

Estimating Mixtures of Regressions

Merrilee HURN, Ana JUSTEL, and Christian P. ROBERT

This article shows how Bayesian inference for switching regression models and their generalizations can be achieved by the specification of loss functions which overcome the label switching problem common to all mixture models. We also derive an extension to models where the number of components in the mixture is unknown, based on the birth-and-death technique developed in recent literature. The methods are illustrated on various real datasets.

Key Words: Bayesian inference; Birth-and-death process; Label switching; Logistic regression; Loss functions; MCMC algorithms; Poisson regression; Switching regression.

1. INTRODUCTION

The *switching regression* model is well known in the econometrics literature; it arises when an observed quantity y depends on a vector of covariates x in a linear way

$$y = x' \beta_i + \sigma_i \epsilon, \quad \epsilon \sim g(\epsilon), \quad (1.1)$$

where the (β_i, σ_i) 's ($i = 1, \dots, k$) vary among a set of k possible values with probabilities p_1, \dots, p_k . In other words, assuming a normal distribution on the perturbation ϵ , the conditional distribution of y given x is a mixture of normal distributions

$$y|x \sim p_1 \mathcal{N}(x' \beta_1, \sigma_1^2) + \dots + p_k \mathcal{N}(x' \beta_k, \sigma_k^2), \quad (1.2)$$

as illustrated by Figure 1 in a simple regression case.

This model was introduced by Goldfeld and Quandt (1976) and has been mainly studied from a likelihood point of view. Its appeal is clear when considering datasets such as those in Figure 2(a), which describes the equivalence ratio, that is, the richness of the air-ethanol mix in an engine, against the concentration of nitrous oxide in exhaust, as presented in

Merrilee Hurn is Lecturer in Statistics, Department of Mathematical Statistics, University of Bath, UK (E-mail: mah@maths.bath.ac.uk). Ana Justel is Profesora Titular de Estadística, Department of Statistics, Universidad Autónoma de Madrid, Spain (E-mail: ana.justel@uam.es). Christian P. Robert is Professor de Statistique, Université Paris Dauphine and Head, CREST, Insee, Paris (E-mail: xian@ceremade.dauphine.fr).

©2003 American Statistical Association, Institute of Mathematical Statistics,
and Interface Foundation of North America

Journal of Computational and Graphical Statistics, Volume 12, Number 1, Pages 1–25
DOI: 10.1198/1061860031329

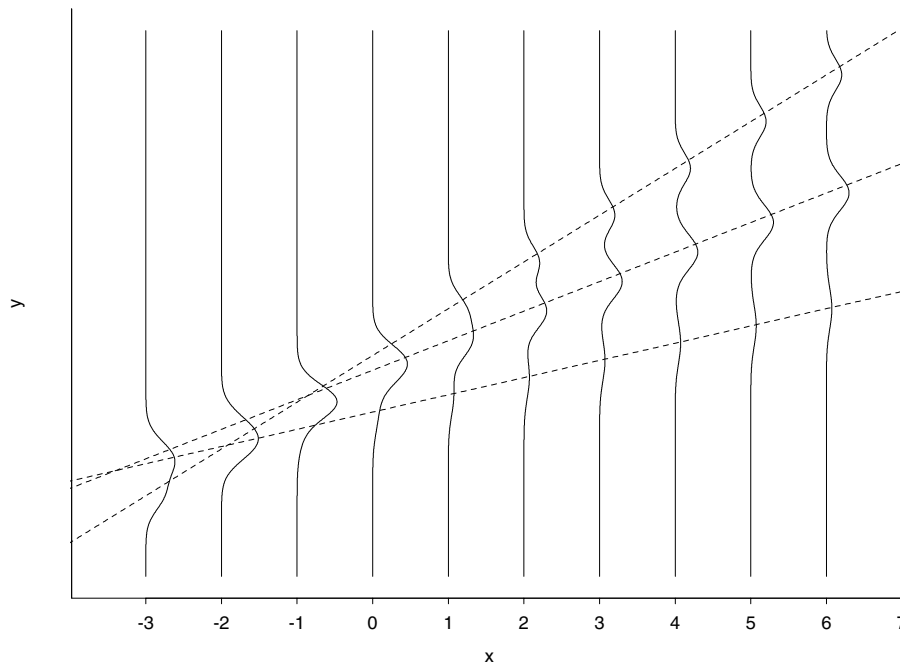


Figure 1. Conditional distributions of y for several values of x , against the regression lines.

Hurvich, Simonoff, and Tsai (1998), and in Figure 2(b), which plots the pairs (GNP,CO₂) for various countries, where GNP denotes the gross national product per capita in 1996 and CO₂ the estimated CO₂ emission per capita the same year. The left graph clearly indicates two different nitrous oxide concentration dependencies. The data in the right graph is quite spread out, with no clear linear trend; however, there do seem to be several groups for which a linear model would be a reasonable approximation. To identify those groups and the corresponding linear models is of interest for low GNP countries as it may help to clarify on which development path they are embarking.

As well as econometrics and chemometrics, there exist many other areas where heterogeneous covariate dependent populations are found and studied, both at the classification level, namely to identify the homogeneous subpopulations, and at the inference level, namely to estimate the corresponding models. See for instance Bar-Shalom (1978), Quandt and Ramsey (1978), or Kiefer (1980). Extensions of (1.1) cover time-dependent models such as hidden Markov models and nonlinear switching regressions; for example, Hamilton (1989), Shumway and Stoffer (1991), or Billio, Monfort, and Robert (1999).

The framework is also wider than just simple linear regressions. We can for instance consider mixtures of logistic regressions:

$$P(y = 1|x) = \sum_{i=1}^k p_i \frac{\exp(x'\beta_i)}{1 + \exp(x'\beta_i)}, \quad (1.3)$$

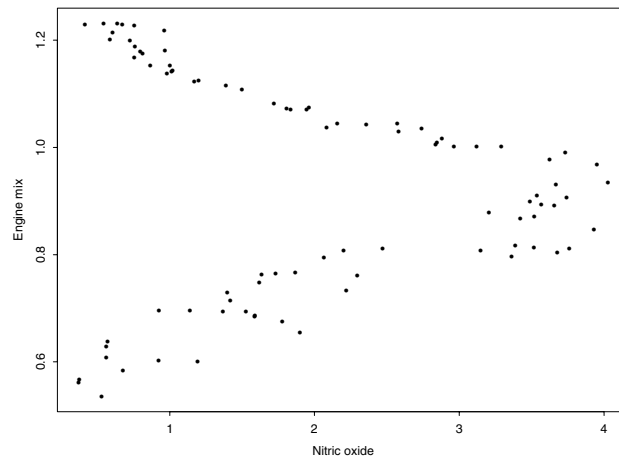
as seen later in an example describing investment choices based on socio-economic covari-

ates, and mixtures of Poisson regressions,

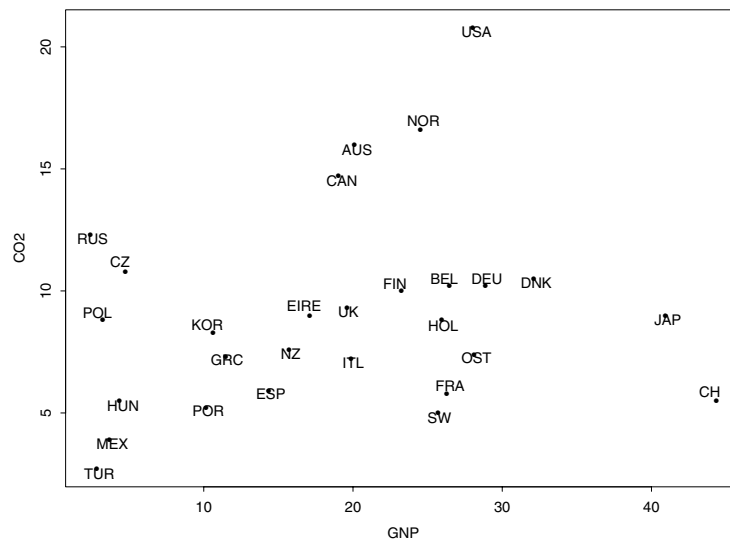
$$y|x \sim \sum_{i=1}^k p_i \mathcal{P}(\exp(x'\beta_i)) , \tag{1.4}$$

as illustrated on two datasets relating numbers of accidents and covariates.

The article is organized as follows: Section 2 describes a Bayesian approach to these models and uncovers the difficulties which arise in their processing, including the



(a)



(b)

Figure 2. (a) Equivalence ratio against exhaust nitric oxide concentration (Source: Hurvich et al., 1998); (b) representation of the GNP and CO₂ emission levels in 1996 for various countries (Source: OECD).

specificities of MCMC algorithms that must be used in such settings and discussing their limitations. Section 3 discusses different choices of loss functions depending on different inferential imperatives and designed to overcome these difficulties. Section 4 presents a natural extension to unknown numbers of components, based on the birth-and-death process technology, for the normal regression as well as the two regression models (1.3) and (1.4). We illustrate those methods on several simulated and real datasets.

