUAL PROPORTION TO THE THREE CLUSTER CONFIGURATION IS ALSO AN EXTREME POINT OF THE SIMPLEX.



observable, and hence we need only to consider the proportions of the sample in the clusters. Consequently, as n tends to infinity, the notion of consistency is now specific to the class of partitions having a limiting configuration $(n_1/n, \dots, n_p/n) \rightarrow (\nu_1, \dots, \nu_p)$, as $n \rightarrow \infty$. The interpretation of model M_{ν_1, \dots, ν_p} is that observations are assigned to clusters in the proportion ν_1, \dots, ν_p .

This also implies that, as $n \rightarrow \infty$, the model $M_{\mathbf{r}_p}$ is not distinguishable from $M_{\mathbf{r}'_p}$, assuming that both partitions \mathbf{r}_p and \mathbf{r}'_p belong to the same configuration class $\mathfrak{R}_{p;n_1, \dots, n_p}$. Consequently, as n tends to infinity, consistency will be examined in the probability space generated by the configuration classes $\{\mathfrak{R}_{p;n_1, \dots, n_p}, n_1 \leq \dots \leq n_p, n_1 + \dots + n_p = n, p = 1, 2, \dots\}$, and in the probability space generated by the cluster classes $\{\mathfrak{R}_p, p = 1, \dots, T\}$. We also note that configuration classes, and their limits, can be identified with points in the simplex $\mathcal{S}_p = \{\nu = (\nu_1, \dots, \nu_{p-1}); 0 \leq \nu_1 \leq \dots \leq \nu_{p-1}, \text{ and } \nu_1 + \dots + \nu_{p-1} < 1\}$, and cluster classes \mathfrak{R}_p with the simplex \mathcal{S}_p . This is illustrated in Figure 3, where we show the cluster classes and models for $p = 3$. The cluster class \mathfrak{R}_3 is the entire simplex, while \mathfrak{R}_2 and \mathfrak{R}_1 are an edge and a vertex, respectively. The model $M_{1/3, 1/3, 1/3}$, the extreme

right vertex of the simplex, is the limit of the uniform prior. Of course, this is extendable to higher dimensional simplexes.

We now use this representation of the Bayes factor to show that a clustering procedure that uses either the uniform prior or the Jensen-Liu prior is inconsistent.

Theorem 3 *Suppose that $p \leq T$ and we use either the uniform prior or the Jensen-Liu prior on the class of all partitions $\mathfrak{R} = \cup_{p=1}^T \mathfrak{R}_p$. Then, when sampling from $M_{\mathbf{r}_1}$ the Bayes procedure is inconsistent in both the cluster classes and the configuration classes. Moreover, in the probability space generated by the cluster classes $\{\mathfrak{R}_p, p = 1, \dots, T\}$, the asymptotic posterior distribution of \mathfrak{R}_p , $p = 1, \dots, T$, is*

$$\lim_{n \rightarrow \infty} [M_{\mathbf{r}_1}] \Pr(\mathfrak{R}_p | \mathbf{y}) = \begin{cases} 1, & \text{if } p = T, \\ 0, & \text{if } p \leq T - 1. \end{cases},$$

thus the largest model is chosen with probability one.

Proof. The proof is given in Appendix A.3 ■

The implications of this theorem are quite interesting, and help explain some of what we had observed in looking at examples (illustrated in Section 8). With priors like the uniform, the answers from the cluster algorithm tend to be partitions with a large number of clusters, and a small number of subjects per cluster.

The situation for the Ewens-Pitman prior and the hierarchical uniform prior is different.

Theorem 4 *Suppose that we use either the Ewens-Pitman prior or the hierarchical uniform prior on the class of all partitions $\mathfrak{R} = \cup_{p=1}^n \mathfrak{R}_p$. Then, when sampling from $M_{\mathbf{r}_1}$, the Bayesian procedure is consistent. That is, in the probability space generated by the cluster classes $\{\mathfrak{R}_p, p = 1, \dots, T\}$,*

$$\lim_{n \rightarrow \infty} [M_{\mathbf{r}_1}] \Pr(\mathfrak{R}_1 | \mathbf{y}) = 1,$$

so the correct model is chosen with probability 1.

Proof. The proof is given in Appendix A.4 ■

Thus, both of these priors exhibit good asymptotic behavior, but if we investigate further we see that, ultimately, the hierarchical uniform prior is the preferred choice. Consider the rate of convergence of the posterior probability of \mathfrak{R}_1 to one. For the Ewens-Pitman prior the convergence rate is $O\left(\frac{\log n}{n^{p-1}}\right)$, while for the HUP it is $O\left(\frac{1}{n^{p-1}}\right)$. This means that the convergence rate is faster when using the HUP than with the Ewens-Pitman prior. The difference is the presence of the $\log n$ term. We also note that the rate with the Ewens-Pitman prior also depends on the value of λ ; the larger this value the slower the consistency under the null.

7 Search Algorithm

Here we develop a hybrid search algorithm, using a Metropolis-Hastings (MH) algorithm that has stationary distribution proportional to the Bayes factor times the prior odds. The hybrid algorithm is a mixture of a random walk and a jump. We use a random walk component to be able to explore locally, and the jump allows escape from regions with small Bayes factors.

In setting up the algorithm there is one immediate computational problem, that of calculating the correct probabilities for the MH ratio. In the random walk piece we solve this problem by using the *biased random walk* of Booth *et al.* (2008). Suppose, for example, that at iteration t we have the partition $r_{p^{(t)}}^{(t)}$. We now generate a candidate partition $r_{p'}^{(t)}$ from a distribution G , and accept the move with probability

$$\min \left\{ \frac{\pi(p, \mathbf{r}_{p'}^{(t)} | n) B_{\mathbf{r}_{p'}, \mathbf{r}_1}}{\pi(p, \mathbf{r}_{p^{(t)}}^{(t)} | n) B_{\mathbf{r}_{p^{(t)}}, \mathbf{r}_1}} \frac{G(r_{p^{(t)}}^{(t)} | r_{p'}^{(t)})}{G(r_{p'}^{(t)} | r_{p^{(t)}}^{(t)})}, 1 \right\}$$

where $\pi(p, \mathbf{r}_{p'}^{(t)} | n)$ is the prior, and the Bayes factor $B_{\mathbf{r}_{p'}, \mathbf{r}_1}$ is given in (13). The computational problem arises in calculating the ratio of candidate probabilities, which could entail summing over an enormous number of partitions. However, the biased random walk has the property that $G(x|y) = G(y|x)$, and thus these terms cancel from the MH ratio. (Details and properties of the biased random walk are discussed in Booth *et al.* (2008), so here we will just give a brief description.)

Biased Random Walk

With the current iteration at $r_{p^{(t)}}^{(t)}$, we generate a candidate $r_{p'}^{(t)}$ as:

1. If $p = 1$, choose an observation at random from all n observations, and move the chosen observation to its own cluster. The new configuration is $r_{p'}^{(t)}$.
2. If $p > 1$, choose an observation at random from all n observations.
 - (a) If the object is a singleton cluster, move it to one of the $p-1$ other clusters with probability $1/(p-1)$.
 - (b) If the object is not a singleton cluster, move it to one of the $p-1$ other clusters, or to its own (new) cluster, each with probability $1/p$.

This is the biased random walk which, although similar to a nearest neighbor random walk, has the property that the probability of the move $r_{p^{(t)}}^{(t)} \rightarrow r_{p'}^{(t)}$ is the same as the probability of the move $r_{p'}^{(t)} \rightarrow r_{p^{(t)}}^{(t)}$, eliminating the necessity for calculating these probabilities in the Metropolis-Hastings algorithm.

Unfortunately, the mixing from the biased random walk is too slow for clustering large, or even medium, data sets. If we are in a good portion of the space then the random walk will explore that portion, but we also need to be able to escape from areas with small Bayes factors. To do so we have a second piece in the search algorithm, a jump based on sampling from the Ewens-Pitman distribution of (8). We can sample from this distribution using a number of algorithms (see, for example, Neal 2000), but we will draw our candidate using the algorithm of Kyung *et al.* (2010), which has been shown to mix better than many of its competitors. (Note that using the Ewens-Pitman distribution to drive a search has nothing to do with the choice of model space priors.)

Jumping with the Ewens-Pitman Distribution

We use the Ewens-Pitman distribution to generate a random jump because it is easy to calculate the Metropolis-Hastings correction. With the current iteration at $r_{p^{(t)}}^{(t)}$, we generate a candidate $r'_{p'}$ as follows: Start with $\mathbf{n}_p = (n_1, \dots, n_p)$ obtained from $r_{p^{(t)}}^{(t)}$, and draw \mathbf{q} from the Dirichlet distribution

$$\mathbf{q} \sim f(\mathbf{q}|\mathbf{n}_p) = \frac{\Gamma(2n)}{\prod_{j=1}^p \Gamma(n_j + 1)} \prod_{j=1}^p q_j^{n_j} \quad (15)$$

Given \mathbf{q} , we draw $\mathbf{n}'_{p'}$ from

$$\mathbf{n}'_{p'} \sim P(\mathbf{n}'_{p'}|\mathbf{q}) \sim \frac{\frac{\Gamma(n)}{\Gamma(n+\lambda)} \lambda^{p'} \prod_{j=1}^{p'} \Gamma(n'_j) \binom{n}{n_1 \dots n'_{p'}} \prod_{j=1}^{p'} q_j^{n'_j}}{\sum_{\mathbf{n}} \frac{\Gamma(n)}{\Gamma(n+\lambda)} \lambda^p \prod_{j=1}^p \Gamma(n_j) \binom{n}{n_1 \dots n_p} \prod_{j=1}^p q_j^{n_j}} \quad (16)$$

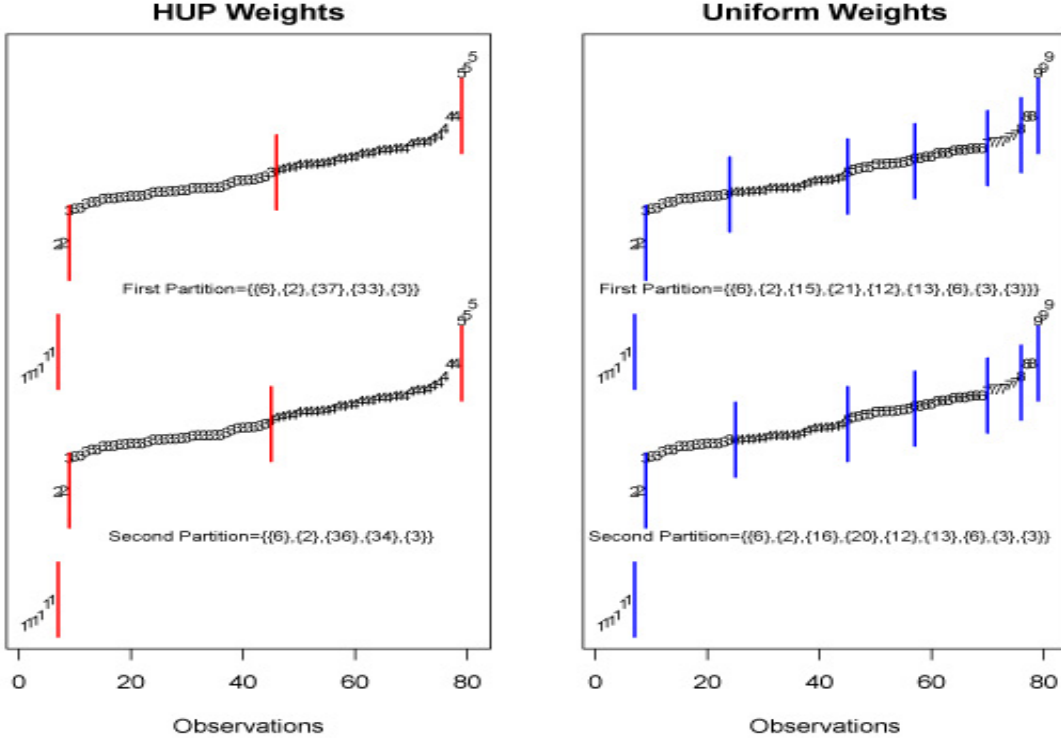
It has been established (see Kyung *et al.* 2010) that the Ewens-Pitman distribution (8) is the stationary distribution of the transition kernel $K(\mathbf{n}_p, \mathbf{n}'_{p'}) = \int_{\mathbf{q}} P(\mathbf{n}'_{p'}|\mathbf{q}) f(\mathbf{q}|\mathbf{n}_p) d\mathbf{q}$. Moreover, to sample from this distribution we can generate a candidate according to the multinomial distribution $\binom{n}{n_1 \dots n_p} \prod_{j=1}^p q_j^{n_j}$ and apply an MH step with ratio $\frac{K(\mathbf{n}_p, \mathbf{n}'_{p'})}{K(\mathbf{n}'_{p'}, \mathbf{n}_p)} = \lambda^{p'-p} \frac{\prod_{j=1}^{p'} \Gamma(n'_j)}{\prod_{j=1}^p \Gamma(n_j)}$.

Finally, we take our search to be a mixture of the biased random walk and the jump, choosing the biased random walk with probability a , set by the user. In our searches we have taken $a = .75$.

8 Examples

In this section we give a number of examples to illustrate the working of the clustering algorithm, and examine the effect of the choice of the prior on model space. We start with the well known Galaxy data (Roeder 1990) as a benchmark. We then look at some simulated examples, where we check that the HUP intrinsic Bayesian procedure gives the highest weight to the correct model, and

Figure 4: GALAXY DATA (ROEDER 1990), 81 OBSERVATIONS. THE LEFT PANEL SHOWS THE TOP TWO CLUSTERS FROM THE SEARCH ALGORITHM WITH HUP WEIGHTS ON THE BAYES FACTORS, WITH RED LINES DELIMITING THE CLUSTERS. THE RIGHT PANEL SHOWS THE TOP TWO CLUSTERS FROM THE SEARCH ALGORITHM WITH UNWEIGHTED BAYES FACTORS, WITH BLUE LINES DELIMITING THE CLUSTERS. IN EACH PANEL THE TOP AND BOTTOM CLUSTERS DIFFER BY ONE OBSERVATION MARKING A SHIFT POINT IN THE LONG MIDDLE STRING. THE STOCHASTIC SEARCH WAS RUN FOR 50,000 ITERATIONS.



then we look at the effectiveness of the search algorithm. We then apply our method to two data sets, and provide comparisons with other algorithms.

8.1 Galaxy Data

The Galaxy data consists of 81 observations on the velocity (km/second) of 81 galaxies in the Corona Borealis Region. It is well accepted that there are between 5 and 7 clusters in the data (Richardson and Green 1997, Jasra *et al.* 2005). Using an intercept-only model, we ran our algorithm with HUP and uniform weights.

The results are displayed in Figure 4, where we see that the search with HUP weights gave five clusters in the top partition with Bayes factor= 6.2×10^{10} , while the second partition, which differs by the switching of one galaxy, has Bayes factor= 4.9×10^{10} . The top 25 partitions in the search all had 5 clusters. In contrast, the unweighted uniform search found partitions with 9 clusters, which by consensus is too many clusters. From the results in Section 6, this performance is expected.