2. A BAYESIAN APPROACH

This section briefly reviews some of the aspects of the Bayesian processing of such models, and describes the inferential difficulties which arise as a result of a lack of identifiability.

2.1 PRIOR SPECIFICATION

The inferential difficulties we describe in Section 2.3 are not particular to a specific prior modeling, but rather appear as a generic issue. We will thus choose for our simulations and applications standard conjugate priors, that is, independent normal priors on the regression coefficients β_i 's, such as $\mathcal{N}(0, 10)$ or $\mathcal{N}(0, 10\sigma_i^2)$, inverse gamma priors on the σ_i^2 's for (1.2), most often an $\mathcal{Exp}(1)$ distribution when dealing with a normalized dataset, and a Dirichlet prior on the weight vector (p_1, \dots, p_k) , $\mathcal{D}_k(1, \dots, 1)$. This particular choice allows us to use a Gibbs sampler but other choices are equally acceptable since they can be implemented via a Metropolis–Hastings algorithm, as described in the next paragraph.

Notice first that improper priors, hence standard flat priors, are not acceptable in this setup, because they lead to undefined posterior distributions, whatever the sample size (see, e.g., Robert 1996), and, second, that exchangeable priors on the components—that is, priors that are invariant under permutations of the component labels—produce posterior distributions with identical marginal distributions on all components. This feature is at the root of the inferential difficulties discussed in Section 2.3.

2.2 MCMC ALGORITHMS FOR SWITCHING REGRESSIONS

2.2.1 Gibbs Sampling

The Gibbs sampler is the most commonly used approach in mixture estimation (Diebolt and Robert 1994) and consists of augmenting the sample $(x_1, y_1), \dots, (x_n, y_n)$ from (1.2) with artificial allocation variables z_i , in order to “de-mix” the observed sample. This allows for simulation of the parameters of each component conditionally on the allocations, that is, taking into account only the observations which have been allocated to this component.

For instance, given a mixture of simple regressions,

$$y_i \sim \sum_{j=1}^k p_j \mathcal{N}(\beta_{0j} + \beta_{1j}x_i, \tau_j^2)$$

the conditional full distributions involved in this case are simply normal on both β_{0j} and β_{1j} and gamma on τ_j^{-2} .

For logistic regressions (1.3), a second level of augmentation is helpful to simulate from the full conditionals since, once the observations have been grouped by allocations, the k densities are of the form

$$\left\{ \prod_{i=1}^m \frac{\exp(y_i x_i' \beta)}{1 + \exp(x_i' \beta)} \right\} e^{-\|\beta\|^2 / 2\tau^2},$$

which is not directly tractable, but simulation from this distribution can be done via a *slice sampler* (see, e.g., Damien, Wakefield, and Walker 1999, or Robert and Casella 1999), leading to truncated normal proposals. In the case of the Poisson regressions (1.4), the posterior distribution is also untractable directly, being

$$\exp \left\{ \left(\sum_i y_i x_i \right)^t \beta - \sum_i e^{x_i^t \beta} \right\} e^{-\|\beta\|^2 / 2\tau^2}, \quad (2.1)$$

and a similar slice sampling solution applies.

2.2.2 Random Walk Metropolis–Hastings Algorithms

The intricacy of using the Gibbs sampler plus possibly a slice sampler, while natural, is not necessary. Indeed, a Metropolis–Hastings algorithm based on a properly calibrated Cauchy perturbation may be used for all types of regression models considered here (see Robert and Casella 1999, for details). For an unconstrained parameterization θ (i.e., after a log-transform on the variances/scale parameters, and a logit-transform on the weights), the proposal is $\tilde{\theta} = \theta^{(t)} + \omega \epsilon_t$, where ϵ_t is multivariate Cauchy and ω calibrated (via a burn-in step) to lead to an acceptance rate of about 0.4. For logistic models, the standard random walk proposal seems to mix faster than the slice sampler (see Altaleb 1999, for examples) and we observed the same phenomenon for the Poisson regression. In addition, the programming effort associated with the Metropolis–Hastings solution is much lower than with the Gibbs sampler: the same program applies for all types of regression, once the likelihood and prior are properly modified (Cappé and Robert 2000).

2.3 INFERENCEAL DIFFICULTIES

Figure 3(a) presents the Metropolis–Hastings output as a projection of the MCMC sample in the (β_0, β_1) plane for the regression coefficients of the simple linear regression $y = \beta_{0i} + x\beta_{1i}$ ($i = 1, 2$) and the ethanol dataset of Figure 2(a). The two clusters corresponding

to the components are clearly identified and lead to a good fit for the regression lines of the model (as they should, given the clear separation between the components in the data). Note: The graphs of this article are also available in color at the StatLib archive of *JCGS*.

Unfortunately this ideal situation is rarely encountered, and the general scenario is closer to the picture given in Figure 3(b), where the regression coefficients for the three components associated with the CO₂ dataset take values in the same area of the (β_0, β_1) space. Clearly, the three components are then undistinguishable as a whole (i.e., over the sequence of iterations), even though they can be separated at any given iteration.

This phenomenon is often called *label-switching* and has been known for a while, since it follows from the lack of identifiability of mixture models (see, e.g., Celeux, Hurn, and

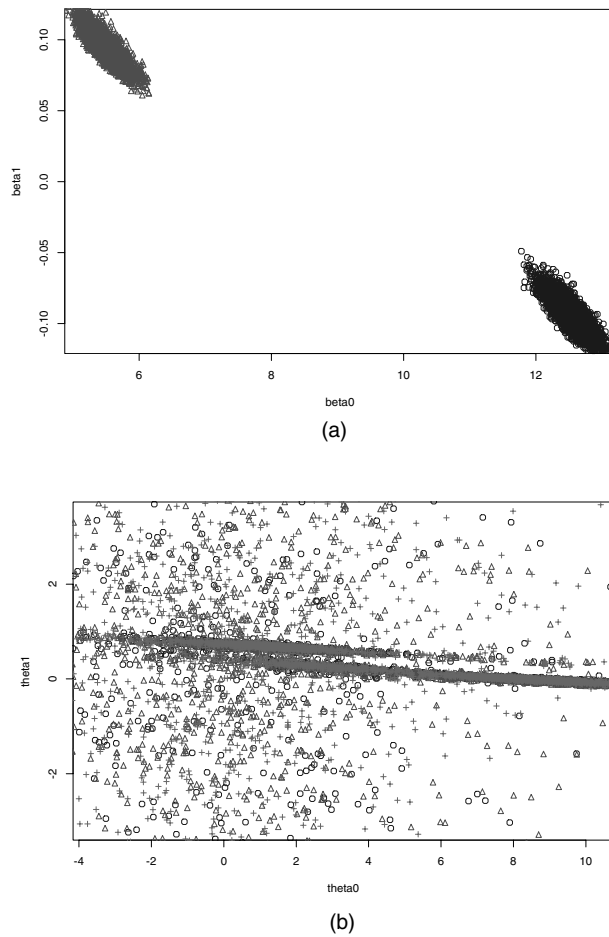


Figure 3. (a) Output of the random walk Metropolis–Hastings sampler projected on the (β_0, β_1) space for the ethanol dataset; (b) output of the Gibbs sampler projected on the (β_0, β_1) space for the CO₂ dataset and a three-component mixture model. (In both graphs, the MCMC output for each component of the mixture is represented by a different symbol and a different gray level.)

Robert 2000). In fact, as mentioned in Section 2.1, the posterior distribution is invariant under permutation of the labels of a mixture model. This means that a given component of a mixture cannot be individually extracted from the likelihood because the component labels i cannot be distinguished from one another, unless some identifiability constraint is imposed a priori [see the discussion of Richardson and Green (1997) for different perspectives on identifying constraints]. When the prior implies exchangeability over the parameters (β_i, p_i) or (β_i, σ_i, p_i) ($i = 1, \dots, k$), the symmetry of the posterior distribution in i , that is, in terms of the labeling of the parameters for a given number k of regression lines, implies that the corresponding marginal distributions of the (β_i, p_i) 's are all equal. Consequently the posterior marginal means $\mathbb{E}[\beta_i|x, y]$ are all equal and do not provide sensible estimators of the regression lines. Obviously, there are some settings where the different modes in the posterior distribution are so well separated that the MCMC sampler only visits one of the modes and label-switching does not occur, as shown in Figure 3(a). In such cases, where, *stricto sensus*, the MCMC sampler does not function correctly, the ergodic averages can provide acceptable answers, as shown by Figure 4(a) for the CO₂ dataset and a two-component mixture modeling.

Celeux et al. (2000) considered this problem for normal and exponential mixtures. They showed that the solution based on identifiability constraints, that is with plug-in estimates for the mixture parameters based on ordered MCMC samples, as in Richardson and Green (1997), does not always provide sensible solutions. This is also true for regression, even when the MCMC samples are well separated against the components. While Figure 4(a) gives the same answers for the unordered and ordered MCMC samplers, Figure 4(b), which corresponds to a simulated dataset of 185 observations from a four-component mixture, given on the upper left hand corner of the graph, shows that the two graphs corresponding to the orderings on β_0 and on β_1 differ from the standard averages, which are closer to the true regression lines. (See Figure 6(a) for a plot of the corresponding MCMC samples.)

2.4 INFORMATION CARRIED BY THE MCMC OUTPUT

The above statement of the problem does not imply that inference cannot be made in this case, since the MCMC output contains sufficient information on the regression lines and the parameters of the model. This is illustrated by Figure 5, which gives a sampling representation of the posterior distribution of the regression lines $y = \beta_{0i} + x\beta_{1i}$ ($i = 1, \dots, k$) for the two datasets of Section 1. This approximation is obtained the following way: at each iteration t of the MCMC sampler, five x 's are selected at random in the range $[x_{\min}, x_{\max}]$ and the points

$$\left(x, \beta_{0i}^{(t)} + x\beta_{1i}^{(t)}\right)$$

are represented on the graph, where the component i is chosen with probability $p_i^{(t)}$. A different color is associated with each component i and Figure 5(b) shows a strong mixing over colors, which means again that MCMC approximations to the posterior expectations $\mathbb{E}[\beta_i|x, y]$ will perform poorly. However, it also exhibits an equally strong stability in the location of the points, which means that the two main regression lines can be located almost

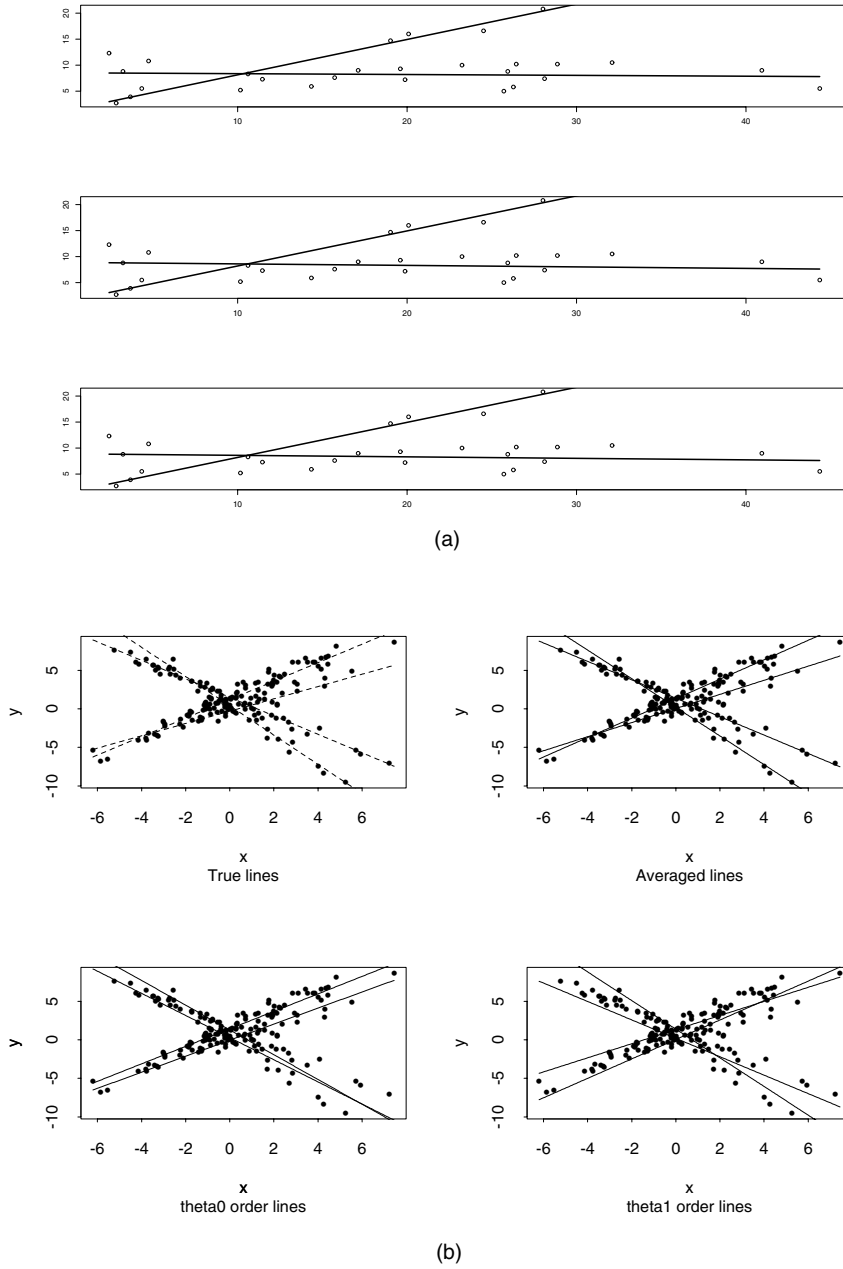


Figure 4. (a) Estimated regression lines for the CO_2 dataset with a two-component mixture model, based on a random walk Metropolis–Hastings algorithm: the upper graph is based on averages component by component, the middle graph on averages after reordering in β_0 and the bottom graph on averages after reordering in β_1 ; (b) same graph for a simulated dataset of 185 observations and a four-component mixture given on the upper left corner, for a random walk Metropolis–Hastings algorithm.

by sight and therefore that inference is still meaningful in this setting. In this regard there is little difference between the upper and lower graphs in Figure 5, even though the MCMC parameter samples are as dissimilar as possible (one being extremely mixed and the other not mixing at all over the labels).

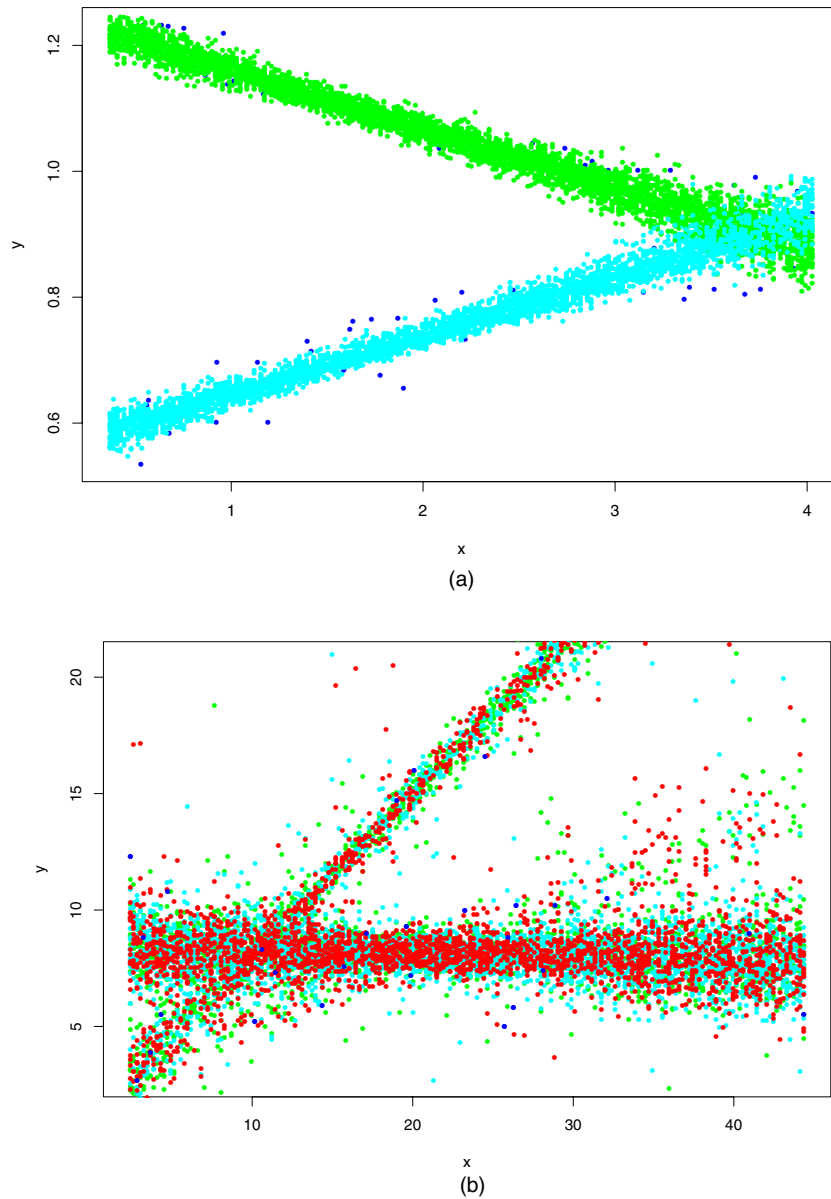


Figure 5. (a) Sampling representation of the posterior distribution on the regression lines $y = \beta_{0i} + x\beta_{1i}$ ($i = 1, 2$) for the ethanol dataset and a two-component mixture model, obtained by selecting points at random on these lines along iterations; (b) same graph for the CO₂ dataset and a three-component mixture model. (In each graph, the dark blue points represent the observations.)

One can notice that a natural byproduct of this sampling graphical device is to provide an evaluation of the variation of the regression lines and, hence, a Monte Carlo confidence band by selecting the 95% of points closest to the lines.

The inferential difficulty thus lies in formalizing the derivation of representative regressions lines from figures such as Figure 5(b). The main goal of this article is therefore to present solutions to this problem that are based on appropriate loss functions.