Table 2: FOR $n = 7$ AND $n = 9$, 25 DATA SETS WERE GENERATED FROM THE INDICATED CONFIGURATION CLASS AND BAYES FACTORS WERE CALCULATED FOR EVERY PARTITION OF $n = 7$ OR $n = 9$ OBSERVATIONS. FOR THE POSTERIOR ODDS COLUMNS, THE NUMBER IS THE AVERAGE PERCENTILE OF THE POSTERIOR ODDS FOR THE TRUE MODEL. FOR K-MEANS WE CALCULATED THE PROPORTION OF TIMES THAT K-MEANS FOUND THE TRUE PARTITION OUT OF TEN TRIES. K-MEANS WAS ALWAYS STARTED AT THE NUMBER OF CLUSTERS IN THE TRUE PARTITION.

Configuration Class	Posterior Odds			K-means
	Uniform	Ewens-Pitman	HUP	
$n = 7$ (7)	0.942	0.999	0.999	---
(3, 4)	0.919	0.971	0.945	0.440
(2, 2, 3)	0.908	0.826	0.919	0.470
(1, 2, 2, 2)	0.787	0.446	0.713	0.390
(1, 1, 1, 2, 2)	0.426	0.131	0.402	0.450
$n = 9$ (9)	0.949	1.000	1.000	---
(3, 3, 3)	0.927	0.922	0.846	0.140
(2, 3, 4)	0.902	0.933	0.909	0.270
(2, 2, 2, 3)	0.967	0.809	0.939	0.160
(1, 2, 3, 3)	0.875	0.755	0.734	0.210
(1, 2, 2, 2, 2)	0.944	0.469	0.891	0.130

8.2 Evaluating the Models

Next we evaluate the ability of the HUP intrinsic Bayes procedure to find the best model, regardless of the search algorithm. We look at two examples, for $n = 7$ and $n = 9$, where we can enumerate all of the models. We generate the data from the model

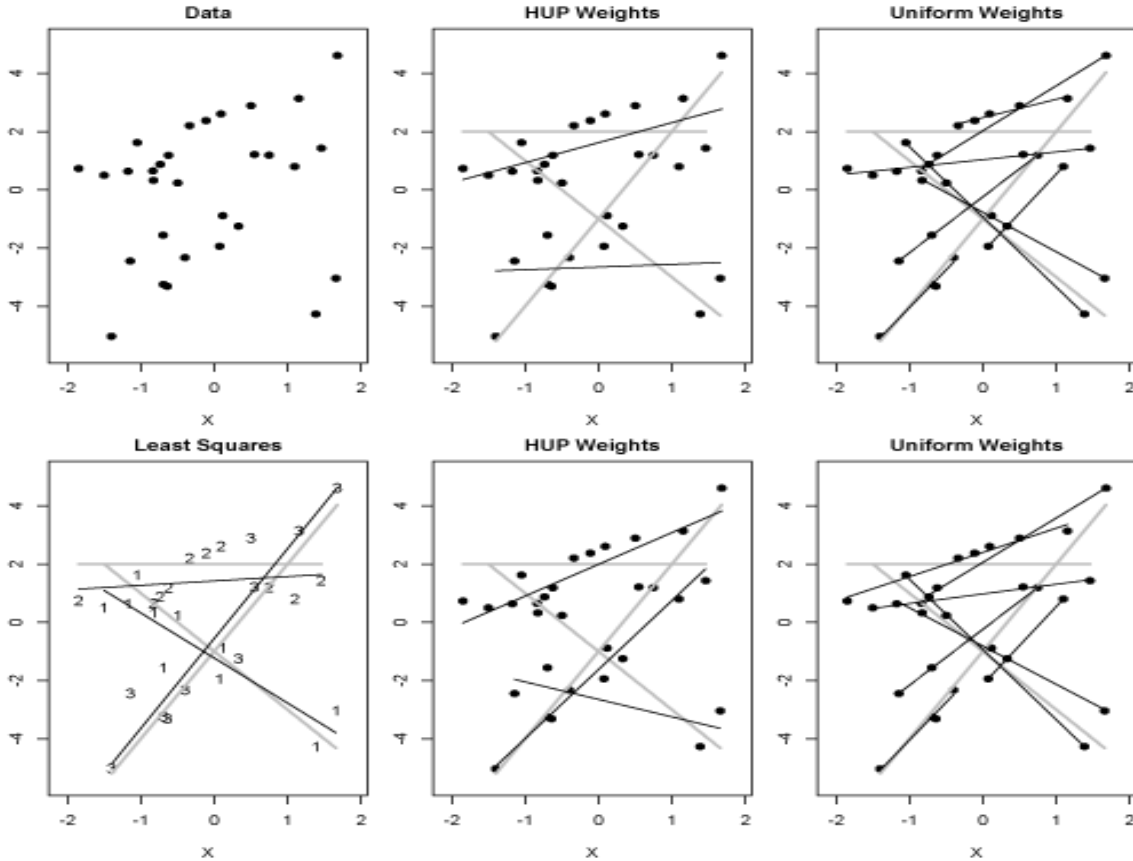
$$Y_{ij} \sim N(\mu_i, 1), \quad j = 1, \dots, n_i, \quad i = 1, \dots, k, \quad \mu_i = i, \quad (17)$$

where $\sum_i n_i = n$, $n = 7$ or 9 .

For a particular configuration of n and n_i we generated 25 data sets. For each data set we calculated all of the Bayes factors. For $n = 7$ there are 877 distinct partitions, and for $n = 9$ there are 21, 147 distinct partitions. For each of the 25 data sets we checked if the posterior odds of a model using the uniform prior, the Ewens-Pitman prior, and the HUP prior, was in the top 10 models

The results are shown in Table 2, where we see excellent performance of the Bayes factors with HUP weights. The Ewens-Pitman weight does well except when there are many small clusters (such as (1, 1, 1, 2, 2), begin dominated by the HUP weights. However, for these cases the unweighted Bayes factor does the best, reflecting its bias toward partitions with a large number of clusters. However,

Figure 5: SIMULATED DATA ($n = 30$) FROM MODEL (18). THE TWO LEFTMOST PANELS SHOW THE DATA (UPPER) WITHOUT CLUSTER LABELS. THE LOWER LEFTMOST PANEL SHOWS THE CLUSTER-IDENTIFIED DATA ALONG WITH THE TRUE MODELS (GREY) AND THE LEAST SQUARES FIT BASED ON KNOWING THE CLUSTER MEMBERSHIP. THE MIDDLE PANELS SHOW TYPICAL RESULTS FROM THE SEARCH WITH HUP WEIGHTS, AND THE RIGHT PANELS SHOW TYPICAL RESULTS FROM UNWEIGHTED SEARCHES. THE STOCHASTIC SEARCH WAS RUN FOR 50,000 ITERATIONS.



these are the only cases where the unweighted Bayes factor does better than the HUP weights.

The comparison with K-means shows that K-means has a very difficult time in identifying the true model. Since K-means only returns one partition, we could not calculate percentiles, but instead gave it ten tries to identify the true model. We also helped out K-means by starting with the correct number of clusters, that is, telling it to find a partition with the same number of clusters as the true cluster. Even with this help its performance was well below that of the Bayes factors.

8.3 Simulated Regression Data

Next we give an example of the search algorithm using data simulated from the regression model

$$Y_i \sim \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i = 1, \dots, n_i, \quad \varepsilon_i \sim N(0, 1) \quad (18)$$

with

$$\begin{aligned} \text{First } n_i \text{ observations} & : \beta_0 = -1, \beta_1 = -2 \\ \text{Second } n_i \text{ observations} & : \beta_0 = 2, \beta_1 = 0 \\ \text{Third } n_i \text{ observations} & : \beta_0 = -1, \beta_1 = 3 \end{aligned}$$

and the x_i are generated uniformly in $(0, 10)$. We actually did a large number of simulations, using data sets of different sizes and different configurations. Here we only present a typical result; the other simulations were very similar.