3. LOSS-BASED INFERENCE

Given the overall unreliability of the ad hoc average and of ordered estimates, we concentrate on the specification of loss functions $L(\xi, \hat{\xi})$ for the various inferential questions which arise in this setting. Following Celeux et al. (2000), these loss functions will not rely on the labeling of the components and so will not be affected by the lack of identifiability. Therefore, the corresponding Bayes estimator

$$\hat{\xi}^* = \arg \min_{\hat{\xi}} \mathbb{E}_{\xi|x,y} L(\xi, \hat{\xi}) \quad (3.1)$$

is equally unaffected by label indexes. [See Stephens (2000b) for a cluster-like approach in the mixture setting.]

For many choices of L , it is not possible to find $\hat{\xi}^*$ explicitly. However, for a large class of losses, it is computationally feasible to use the two-step approach of Rue (1995): The first step is to use MCMC ergodic averages to approximate $\mathbb{E}_{\xi|x,y} L(\xi, \hat{\xi})$ for a given $\hat{\xi}$. The second step is to perform the optimization of the estimated expected losses over $\hat{\xi}$. In order for this approach to be computationally viable, a practical restriction on the loss function $L(\xi, \hat{\xi})$ is that it is possible to separate the ξ and $\hat{\xi}$ terms in $L(\xi, \hat{\xi})$ in the sense that all evaluations of the estimated $\mathbb{E}_{\xi|x,y} L(\xi, \hat{\xi})$ with respect to different values of $\hat{\xi}$ can be performed using a single common MCMC sample from the posterior of $\xi|x$. We will restrict ourselves to loss functions for which this requirement holds, and explain how the estimate is obtained in each of the three estimation problems we consider.

3.1 ESTIMATING THE MIXTURE OF REGRESSION LINES

A natural choice for $L(\xi, \hat{\xi})$ is to use a loss function that considers the difference between two density functions, one using the parameters ξ and the other the parameters $\hat{\xi}$. The set of k regression lines defines a normal mixture density at any particular x value (as illustrated in Figure 1). Given that we do not make assumptions on the distribution of the regressors x , a natural solution advocated in Dupuis and Robert (2001) is to sum the distances between the conditional distributions, $d(f(\cdot|x, \xi), f(\cdot|x, \hat{\xi}))$, over the observed x values. For the density difference itself, the symmetrized Kullback–Leibler distance can be used

$$L(\xi, \hat{\xi}) = \sum_x \int_{\mathcal{R}} \left[f_{\xi}(y) \log \frac{f_{\xi}(y)}{f_{\hat{\xi}}(y)} + f_{\hat{\xi}}(y) \log \frac{f_{\hat{\xi}}(y)}{f_{\xi}(y)} \right] dy, \quad (3.2)$$

where f_ξ denotes the (conditional) density of the k -component regression mixture (the regressor variables x are omitted for ease of notation). Notice that this loss function does not involve the allocation variables which appear only in the completed density used by the Gibbs sampler. In the case of the qualitative regression models (1.3) and (1.4), the Kullback–Leibler distance (3.2) is modified as

$$L(\xi, \hat{\xi}) = \sum_x \sum_y \left[f_\xi(y) \log \frac{f_\xi(y)}{f_{\hat{\xi}}(y)} + f_{\hat{\xi}}(y) \log \frac{f_{\hat{\xi}}(y)}{f_\xi(y)} \right], \quad (3.3)$$

to account for the discrete nature of the observations y .

The choice of the Kullback–Leibler distance is defended for several reasons in Bernardo and Smith (1994), reasons ranging from information theory to asymptotics, including the fact that this loss function is invariant under reparameterization, unlike the L_2 distance for instance. An additional incentive for choosing this distance is that it is invariant under relabeling. Moreover, it lends itself to the two-step “estimation-then-maximization” approach in that we can decompose the posterior expected loss as follows:

$$\begin{aligned} & \mathbf{E}_{\xi|x,y} \left[\sum_x \int_{\mathcal{R}} \left(f_\xi(y) \log \frac{f_\xi(y)}{f_{\hat{\xi}}(y)} + f_{\hat{\xi}}(y) \log \frac{f_{\hat{\xi}}(y)}{f_\xi(y)} \right) dy \right] \\ &= \sum_x \int_{\mathcal{R}} \left(\mathbf{E}_{\xi|x,y} [f_\xi(y) \log f_\xi(y)] - \log f_{\hat{\xi}}(y) \mathbf{E}_{\xi|x,y} [f_\xi(y)] \right. \\ & \quad \left. + f_{\hat{\xi}}(y) \log f_{\hat{\xi}}(y) - f_{\hat{\xi}}(y) \mathbf{E}_{\xi|x,y} [\log f_\xi(y)] \right) dy, \end{aligned} \quad (3.4)$$

assuming that the order of integration may be interchanged. Considering the four terms on the right hand side of this equation individually, the first can be seen to be irrelevant for our purposes since it does not depend on $\hat{\xi}$. The second and fourth terms both involve posterior expectation, and we propose using the Markov chain Monte Carlo step to estimate these values on a series of y -grids (one for each of the x_i 's in the summation). We can then use numerical integration for the second, third, and fourth terms on these same grids of y 's to evaluate the estimated expected posterior loss (less a constant independent of $\hat{\xi}$) for any value of $\hat{\xi}$, that is, each evaluation requires n numerical integrations. The optimization step cannot be performed explicitly, and we have to resort to simulated annealing as in Celeux et al. (2000). The computational burden grows with the number of distinct observed x values but is feasible for all of our examples.

As a first illustration, Figure 6 depicts the MCMC sample corresponding to the simulated dataset of Figure 4(b), plus the comparison between the estimated and the true regression lines. As can be seen on Figure 6(b), the fit is particularly good, but Figure 6(a) shows that the MCMC clusters hardly overlap.

Moving on to the CO₂ dataset introduced in the Introduction, Figure 7(a) compares the estimators obtained when fitting a two-component and a three-component model to the data. The two-component model is more satisfactory, in the sense that it corresponds to the posterior spread of the regression lines in Figure 7(b), obtained by the sampling device explained in Section 1, while the estimated weight of the component with the negative slope is 0.007.

3.2 ESTIMATING CLUSTERS

There are applications where the primary information of interest is not the regression lines, but *clusters*, that is, the groupings of allocations of observations. Here again we need a loss function for which the arbitrariness of the component labels is immaterial. Recalling that the notation $z_i = l$ denotes that observation i is allocated to line l , we suggest using a function which compares configurations z and \hat{z} by noting whether pairs of observations

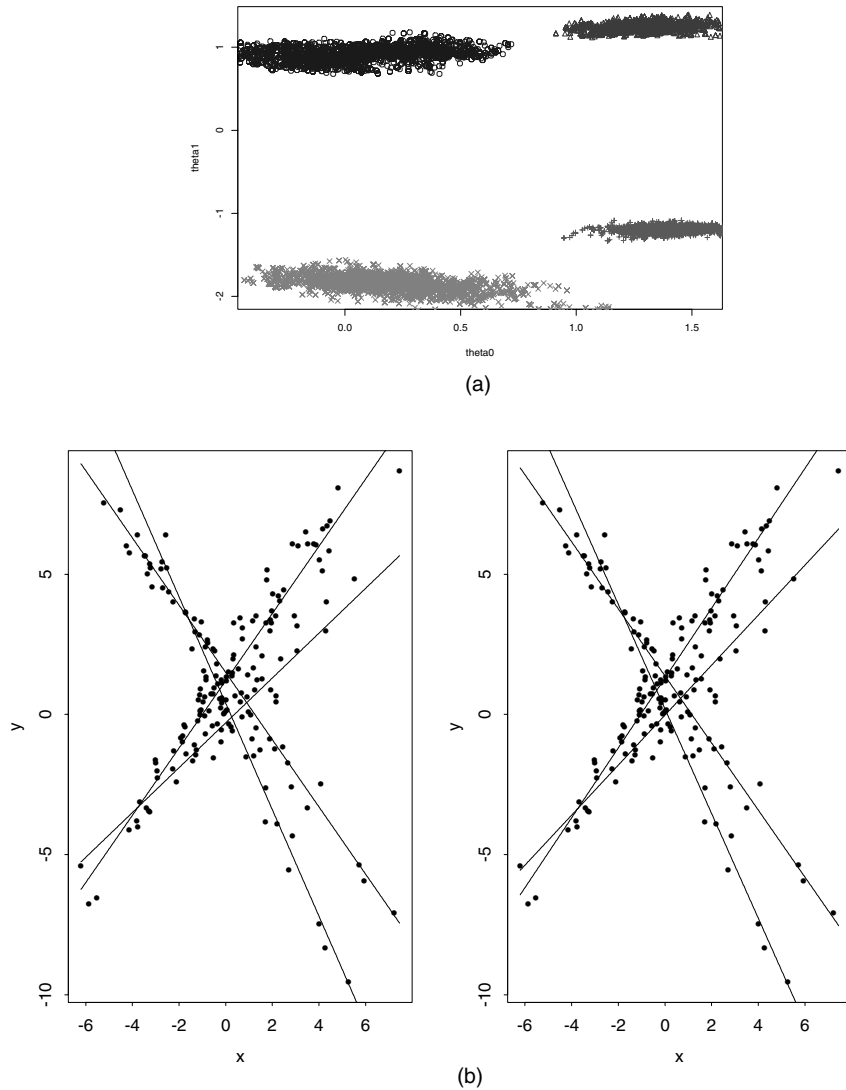
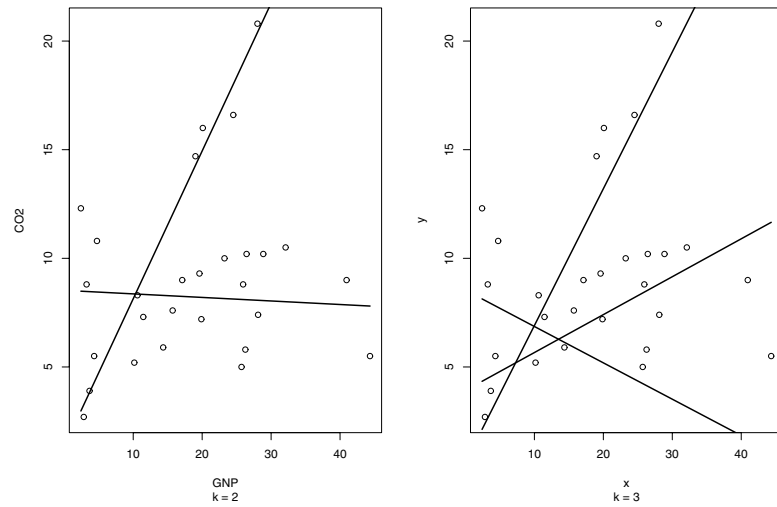


Figure 6. (a) Random walk MCMC sample associated with a simulated dataset of 185 observations in the (β_0, β_1) space for a four-component mixture model; (b) corresponding estimates of the regression lines obtained by the loss function (b) along with the true lines (left) (the dataset is represented on each plot).

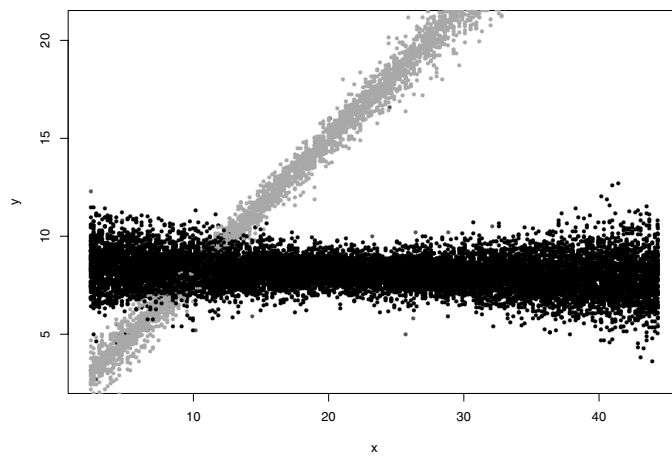
i and j are allocated to the same component under one set of allocations but to different components under the other set:

$$L(z, \hat{z}) = \sum_{i < j} \left\{ \mathbb{I}_{[z_i = z_j]} (1 - \mathbb{I}_{[\hat{z}_i = \hat{z}_j]}) + \mathbb{I}_{[\hat{z}_i = \hat{z}_j]} (1 - \mathbb{I}_{[z_i = z_j]}) \right\}. \quad (3.5)$$

Since the posterior expectation of $\mathbb{I}_{[z_i = z_j]}$ is the posterior probability of this event,



(a)

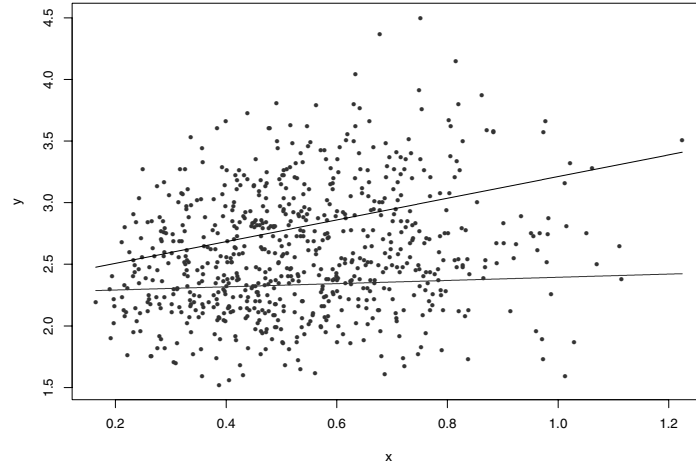


(b)

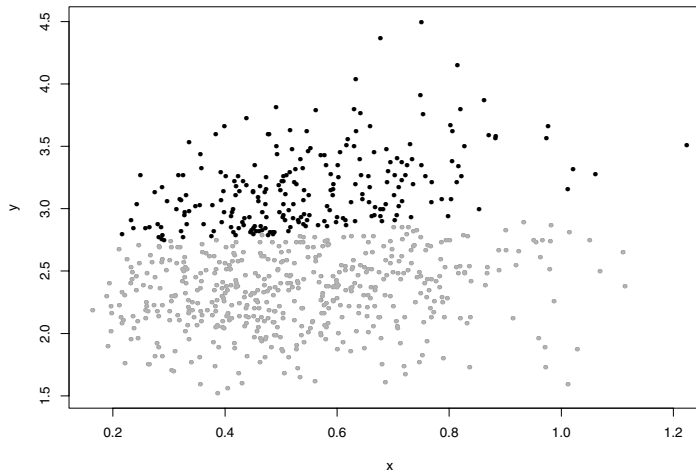
Figure 7. Comparisons of the loss estimates for two- and three-component mixtures for the CO₂ dataset and a random walk MCMC sample (left and right, respectively); (b) repartition of the derived sample of the regression lines for the two-component mixture model.

$$\mathbb{E}_{z|x,y}L(z, \hat{z}) = \sum_{i < j} \{P_{z|x,y}(z_i = z_j)(1 - \mathbb{I}_{[\hat{z}_i = \hat{z}_j]}) + \mathbb{I}_{[\hat{z}_i = \hat{z}_j]}P_{z|x,y}(z_i \neq z_j)\} . \quad (3.6)$$

The z_i 's are automatically generated as part of the Gibbs sampler steps, but it is straightforward to add an extra data augmentation step to the other samplers, without jeopardizing the stability properties of both the MCMC and the z_i chains. Notice moreover that this data augmentation step is fairly generic and hence costless, as opposed to the Gibbs parameter



(a)



(b)

Figure 8. (a) Dataset and estimated regression lines for the dataset of Fan and Chen (1999) for a two-component model via the loss function of Section 3.1; (b) representation of the allocations of the observations to components 1 (dark) and 2 (gray) via the loss function of Section 3.2.

step, which always involves specific programming. (All the C programs used in this article are available upon request from the authors.)

Consider, as an illustration of both loss functions presented so far, a dataset from Fan and Chen (1999), which relates daily measurements of the nitrogen pollutant NO₂ in Hong Kong in 1994–1995 and the number of daily hospital admissions for circulation and respiration problems. For two components, the MCMC clusters are clearly separated and the posterior distribution of the regression lines gives two distinct bands. The loss analysis of the MCMC sample is in agreement with this empirical analysis: the estimated lines in Figure 8(a) correspond to these two bands, while the dataset is neatly separated into two groups, as shown in Figure 8(b), with a division which is close to the average of the two regression lines. Notice that the estimated weights of the two components are 0.553 and 0.447, respectively. (The fact that one regression line is almost flat is quite interesting in that it suggests there are days when the NO₂ level does not directly affect hospital admissions.)

3.3 ESTIMATING THE LINES PLUS ALLOCATED OBSERVATIONS

A more likely scenario is that the question of interest is actually to identify the lines to which groups of observations belong. By this we mean that we wish to state which subsets of the observations are associated with a regression line with parameters β, σ^2 . Notice that in this scenario, there is a redundancy in estimating both the allocations for each observation and the weights associated with each line. We therefore consider only the allocations and the parameters of the different regressions. The parameter of interest is thus η defined to be the set of k $\{\beta, \sigma^2\}$'s, together with (z_1, \dots, z_n) .

We again propose a form of Kullback–Leibler distance as a loss function, but now considering the allocation of observation (x_i, y_i) , say $z_i = l$, the density in question at x_i is the regression with parameters $\{\beta_l, \sigma_l^2\}$ rather than the mixture regression. The difference between these densities is illustrated by comparing Figures 1 and 9, the former showing the full regression mixture at each x_i , the latter the normal densities corresponding to the z_i -th regression at each labeled x_i . The loss function is therefore

$$L(\eta, \hat{\eta}) = \sum_x \int_{\mathcal{R}} \left[f_{\eta}(y) \log \frac{f_{\eta}(y)}{f_{\hat{\eta}}(y)} + f_{\hat{\eta}}(y) \log \frac{f_{\hat{\eta}}(y)}{f_{\eta}(y)} \right] dy, \quad (3.7)$$

where f_{η} now denotes the density of the regression line corresponding to the allocation of the corresponding x . In the vocabulary of the EM algorithm, we thus work with the completed likelihood rather than with the observed likelihood. We employ the same Markov chain Monte Carlo ideas as in Section 3.1 to estimate expected posterior values on a grid, and numerical integration and simulation for the optimization step.