Our example has $n_i = 10$, and we use the data shown in Figure 5. The algorithm was run for 50,000 iterations and representative partitions, from the top 25 using HUP weights, and the top 25 using unweighted Bayes factors, are shown in the figure.

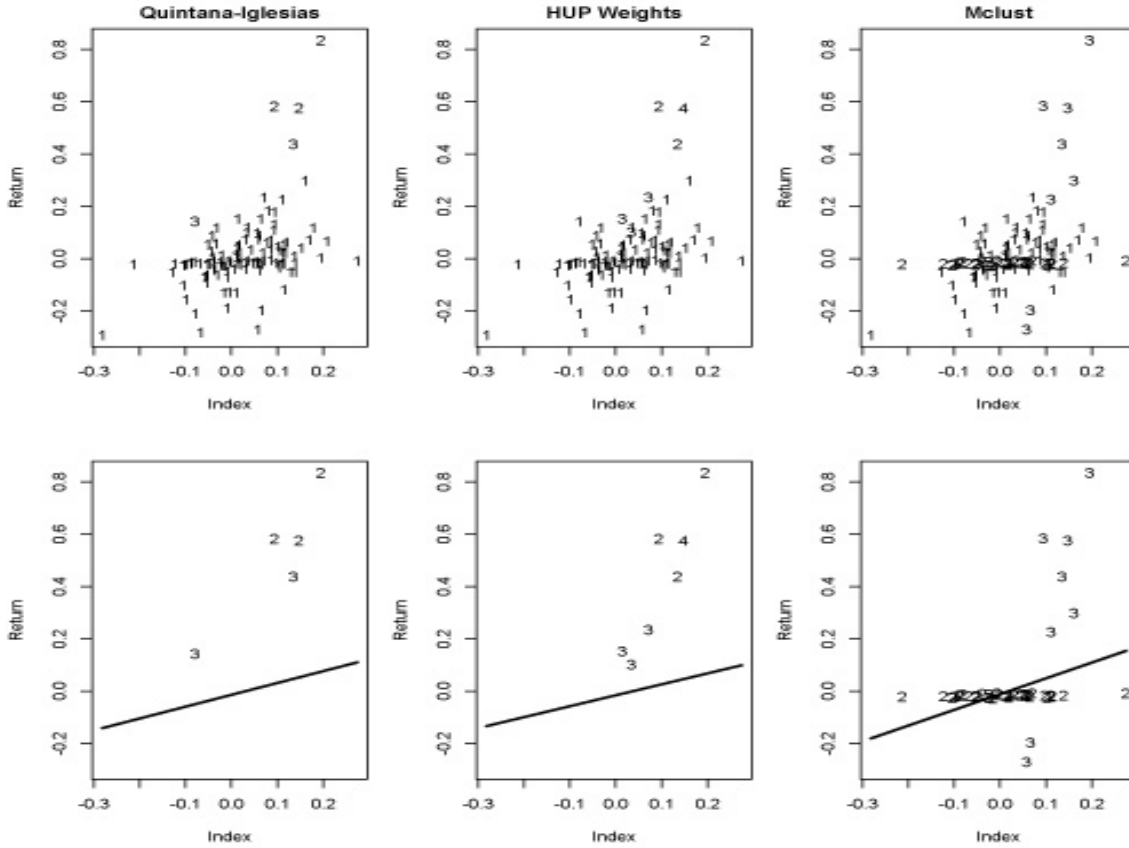
First note that when looking at the data without clusters identified (upper left panel) it is quite difficult to discern what the true clusters might be. The lower left panel shows that least squares, with knowledge of the true clusters, does reasonably well. In this light, the performance of the Bayes factor with HUP weights is quite remarkable. The top 25 partitions all had either 2 or 3 clusters, and the two partitions that we show are representative. The three cluster partition, in particular, does a very good job of recovering the underlying structure.

Similar to what we saw with the galaxy data, this example shows that searching large data sets without prior weights on the partitions leads to finding partitions with too many clusters. The rightmost panels in Figure 5 are representative of the top 25 partitions from the unweighted Bayes factor search, all of which had seven clusters. The underlying structure is not recovered. As mentioned before, one reason why we are doing cluster analysis is to find partitions with a small number of meaningful clusters. As of now it seems that the best way to accomplish this is to have HUP weights on the partitions.

8.4 Analyses and Comparisons

In this section we look at two different data sets, from Economics and Biostatistics. We see that in all cases the HUP Intrinsic Bayes procedure performs extremely well, not only giving reasonable answers in its own right, but also comparing favorable with other approaches.

Figure 6: RESULTS FROM THE ANALYSIS OF THE *Concha y Toro* DATA. FROM LEFT TO RIGHT THE PANELS CORRESPOND TO QUINTANA AND IGLESIAS, INTRINSIC BAYES WITH HUP WEIGHTS, AND MCLUST, RESPECTIVELY. THE TOP PANELS SHOW SCATTERPLOTS OF THE STOCK INDEX AGAINST GROWTH, WITH THE PLOTTING CHARACTER CORRESPONDING TO THE CLUSTER IDENTIFIER. THE LOWER PANELS SHOW THE REGRESSION FIT TO THE LARGEST CLUSTER, PLOTTING THE REMAINING POINTS AS OUTLIERS.



8.4.1 Chilean Stock Market

Quintana and Iglesias (2003) (QI) analyze economic data pertaining to the winemaker *Concha Y Toro*. This is simple linear regression data, using a model of the form

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i,$$

where y = the *Concha Y Toro* stock return, and x = a stock market index, similar to the US Dow-Jones Index. The data are fully described by QI.

QI use a version of their full PPM model set up for outlier detection, as they are interested in seeing if the *Concha Y Toro* returns follow the market. The model they fit has a common slope parameter, and they use individual intercepts to create the clusters. They use a PPM Gibbs sampler with the algorithm of Bush and MacEachern (1996). For a variety of parameter choices their analysis supports partitioning the data into a small number of clusters, three or four, where one cluster is

very large (signifying the data without outliers) and the other clusters are very small, containing the outliers. The partition that they obtained, assuming normal errors, is displayed in Figure 6.

We ran the data with our algorithm using only default settings and HUP weights, that is, we did not tune the model for outliers. The results from that analysis are also shown in Figure 6. We found four clusters, one large one containing the data without outliers, and three others which can clearly be considered outliers. Thus, our default analysis give results that were very similar to those of QI, with the exception that our algorithm was only set up to find clusters, not specifically to find outliers.

Finally, we also ran Mclust on the data, which also found three clusters. However, it did not find one large one and two small ones but rather two large ones and one small one. The larger cluster from Mclust had a slope similar to the large clusters found by QI and the HUP Bayes procedure, but the second largest cluster found by Mclust was not found by the others; both QI and HUP Bayes put those observations into the first cluster.