4. UNKNOWN NUMBERS OF COMPONENTS

When the degree of heterogeneity in the data is unknown, as is the case for the CO₂ dataset of Section 1, it is unreasonable to postulate a fixed number k of components in the

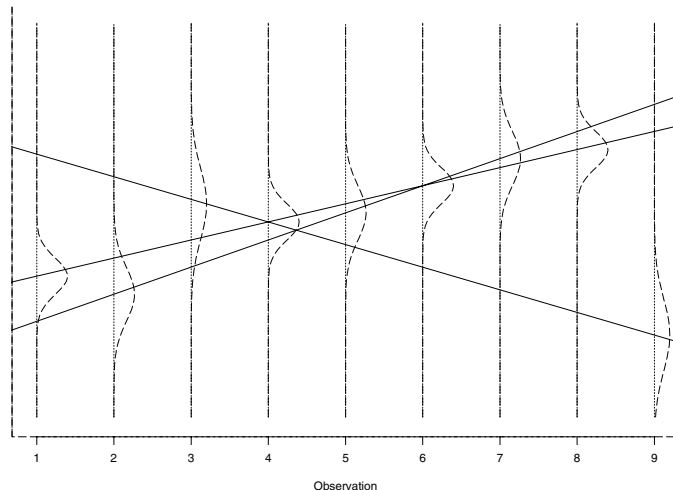


Figure 9. The three-component labeled regression mixture conditional distributions of $y|x, z$.

mixtures (1.2), (1.3), and (1.4), and k must be estimated as well. From a Bayesian point of view, this implies using a prior distribution on k : for illustration purposes as well as simplicity's sake, we choose a Poisson $\mathcal{P}(\lambda_1)$ distribution for this prior with $\lambda_1 = 2$ in most of the following examples. While this complex estimation problem has been addressed for standard mixtures through reversible jump techniques (Richardson and Green 1997), we opt for the recent alternative proposed by Stephens (2000a) and based on birth-and-death processes. Both approaches are equally valid on theoretical grounds (since they both satisfy the required detailed balance conditions) and our choice is based solely on practical considerations; the birth-and-death process avoids the specification of well-calibrated moves, such as the splits-and-merges of Richardson and Green (1997) and the corresponding computation of the Jacobians of the one-to-one transforms, plus it allows for easier changes of both prior distributions and parameterization. Another relevant point is that the implementation of the birth-and-death procedure is completely independent of the MCMC algorithm chosen for the fixed k steps. In fact, it works even without a supplementary MCMC algorithm. This property allows very modular programming when implementing the method and thus improves portability. [Notice, however, that the reversible jump algorithm, when restricted to birth and death moves with birth proposals based on the prior distribution, enjoys similar properties. See Cappé, Robert, and Rydén (in press), for further connections between the approaches.]

4.1 BIRTH-AND-DEATH PROCESSES

We will not go into details about the birth-and-death procedure since this is a straightforward adaptation of Stephens' (2000a) algorithm. Figure 10 gives a description of the method in an program-like manner. The method depends on a single parameter λ_0 , which

For $v = 0, 1, \dots, V$
 $t \leftarrow v$

Run till $t > v + 1$

1. Compute $\delta_j(\beta) = \frac{L(y|\beta_j)}{L(y)} \frac{\lambda_0}{\lambda_1}$
2. $\delta(\beta) \leftarrow \sum_{j=1}^k \delta_j(\beta)$, $\mu \leftarrow \lambda_0 + \delta(\beta)$, $u \sim \mathcal{U}([0, 1])$
3. $t \leftarrow t - u \log(u)$
4. With probability $\delta(\beta)/\mu$
 - Remove component j with probability $\delta_j(\beta)/\delta(\beta)$
 - $k \leftarrow k - 1$
 - $p_\ell \leftarrow p_\ell / (1 - p_j)$ ($\ell \neq j$)

Otherwise,

- Add component j from the prior $\pi(\beta)$
- $p_j \sim \mathcal{B}e(\gamma, k\gamma)$
- $p_\ell \leftarrow p_\ell (1 - p_j)$ ($\ell \neq j$)
- $k \leftarrow k + 1$

5. Run I MCMC(k, β, p)

Figure 10. Algorithmic representation of the birth-and-death process. The parameter λ_1 is the mean of the Poisson prior on k and γ is the parameter of the Dirichlet prior on p . The MCMC(k, β, p) function is a generic MCMC algorithm for a fixed k .

is the birth rate in the point process. We found in practice that better mixing was associated with higher values of λ_0 , as should be expected, but this is at the expense of running more simulations between two monitoring (integer) instants. The influence of the number of iterations I where an MCMC sampler is run for the current value of k (Step 5) is small and we chose $I = 10$ in our experiments.

The algorithms perform quite satisfactorily in our case and produces output where label switching is naturally very high (since the meaning of component 1, say, is vacuous, given that the number of components varies from one iteration to one another). The inferential issues addressed in Section 3 are thus reproduced in the current setup and we study the implementation of the above loss functions in the case of varying k .

Notice first that, from an inferential point of view, there is little difference from the fixed k setting. Indeed, the inference must be conditional on k , as argued by Robert (1997), because a general principle in Bayesian model choice is that parameters appearing in different models must be considered as separate entities. Therefore, the analysis for each value of k is as in the fixed k case. Obviously, running the analysis for varying k also leads to inference about the value of k supported by the data, both through the posterior probabilities and through the predictive plots of the regression lines, approximated as in the fixed k case.

For instance, the most probable value of k for the simulated dataset of 185 observations is $k = 4$.

4.2 NORMAL REGRESSION

We reanalyzed the CO₂ dataset, with the conclusion shown in Figure 11 that $k = 2$ is the solution most favored by the data (using a Poisson prior on k with mean 2). It is

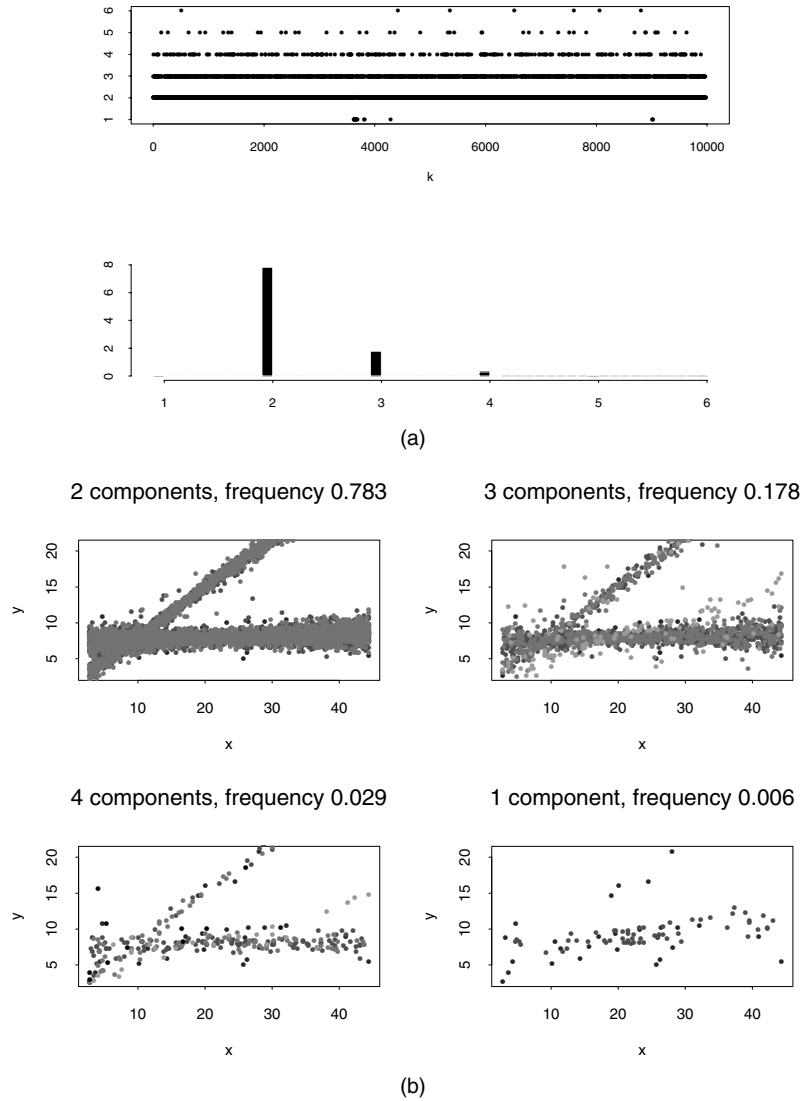


Figure 11. (a) Rawplot of the successive values of k along 10,000 consecutive iterations of the birth-and-death algorithm and corresponding histogram for the CO₂ dataset; (b) representation of the posterior distributions of the regression lines for the four most likely values of k .

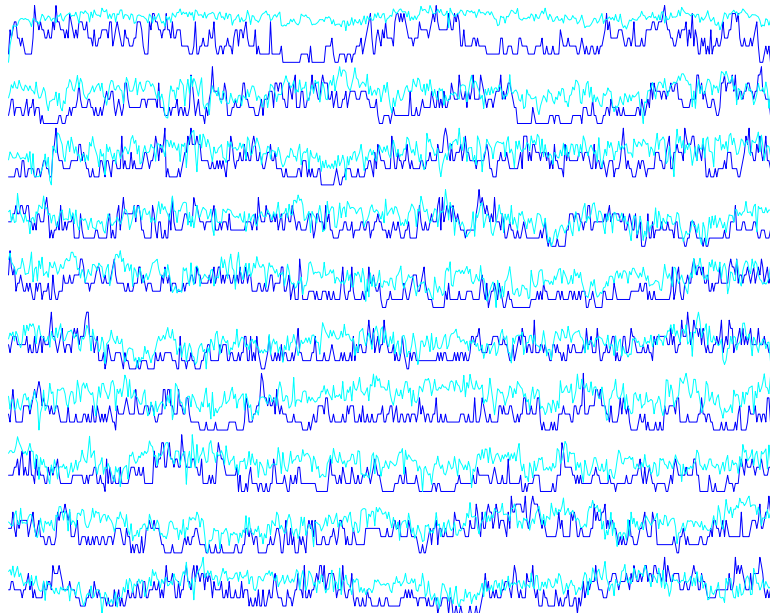


Figure 12. Plot of the successive values of k along 5,000 consecutive iterations of the birth-and-death algorithm for the INSEE dataset, with in superposition the corresponding value of the log-likelihood.

interesting to see that the birth-and-death process creates enough mixing to ensure that k moves at a good rate. Notice that, in the case $k = 3$, there is a stable third line appearing, which is consistent with the findings of Figure 7 when k is fixed. However, the weight associated with the third regression line is estimated by 0.0069 using the loss of Section 3 and is thus negligible. (From Figure 11(a), it would seem that $k = 5$ is more likely than $k = 1$, but this is an artifact, due to the fact that the chain does not change for a large number of iterations in the few cases $k = 1$.)

4.3 LOGISTIC REGRESSION

A mixture of logistic regression models (1.3) is appropriate for the analysis of a dataset from the INSEE (Institut National de la Statistique et des Etudes Economiques) 1992 survey on financial assets. This dataset of 1,278 individuals relates the possession of certain financial assets (SICAV) to some socioeconomic covariates, from which we selected the inactivity dummy variable (1 if inactive, 0 if active), the age group (from 1 to 6), the yearly income, and the overall wealth (both variables have been standardized). The prior on the β_{ij} 's is a normal $\mathcal{N}(0, 100)$ distribution.

The trace of the output of the Gibbs sampler (not represented here) shows that the label-switching phenomenon occurs quite intensely. When k varies, the mixing over k induced by the birth-and-death algorithm is quite high, as shown by Figure 12, which also contains the sequence of the log-likelihoods at the successive values of the parameters. The shapes of the posterior distributions of the β_{ij} 's, when averaging over k , are very similar to those

Table 1. Loss Estimates of the Logistic Coefficients for the INSEE Dataset Produced by the Procedure of Section 3.1 for the Most Probable Values of k (the numbers of occurrences of the corresponding k 's are indicated in parentheses, out of a total number of 5,000 simulations)

k	p_i	β_{0i}	β_{1i}	β_{2i}	β_{3i}
2	0.750707	4.947922	12.030713	-7.071619	4.142905
(1204)	0.249293	15.330934	-17.102980	4.523014	-1.444494
3	0.713159	6.714042	13.041261	-5.963821	3.590483
(1294)	0.214304	-5.907764	-15.258550	-1.879069	2.081183
	0.072536	1.926905	0.032500	1.811353	0.305916
4	0.510735	5.131563	13.833825	-2.484054	8.002322
(1064)	0.256141	10.304051	3.565791	2.601400	7.830500
	0.195479	-22.490153	-7.773485	-4.394919	-16.251573
	0.037574	0.768561	-24.482456	-8.485442	-1.261904

obtained for $k = 4$. Notice also how the log-likelihood follows more accurately the shape of the k curve in the last part of the iterations, which may indicate a better fit for all values of k due to a stabilization of the process. Table 1 provides the estimates of the logistic coefficient obtained via the loss function of Section 3.1. It shows a persistence of the first component for the three values of k , with a negative effect of the yearly income which can be explained as a tendency for people with higher incomes to choose more risky but more profitable financial assets. A similar analysis can be made for the second component when $k = 3$ and the third group when $k = 4$: younger active people with higher incomes are less likely to own the SICAV assets. Notice also that the last component for $k = 4$ has a fairly low probability.

4.4 POISSON REGRESSION

When considering the Poisson regression mixture (1.4), the setting is very similar to the previous paragraph. We first check the performances of the birth-and-death algorithm on a dataset analyzed by Lenk (1999), which relates the monthly unemployment rate with the monthly number of accidents (in thousands) in Michigan, from 1978 to 1987. The model is thus

$$N|\varrho \sim \sum_i p_i \mathcal{P}(\exp\{\beta_{0i} + \beta_{1i} \log(\varrho)\}),$$

where N denotes the number of accidents and ϱ the corresponding unemployment rate. Since Lenk (1999) noted a possible seasonal effect on the dependence, we postulated a $\mathcal{P}(4)$ prior on the number of components.

The result of this study is that the series of the $k^{(t)}$'s produced by the birth-and-death algorithm is mostly composed of 1's, with a frequency approximately of 0.91, and of 2's, with a frequency approximately of 0.09. Examination of the log-likelihood also shows that the Markov chain does not remain in the neighborhood of a local maximum but explores (to a certain extent) the posterior surface. When $k = 1$, the standard posterior mean is acceptable as estimate and the "plug-in" estimate of the regression function provides a reasonable fit

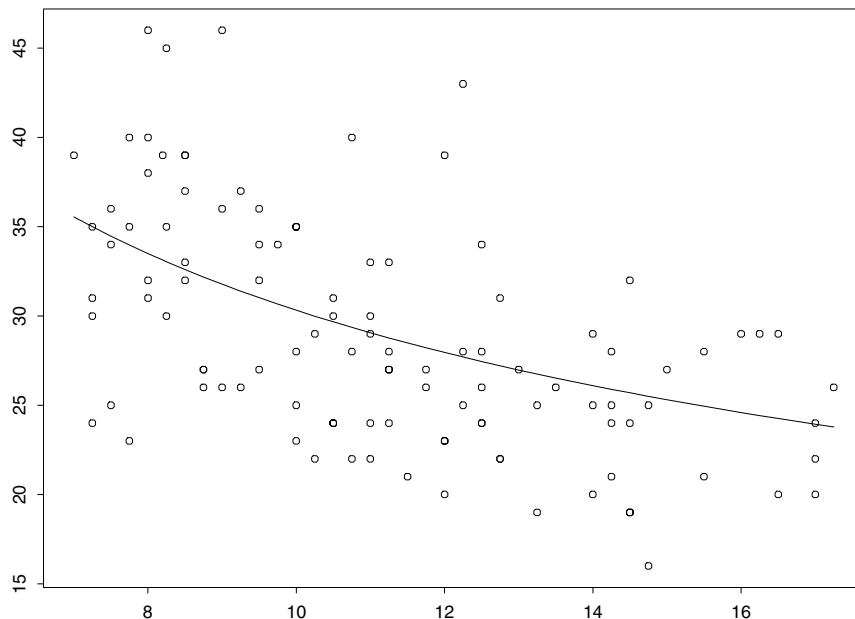


Figure 13. Regression curve with plug-in loss estimate for the Poisson modeling of the Michigan accident dataset.

of the dataset which cannot be distinguished from the loss estimate derived from Section 3.1 in Figure 13. For $k = 2$, the loss estimate produces a second component with a very small weight, 0.001399, which can be ignored.

We now turn to a second accident dataset, provided by INRETS (Institut National de Recherche sur les Transports et leur Sécurité), which relates the number of accidents at 107 selected crossroads with the average traffic on both the primary (t_p) and the secondary (t_s) roads connecting through these crossroads (Brenac 1994). The expectation of the Poisson distribution is then

$$\exp\{\beta_0 + \beta_1 \log(t_s) + \beta_2 \log(t_p)\}. \quad (4.1)$$

As in Viallefont, Richardson, and Green (2002), whose analysis uses reversible jump techniques, the most probable value of k is 3, although $k = 4$ comes a close second. However, the MCMC sample conditional on $k = 4$ accumulates values near 0 and this case does not lead to higher values of the likelihood. Moreover, for the three-component case, the third component corresponds to the high-risk crossroads, with an increasing risk factor (or *elasticity*) β_1 due to the secondary road and a decreasing elasticity for the primary road, while, for the four-component case, the elasticities are too high to be realistic from a traffic safety point of view.