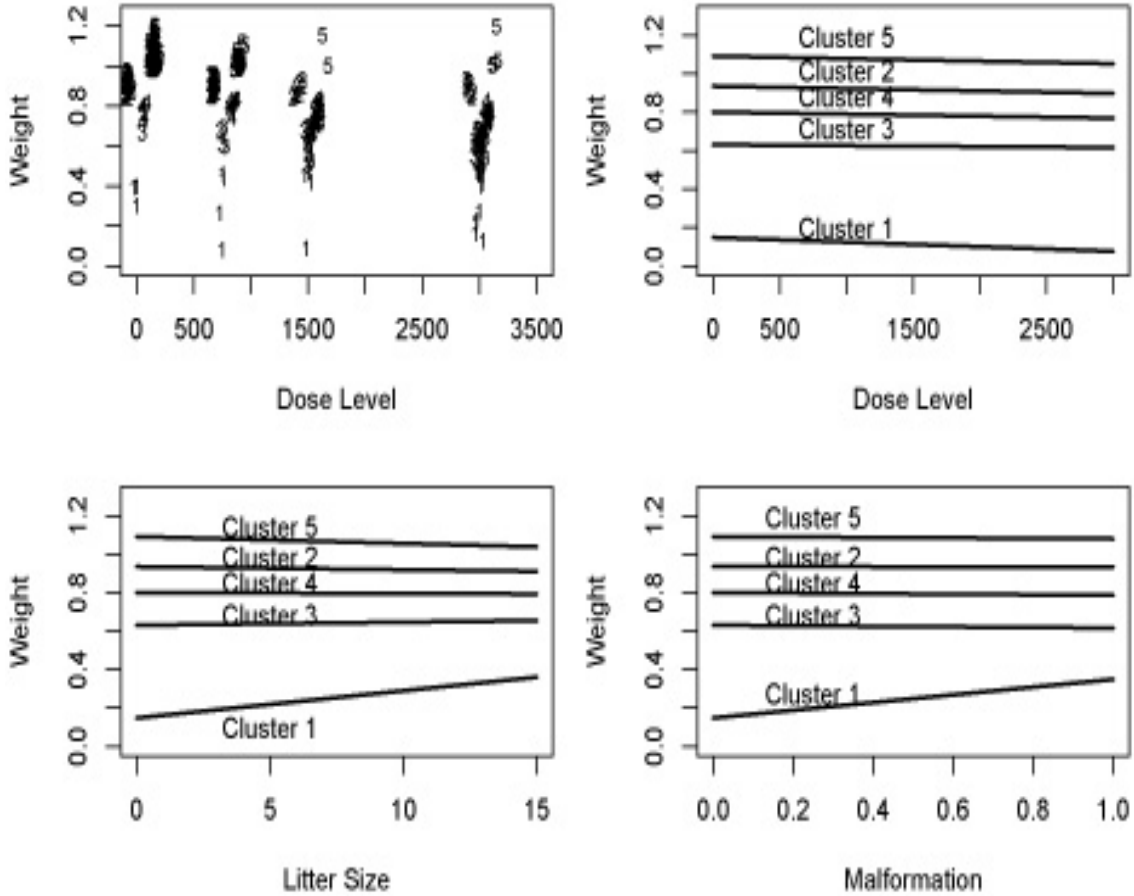
8.4.2 Developmental Toxicology Data

We look at the data analyzed by Dunson *et al.* (2003), and many other authors (see the references therein). It is data of a developmental toxicity study of ethylene glycol in mice conducted by the National Toxicology Program, and first reported in Price *et al.* (1985). During pregnancy, mice were exposed to four levels of ethylene glycol (0, 0.75, 1.5, and 3 mg/kg). The response of interest is the fetal weight of the babies. Other covariates were measured but here, like Dunson *et al.* (2003), we focus on two others (in addition to dose level): litter size, and a 0 – 1 indicator of malformations.

After removing some observations with missing data, the remaining data set had $n = 1048$ observations. We ran a stochastic search for 50,000 iterations using the Bayes factors with HUP weights; we did not use the unweighted Bayes factor due to its previous poor performance. The results are shown in Figure 7, which shows a partition that is typical of the top 10. (All of the top 10 partitions had five clusters, and there was little difference among them.)

The five clusters are quite interesting, showing that the partition aligns the mice on the intercepts, with fetal weight increasing as we move from Cluster 1 to Cluster 5 and with the slopes having little effect. The effect of dose level is similar in the clusters (parallel lines), decreasing the fetal weight slightly at higher doses. However, it is clear that in this partition the effect of the ethylene glycol dose is independent of the fetal weight. Also, the effect of litter size on fetal weight is also minimal in the partition. The only substantial slope effects are in Cluster 1 for the litter size and malformation. There we see an effect of increasing weight with litter size, and that zero malformations align with the lower fetal weight mice, and increased malformations as associated with higher fetal

Figure 7: RESULTS FROM THE CLUSTER ANALYSIS OF THE DEVELOPMENTAL TOXICOLOGY DATA. SHOWN ARE THE RESULTS OF THE PARTITION WITH THE HIGHEST POSTERIOR ODDS; THE TOP 10 PARTITIONS ALL HAD FIVE CLUSTERS. THE UPPER LEFT PANEL SHOWS THE DATA IDENTIFIED BY CLUSTER ASSIGNMENT, AND THE REMAINING PANELS SHOW THE SLOPES OF THE CLUSTERS FOR EACH OF THE THREE EXPLANATORY VARIABLES. THE STOCHASTIC SEARCH WAS RUN FOR 50,000 ITERATIONS. THE CLUSTER SIZES ARE (23, 341, 187, 301, 196). (FOR CLARITY, NOT ALL POINTS ARE PLOTTED IN THE CLUSTERS, AND THEY ARE JIGGLED.)



weights. However, in the other clusters, where all fetal weights are higher, there is no effect due to malformation.

9 Discussion

We have presented an objective Bayesian analysis for clustering based on product partition models using a model selection approach. Our major finding concerns the sensitivity of the procedure to the choice of the prior on model space, and we have seen that the “default” uniform prior on models leads not only to an inconsistent procedure, but also to small sample performance that is not desirable.

Our preferred prior for models, the hierarchical uniform prior, arises from a decomposition of the set of partitions of the sample in classes dictated by the number of clusters. Each class has also been split in subclasses where in each subclass we group all the partitions that only differ by a permutation of the components of the sample, and thus come from the same sampling model. This alleviates the difficulty of assigning a prior to the partitions inside the classes. A truncated Poisson-Intrinsic prior has been chosen for the number of clusters; it gives decreasing probability to partitions with higher numbers of clusters, and has performed well in our evaluations.

We also note the following about clustering priors:

- Cluster analysis is only useful, and will only result in useful inferences, when the answer contains a relatively small number of clusters. The prior should move us toward partitions with a small number of clusters, so the clusters themselves are large.
- Even if the truth is that there are a large number of clusters (say 500 observations have 70 true clusters) this results in a useless inference. In such a case it is better to find partitions with a small number of clusters that explain a large portion of the variability.
- The inconsistency of the uniform prior on the space of models, and its less than desirable small sample performance, eliminates it from serious consideration in clustering problems. It ends up driving the search toward partitions with many small clusters. The other two priors that we looked at, the Ewens-Pitman and Jensen-Liu are clearly preferable to the uniform prior, but still cannot be preferred over the HUP.
- The limit results of Theorem 2 also point out some undesirable behavior. The fact that only the HUP converges to a proper prior tells us that it is correctly compensating for the increasing number of models, which the other priors do not do.

Other points that we would like to emphasize are:

- The findings in our examples are consistent with the theory. The HUP produces partitions with a small number of clusters, while the unweighted Bayes factors almost always return a configuration with a large number of small clusters.
- We can apply some of our previous results (Casella *et al.* 2009, Moreno *et al.* 2010) to show that our procedure is consistent for choosing between two nested cluster models, as for instance $M_{\mathbf{r}_p}$ and $M_{\mathbf{r}_1}$, assuming that one of them is the true one. In such a pairwise comparison the prior on model space plays no role, as there are only two models. Pairwise consistency holds when the number of clusters grows at the rate $p = O(n^\alpha)$ for $\alpha < 1$, and for $\alpha = 1$ when the models $M_{\mathbf{r}_1}$ and $M_{\mathbf{r}_p}$ are not too close.

- Lemma 1 shows that for large n , when sampling from $M_{\mathbf{r}_1}$, the Bayes factor $B_{\mathbf{r}_p\mathbf{r}_1}(\mathbf{y})$ is constant across partitions \mathbf{r}_p in the class \mathfrak{R}_p . This asymptotic result is not limited to intrinsic Bayes factors; for example we know that BIC is asymptotically equivalent to the intrinsic Bayes factor. Moreover, Casella *et al.* (2009) show that the approximation holds for a wide class of priors.

10 References

- Aldous, D. J. (1985). Exchangeability and Related Topics. *É cole d'Été de Probabilités de Saint-Flour XII, Springer Lecture Notes in Mathematics, Volume 1117*. P. L. Hennequin, ed. New York: Springer-Verlag.
- Barry, D. and Hartigan, J.A. (1992). Product partition models for change point problems, *Annals of Statistics*, **20**, 260–279.
- Berger, J.O., Bernardo, J.M., and Sun, D. (2009). The formal definition of reference priors. *Annals of Statistics*, **37**, 905–938.
- Berger, J.O. and Mortera, J. (1999). Default Bayes factors for one-sided hypothesis testing. *J. Amer. Statist. Assoc.*, **94**, 542–554.
- Berger, J.O. and Pericchi, L.R. (1996a). The intrinsic Bayes factor for model selection and prediction. *J. Amer. Statist. Assoc.*, **91**, 109–122.
- Berger and Pericchi (1996b). The intrinsic Bayes factor for linear models. In *Bayesian Statistics 5*, Oxford University Press, pp. 23–42.
- Booth, J.G., Casella, G. and Hobert, J.P. (2008). Clustering using objective functions and stochastic search. *J. R. Statist. Soc. B*, **70**, 1, 119–139.
- Bush, C.A., and MacEachern, S.N. (1996). A Semi-parametric Bayesian Model for Randomised Block Designs. *Biometrika* **83** 275–286.
- Casella, G., Girón, F.J., Martínez, M.L. and Moreno E. (2009). Consistency of Bayesian procedure for variable selection. *Annals of Statistics*, **37**, 3, 1207–1228.
- Casella, G. and Moreno E. (2005). Intrinsic meta analysis of contingency tables. *Statistics in Medicine*, 24, 583–604.
- Casella, G. and Moreno E. (2009). Assessing Robustness of Intrinsic Tests of Independence in Tway Contingency Tables. *J. Amer. Statist. Assoc.*, 104, 1261–1271.
- Crowley E. (1997). Product partition models for normal means, *J. Amer. Statist. Assoc.*, **92**, 192–198.
- Dunson, D.B., Chen, Z. and Harry, J. (2003). A Bayesian Approach for Joint Modeling of Cluster Size and Subunit-Specific Outcomes. *Biometrics* **59** 521–530.
- Ewens, W. J. (1972). The Sampling Theory of Selectively Neutral Alleles. *Theor. Popul. Biol* **3** 87–112
- Fraley C. and , Raftery A. E. (2002) Model-based clustering, discriminant analysis, and density estimation. *Journal of the American Statistical Association*, **97** 611–631.
- Fraley C. and , Raftery A. E. (2007) Bayesian regularization for normal mixture estimation and model-based clustering. *Journal of Classification* **24** 155–181.

- Fraley C. and , Raftery A. E. (2009). MCLUST Version 3 for R: Normal Mixture Modeling and Model-based Clustering. Technical Report No. 504, Department of Statistics, University of Washington, September 2006 (revised December 2009).
- Girón, F.J., Martínez, M.L., Moreno, E. and Torres, F. (2006). Objective testing procedures in linear models. Calibration of the p-values. *Scandinavian Journal of Statistics* **33** 765-787.
- Girón, F.J., Moreno, E., Casella, G. and Martínez, M.L. (2010). Consistency of objective Bayes factors for nonnested linear models and increasing model dimension. *RACSAM* **104** (1), 61–71.
- Hartigan, J. A. (1990). Partition Models. *Communications in Statistics - Theory and Methods* **19** 2745–2756.
- Ishwaran, H. and Zarepour, M. (2002). Dirichlet Prior Sieves in Finite Normal Mixtures. *Statistica Sinica* **12** 941-963.
- Jasra, A., Holmes, C. C, and Stephens, D. A. (2005). Markov Chain Monte Carlo Methods and the Label Switching Problem in Bayesian Mixture Modeling. *Statistical Science* **20** 50-67.
- Jensen S. T. and Liu J. S. (2008). Bayesian Clustering of Transcription Factor Binding Motifs. *J. Amer. Statist. Assoc.* **103** 188-200
- Kim, S. and Sun D. (2000). Intrinsic priors for model selection using an encompassing model. *Life Time Data Analysis*, 6, 251-269.
- Kyung, M., Gill, J. and Casella, G. (2010). Estimation in Dirichlet Random Effects Models. *Ann. Statist* **38** 979-1009
- Lau, J. W. and Green, P. J. (2007). Bayesian Model-Based Clustering Procedures. *Journal of Computational and Graphical Statistics* **16** 526-558.
- McCullagh, P. and Yang, J. (2008). Stochastic Classification Models. *International Congress of Mathematicians III* 669-686.
- Moreno, E. (1997). Bayes factor for intrinsic and fractional priors in nested models: Bayesian robustness. In *L₁-Statistical Procedures and Related Topic*, Institute of Mathematical Statistics: Hayward, CA, pp. 257-270.
- Moreno, E., Bertolino, F. and Racugno, W. (1998). An intrinsic limiting procedure for model selection and hypothesis testing. *J. Amer. Statist. Assoc.*, 93, 1451-1460.
- Moreno, E. (2005). Objective Bayesian analysis for one-sided testing. *Test*, 14, 181-198.
- Moreno, E., Casella, G. and García-Ferrer, A. (2005). An objective Bayesian analysis of the changepoint problem. *SERRA* **19** 191-204.
- Moreno, E., Girón, F.J. and Casella, G. (2010). Consistency of objective Bayesian tests as the model dimension increases. *Annals of Statistics* **38** 1937-1952
- Neal, R. M. (2000). Markov chain sampling methods for Dirichlet process mixture models. *Journal of Computational and Graphical Statistics* **9** 249-265.
- Pitman, J. (1996). Some Developments of the Blackwell-MacQueen Urn Scheme. *Statistics, Probability and Game Theory* 245-267. IMS Lecture Notes Monograph Series, 30, Institute of Mathematical Statistics, Hayward, CA.
- Price, C. J., Kimmel, C. A., Tyl, R. W., and Marr, M. C.(1985). The developmental toxicity of Ethylene Glycol in Rats and Mice. *Toxicology and Applied Pharmacology* **81** 113-127.

Quintana, F.A. and Iglesias, P.L. (2003). Bayesian clustering and product partition models. *J. Roy. Statist. Soc. Ser. B* **65** 557-574.

Richardson, S. and Green, P. J. (1997). On Bayesian analysis of mixtures with an unknown number of components. *J. Roy. Statist. Soc. Ser. B* **59** 731-758.

Roeder, K. (1990). Density estimation with confidence sets exemplified by superclusters and voids in the galaxies. *J. Amer. Statist. Assoc.* **85** 617-624.

Stephens, M. (2000). Dealing with label-switching problem in mixtures models. *J. Roy. Statist. Soc. Ser. B*, 62, 795-809.

A Technical Details

In this section we prove the asymptotic results in Theorems 3, 4, and 2. We begin with some preliminary lemmas that are needed, and then give detailed proofs of the theorems.

A.1 Preliminary Lemmas

Lemma 2 For every $p = 1, \dots$, under the conditions $n_1 \leq \dots \leq n_p$, and $\sum_{i=1}^p n_i = n$, the sum of the series

$$\sum \binom{n}{n_1 \dots n_p} (\lambda + 1)^{-n_1} \dots (\lambda + p)^{-n_p} \approx O\left(\frac{\lambda_p^H}{p}\right)^{-n},$$

where λ_p^H is the harmonic mean of $(\lambda + 1)^{-1}, \dots, (\lambda + p)^{-1}$.

Proof. To evaluate the sum, multiply and divide by $(\sum_{i=1}^p \frac{1}{\lambda+i})^n$. Sum the resulting multinomial to get

$$\sum \binom{n}{n_1 \dots n_p} (\lambda + 1)^{-n_1} \dots (\lambda + p)^{-n_p} = \left(\sum_{i=1}^p \frac{1}{\lambda+i}\right)^n \approx O\left(\frac{\lambda_p^H}{p}\right)^{-n}.$$

■

Lemma 3 For every p the sum

$$S = \sum_{C_p} \frac{1}{n_1 \times \dots \times n_p},$$

where $C_p = \{(n_1, \dots, n_p) : n_1 \leq \dots \leq n_p, \sum_{i=1}^p n_i = n\}$, is of order $O(n^{-1}(\log n)^{p-1})$.

Proof. Denote $x_i = n_i/n$, for $i = 1, \dots, p-1$. Then we can write the sum S as

$$S = n^{-p} \sum \frac{1}{x_1 \cdot x_2 \cdot \dots \cdot (1 - x_1 - x_2 - \dots - x_{p-1})},$$

where the multiple sum is extended for $x_i = 1/n, \dots, 1 - 1/n$ in steps of size $1/n$, for $i = 1, \dots, p-1$. Now, the sum can be approximated, for large values of n , by the multiple integral

$$n^{p-1} \int_{1/n}^{1-1/n} \dots \int_{1/n}^{1-1/n} \frac{1}{x_1 \cdot x_2 \dots (1 - x_1 - x_2 - \dots - x_{p-1})} dx_1 \dots dx_{p-1} = n^{p-1} \mathbf{I},$$

where the factor n^{p-1} is the adjustment due to the $1/n$ step for each variable, and \mathbf{I} is the integral. Thus, the original sum is

$$\sum \frac{1}{n_1 \dots n_p} \approx n^{-1} \mathbf{I}.$$

To evaluate the integral, consider the change of variables

$$\theta_1 = x_1, \quad \theta_i = \frac{x_i}{1 - \sum_{j=1}^{i-1} x_j}, \quad i = 2, \dots, p-1.$$

As the Jacobean J of the variables (x_1, \dots, x_{p-1}) with respect to the new variables $(\theta_1, \dots, \theta_{p-1})$ is $J = \prod_{i=1}^{p-1} (1 - \sum_{j=1}^{i-1} x_j)$, the integral \mathbf{I} in terms of the new variables can be written as

$$\mathbf{I} = \int_{\frac{1}{n}}^{1-\frac{1}{n}} \dots \int_{\frac{1}{n-p+2}}^{1-\frac{1}{n-p+2}} \frac{1}{\theta_1 \cdot \theta_2 \dots \theta_{p-1}} d\theta_1 \dots d\theta_{p-1}$$

which is equal to

$$\mathbf{I} = \prod_{i=1}^{p-1} \int_{\frac{1}{n-i+1}}^{1-\frac{1}{n-i+1}} \frac{1}{\theta_i} d\theta_i = \prod_{i=1}^{p-1} \log(n-i) \sim O((\log n)^{p-1}).$$

■

A.2 Proof of Theorem 2

For part(a), the proof is immediate as the distribution of \mathbf{X}_n on the simplex \mathcal{S}_p is a discrete uniform distribution on the points of the simplex of the form $(\frac{n_1}{n}, \dots, \frac{n_p}{n})$, where $n_1 \leq \dots \leq n_p$ and $n_1 + \dots + n_p = 1$. From the definition of the hierarchical uniform prior,

$$\Pr(\mathbf{X}_n = (\frac{n_1}{n}, \dots, \frac{n_p}{n}) | \mathfrak{R}_p) = \Pr(\mathfrak{R}_{p;n_1, \dots, n_p} | \mathfrak{R}_p) = \frac{1}{b(n, p)} \approx \frac{p!}{n^{p-1}},$$

and this discrete uniform prior on the lattice

$$L_p = \left\{ \left(\frac{n_1}{n}, \dots, \frac{n_p}{n} \right), \quad n_1 \leq \dots \leq n_p, \quad n_1 + \dots + n_p = 1, \right.$$

converges *in distribution* to a continuous uniform distribution on the simplex, which is a Dirichlet $\mathcal{D}(1, \dots, 1)$ truncated at the simplex \mathcal{S}_p .

For part (b), for the Ewens-Pitman prior, it follows that the prior probability of $\mathfrak{R}_{p;n_1, \dots, n_p}$, given p and λ , is

$$\Pr(\mathfrak{R}_{p;n_1, \dots, n_p} | \mathfrak{R}_p) \propto \binom{n}{n_1 \cdots n_p} \times \prod_{i=1}^p \Gamma(n_i) = \frac{1}{n_1 \cdots n_p}.$$

Thus, the distribution of \mathbf{X}_n on the simplex \mathcal{S}_p is

$$\Pr\left(\mathbf{X}_n = \left(\frac{n_1}{n}, \dots, \frac{n_p}{n}\right) \middle| \mathcal{S}_p\right) \propto \frac{1}{\frac{n_1}{n} \cdots \frac{n_p}{n}}.$$

It is clear from the form of the probability mass function of the discrete random vector \mathbf{X}_n , that the limiting distribution is given by the function (2).

For part (c), we can discard the redundancy term $R(n_1, \dots, n_p)$ for large n , as it is of an order of magnitude much smaller than $\binom{n}{n_1, \dots, n_p}$. We then write the prior probability of $\mathfrak{R}_{p;n_1, \dots, n_p}$ as

$$\Pr(\mathfrak{R}_{p;n_1, \dots, n_p} | \mathfrak{R}_p) \propto \binom{n}{n_1 \cdots n_p}.$$

This implies that

$$\Pr(\mathbf{Z}_n = (n_1, \dots, n_p) | \mathfrak{R}_p) \propto \binom{n}{n_1 \cdots n_p},$$

But, because of the proportionality symbol, we can write the preceding as

$$\Pr(\mathbf{Z}_n = (n_1, \dots, n_p) | \mathfrak{R}_p) \propto \binom{n}{n_1 \cdots n_p} \left(\frac{1}{p}\right)^n = \binom{n}{n_1 \cdots n_p} \left(\frac{1}{p}\right)^{n_1} \cdots \left(\frac{1}{p}\right)^{n_p},$$

and this means that the unrestricted \mathbf{Z}_n follows a multinomial distribution $\mathcal{M}(n; 1/p, \dots, 1/p)$.

For large n , because of the restriction $n_1 + \dots + n_p = n$, this multinomial distribution can be approximated by a multivariate normal distribution with the same mean vector and covariance matrix, that is,

$$\mathbf{Z}_n \approx N_p(n(1/p, \dots, 1/p)^t, \Sigma_p),$$

where the covariance matrix $\Sigma_p = (n/p^2)(p\mathbf{I} - \mathbf{J})$, where \mathbf{J} is a matrix of ones. Thus, the distribution of \mathbf{X}_n can be approximated by the following multivariate normal distribution

$$\mathbf{X}_n \approx N_p\left((1/p, \dots, 1/p)^t, \frac{1}{n^2} \Sigma_p\right),$$

which converges to the degenerate random variable at mean vector $(\frac{1}{p}, \dots, \frac{1}{p})$, as the covariance matrix of \mathbf{X}_n converges to the null matrix.

For part (d), from the expression of the Jensen-Liu prior (9), it follows that the prior probability of $\mathfrak{A}_{p;n_1,\dots,n_p}$, given p and λ , can be written

$$\Pr(\mathbf{Z}_n = (n_1, \dots, n_p) | \mathfrak{A}_p) = \Pr(\mathfrak{A}_{p;n_1,\dots,n_p} | \mathfrak{A}_p) \propto \binom{n}{n_1 \dots n_p} \prod_{i=1}^p \left(\frac{(\lambda + i)^{-1}}{\sum_{i'} (\lambda + i')^{-1}} \right)^{n_i}.$$

Thus, the distribution of \mathbf{Z}_n is the following multinomial distribution

$$\mathbf{Z}_n \sim \mathcal{M} \left(n; \frac{(\lambda + 1)^{-1}}{\sum_{i'} (\lambda + i')^{-1}}, \dots, \frac{(\lambda + p)^{-1}}{\sum_{i'} (\lambda + i')^{-1}} \right),$$

and, for large n , the distribution of \mathbf{X}_n can be approximated by the following multivariate normal distribution

$$\mathbf{X}_n \approx N_p \left(\frac{(\lambda + 1)^{-1}}{\sum_{i'} (\lambda + i')^{-1}}, \dots, \frac{(\lambda + p)^{-1}}{\sum_{i'} (\lambda + i')^{-1}} \right)', \frac{1}{n^2} \Sigma_{p,\lambda} \right),$$

which converges to a random variable degenerate at the mean vector as the covariance matrix of \mathbf{X}_n converges to the null matrix.

Note. The Jensen-Liu prior restricted to the simplex shows a similar behavior to the uniform prior on the set of all cluster models with at most T clusters, and both are inconsistent. Indeed, when λ goes to infinity, the Jensen-Liu prior converges to the Uniform prior on the set of all cluster models for all $p \leq T$.

A.3 Proof of Theorem 3

Uniform Prior First, we will prove that the posterior probability of \mathfrak{A}_T converges to 1, when n goes to infinity. The uniform prior on the set of all models results in the prior distribution on the cluster classes $Pr(\mathfrak{A}_p) \propto S_n^{(p)}$, which can be approximated for large n by $Pr(\mathfrak{A}_p) \approx \frac{p^n}{p!}$.

Applying Bayes theorem, we have that, for $p = 1, \dots, T$

$$\Pr(\mathfrak{A}_p | \mathbf{y}) \propto Pr(\mathfrak{A}_p) \times B_{p,1}(\mathbf{y}) \propto \frac{p^n}{p!} \left(\frac{kp+1}{n} \right)^{\frac{k(p-1)}{2}}.$$

Therefore, normalizing the Bayes factors, we have that as $n \rightarrow \infty$,

$$\Pr(\mathfrak{A}_p | \mathbf{y}) = \frac{\frac{p^n}{p!} \left(\frac{kp+1}{n} \right)^{\frac{k(p-1)}{2}}}{\sum_{p=1}^T \frac{p^n}{p!} \left(\frac{kp+1}{n} \right)^{\frac{k(p-1)}{2}}} \rightarrow \begin{cases} 0 & \text{for } p = 1, \dots, T-1, \\ 1 & \text{for } p = T. \end{cases}$$

To show inconsistency in the configuration classes, we have to show that within \mathfrak{A}_T the posterior distribution of the models does not converge in distribution to the degenerate distribution

corresponding to the vertex \mathfrak{R}_1 . In fact, from Theorem 2, under the uniform prior on the set of all models, this posterior distribution converges to the degenerate distribution on the equal size T clusters $\delta_{(\frac{1}{T}, \dots, \frac{1}{T})} \in \mathfrak{R}_T$. This completes the proof for the uniform distribution

Jensen-Liu From (9) it is easy to see that the marginal prior of p or, equivalently \mathfrak{R}_p , is

$$\Pr(\mathfrak{R}_p) \propto \lambda^{p-1}(\lambda + p) \sum \binom{n}{n_1 \dots n_p} (\lambda + 1)^{-n_1} \dots (\lambda + p)^{-n_p}.$$

Using the asymptotic approximation of Lemma 2 the prior can be approximated by

$$\Pr(\mathfrak{R}_p) \propto \lambda^{p-1}(\lambda + p) \left(\frac{\lambda_p^H}{p} \right)^{-n}.$$

Recalling the expression for the Bayes factor given in Lemma 1, we have from Bayes Theorem, for $p = 1, \dots, T$,

$$\Pr(\mathfrak{R}_p | \mathbf{y}) \propto \Pr(\mathfrak{R}_p) \times B_{p1}(\mathbf{y}) \propto \lambda^{p-1}(\lambda + p) \left(\frac{\lambda_p^H}{p} \right)^{-n} \left(\frac{k p + 1}{n} \right)^{\frac{k(p-1)}{2}}.$$

The leading term of the posterior probability is $(\lambda_p^H/p)^{-n}$, which is an increasing function of p for every n . Normalizing the posterior probabilities we finally get

$$\Pr(\mathfrak{R}_p | \mathbf{y}) \rightarrow \begin{cases} 0 & \text{for } p = 1, \dots, T - 1, \\ 1 & \text{for } p = T, \end{cases}$$

It is also the case that within the cluster class \mathfrak{R}_T the posterior distribution of the models does not converge in distribution to the degenerate distribution corresponding to the vertex \mathfrak{R}_1 . In fact, from Theorem 2, this posterior distribution converges to the degenerate distribution of the model

$M_{\left(\frac{\lambda_T^H}{p(\lambda+T)}, \dots, \frac{\lambda_T^H}{p(\lambda+1)}\right)} \in \mathfrak{R}_T$. This completes the proof.

Note that the limiting distribution in this case depends on the value of the parameter λ of the Jensen-Liu prior. Further, as λ increases, the limiting distribution converges to the T -cluster model with equal size clusters.

A.4 Proof of Theorem 4

Hierarchical Uniform Prior We will prove that the posterior distribution of \mathfrak{R}_1 converges to 1, when n goes to infinity. By Bayes theorem, and the fact that the prior over the cluster classes

$\mathfrak{R}_1, \dots, \mathfrak{R}_p$ is uniform, we have that, for $p = 1, \dots, T$,

$$\Pr(\mathfrak{R}_p|\mathbf{y}) \propto Pr(\mathfrak{R}_p) \times B_{p1}(\mathbf{y}) \propto B_{p1}(\mathbf{y}) = \left(\frac{kp+1}{n}\right)^{\frac{k(p-1)}{2}}.$$

Therefore, normalizing the Bayes factors, we have that, as $n \rightarrow \infty$,

$$\Pr(\mathfrak{R}_p|\mathbf{y}) = \frac{\left(\frac{kp+1}{n}\right)^{\frac{k(p-1)}{2}}}{\sum_{p=1}^T \left(\frac{kp+1}{n}\right)^{\frac{k(p-1)}{2}}} \rightarrow \begin{cases} 0 & \text{for } p = 2, \dots, T, \\ 1 & \text{for } p = 1. \end{cases}$$

Ewens-Pitman Prior From (8), the joint prior for p and \mathbf{r}_p , we can calculate the marginal prior of p or, equivalently \mathfrak{R}_p , as

$$\Pr(\mathfrak{R}_p) \propto \lambda^{p-1} n^{p-1} \sum \frac{1}{n_1 \cdots n_p}.$$

Again recalling Lemma 1, Bayes Theorem yields, for $p = 1, \dots, T$,

$$\Pr(\mathfrak{R}_p|\mathbf{y}) \propto \Pr(\mathfrak{R}_p) \times B_{p1}(\mathbf{y}) \propto \lambda^{p-1} n^{p-1} \sum \frac{1}{n_1 \cdots n_p} \left(\frac{kp+1}{n}\right)^{\frac{k(p-1)}{2}}.$$

Using Lemma 3, we can approximate the posterior by

$$\Pr(\mathfrak{R}_p|\mathbf{y}) \approx C_p \lambda^{p-1} n^{-1} (\log n)^{p-1} \left(\frac{kp+1}{n}\right)^{\frac{k(p-1)}{2}} \approx K_{p,k,\lambda} n^{-1} (n^{-k/2} \log n)^{p-1},$$

where C_p is a finite positive constant depending on p and $K_{p,k,\lambda}$ is a positive constant depending on p , k and λ . It is now clear that, as $n \rightarrow \infty$,

$$\Pr(\mathfrak{R}_p|\mathbf{y}) \rightarrow \begin{cases} 0 & \text{for } p = 2, \dots, T, \\ 1 & \text{for } p = 1. \end{cases}$$