Figure 14(a) gives a contour plot of the posterior expectations of the mixtures of the averages (4.1), which do not depend on the labeling and is obtained by averaging over the MCMC iterations, producing the expected increase of the number of accidents in both t_p and t_s , while Figure 14(b) inserts a barplot average predictive distribution on the number of

accidents for each observed couple (t_s, t_p) with, logically, a wider spread for larger values of t_s and t_p . This figure convincingly illustrates the potential of the MCMC sample in providing various interpretations of the estimated model.

As detailed in the technical report associated with this article, it is also possible to use the MCMC output to represent the loss estimates given in Table 2 through contour lines derived from the plug-in expectations (4.1) and color allocations of the 107 observations, for different values of k .

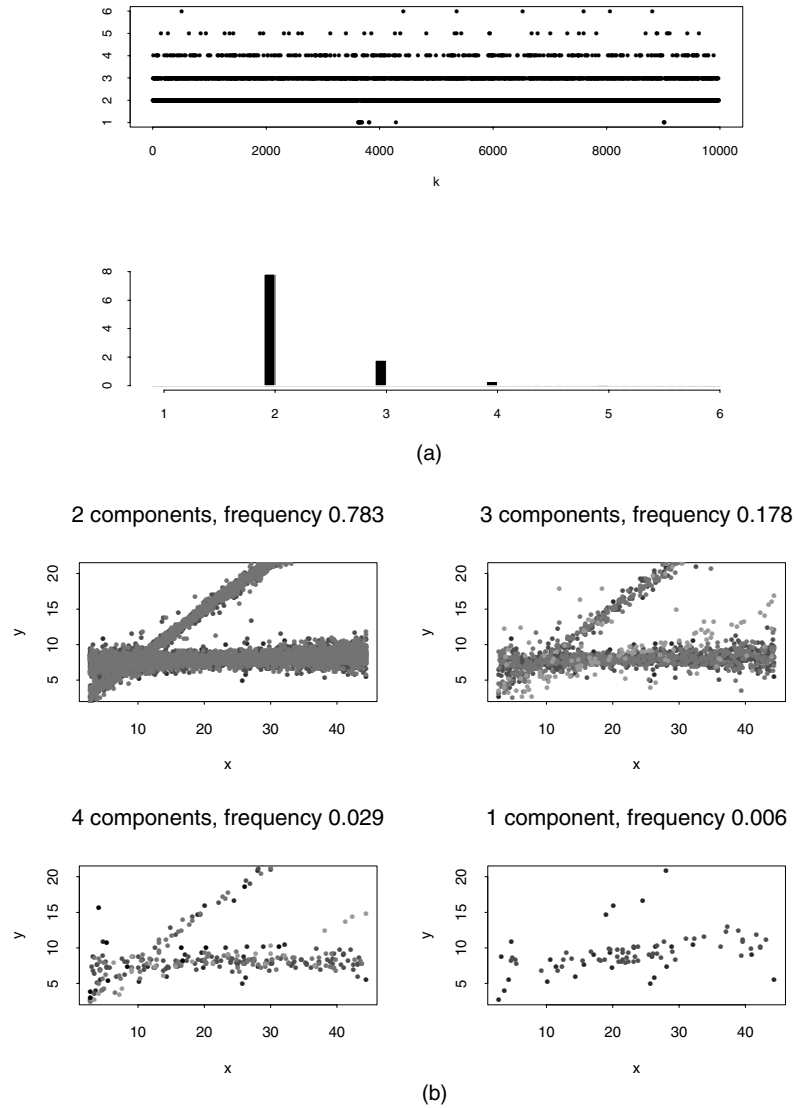


Figure 14. (a) Contour plot of the posterior expectations $\exp(\beta_0 + \beta_1 \log(t_s) + \beta_2 \log(t_p))$ for the INRETS accident dataset; (b) corresponding barplot of the average predictive distribution for each observed crossroad (t_s, t_p) .

Table 2. Estimated parameters of model (4.1) under the loss function (3.7), for the two most probable values of k and the INRETS accident dataset (the numbers of occurrences of the corresponding k 's are indicated in parentheses, out of a total number of 5,000 simulations)

<i>Number of components</i>	β_0	β_1	β_2
3 (2530)	0.271051	0.248952	1.178467
	1.619090	0.223184	0.405672
	3.025534	0.430442	-0.154674
4 (1856)	-0.283464	0.452226	1.809185
	1.019317	0.405510	0.617267
	1.425445	0.269155	2.783191
	1.550998	2.567553	-2.578583

5. CONCLUSION

Although decision theory is often criticized on a practical basis, on the ground that actual decision makers cannot build realistic loss functions, we have demonstrated in this article that formal loss functions based on the Kullback–Leibler functional distance can produce valuable answers in the nonidentifiable setting of mixtures of regressions, where standard empirical measures such as the posterior mean do not work. While the simulation steps used here to produce a sample from the posterior distribution are standard (even though we are not aware of previous uses of Metropolis–Hastings procedures in such settings), the computational part invoked by the minimization of the posterior loss is more intense, but feasible.

We have also shown that inference can be drawn about the number of components in the mixture, coupling the simple procedure of Stephens (2000a) with the MCMC algorithms developed here, and that, besides, the loss derivation of estimates also applies in this variable dimension setting by conditioning on k . Examples demonstrated that the mixing behavior of the samplers was satisfactory, even for several regressors. These various tools, as well as the graphical representations introduced in this article, can be used for a wide variety of regression problems and of prior distributions, and this at a minimal programming cost. This study indicated in addition that the prior distribution on k is not without effect on the resulting inference: a Poisson $\mathcal{P}(1)$ prior does provide a posterior distribution on k that differs widely from the posterior associated with a Poisson $\mathcal{P}(10)$ prior. This feature is to be expected, given the ill-posed nature of the mixture models. It has, however, a limited influence on the predictive distribution, in the sense that curves such as those in Figure 14 do not vary between priors.

ACKNOWLEDGMENTS

The authors are thankful to the participants to the second TMR workshop in Crete, May 1999, for their comments and in particular to Arnaldo Frigessi (Oslo) for his support. The INSEE dataset was kindly provided by Denis Fougères (CREST) and the INRETS dataset by Valérie Viallefont (CREST), to whom we are extremely

grateful, as well as to T. Brenac and S. Lassare (INRETS), who gave us permission to use this second dataset and commented on a previous version. Discussions with M.T. Wells (Cornell), and suggestions from the associate editor and two referees, were equally helpful. Ana Justel was partially supported by CREST during a visit in May–June 1999 and by DGES Grant BFM2001-0169. Merrilee Hurn and Christian P. Robert were supported by EU TMR network ERB–FMRX–CT96–0095 on “Computational and Statistical Methods for the Analysis of Spatial Data.”

[Received September 2000. Revised June 2001.]

REFERENCES

- Altaleb, A. (1999), “Contributions à la Théorie des Algorithmes MCMC,” unpublished Ph.D. thesis, Université de Rouen.
- Bar-Shalom, Y. (1978), “Tracking Methods in a Multi-Target Environment,” *IEEE Transactions on Automatic Control*, 23, 618–626.
- Bernardo, J. M., and Smith, A. F. M. (1994), *Bayesian Theory*, New York: Wiley.
- Billio, M., Monfort, A., and Robert, C. P. (1999), “Bayesian Estimation of Switching ARMA Models,” *Journal of Econometrics*, 93, 229–255.
- Brenac, T. (1994), “Accidents en Carrefour sur Routes Nationales, Modélisation du Nombre d’Accidents Prédictible sur un Carrefour et Exemples d’Applications,” Rapport 185, INRETS, Arcueil.
- Cappé, O., and Robert, C. P. (2000), “Markov Chain Monte Carlo: Ten Years and Still Running!” *Journal of the American Statistical Association*, 95, 1282–1286.
- Cappé, O., Robert, C. P., and Rydén, T. (in press), “Reversible Jump, Birth-and-Death, and More General Continuous Time MCMC Samplers,” *Journal of the Royal Statistical Society, Ser. B*.
- Celeux, G., Hurn, M., and Robert, C. P. (2000) “Computational and Inferential Difficulties with Mixture Posterior Distributions,” *Journal of the American Statistical Association*, 95, 957–970.
- Damien, P., Wakefield, J., and Walker, S. (1999), “Gibbs Sampling for Bayesian Non-conjugate and Hierarchical Models by Using Auxiliary Variables,” *Journal of the Royal Statistical Society, Ser. B*, 61, 331–344.
- Diebolt, J., and Robert, C. P. (1994), “Estimation of Finite Mixture Distributions by Bayesian Sampling,” *Journal of the Royal Statistical Society, Ser. B*, 56, 363–375.
- Dupuis, J., and Robert, C. P. (2001), “Bayesian Variable Selection in Qualitative Models by Kullback–Leibler Projections,” *Journal of Statistical Planning and Inference*, 111, 77–94.
- Fan, J., and Chen, J. (1999), “One-Step Local Quasi-Likelihood Estimation,” *Journal of the Royal Statistical Society, Ser. B*, 61, 927–943.
- Goldfeld, S. M., and Quandt, R. E. (1976), “A Markov Model for Switching Regression,” *Journal of Econometrics*, 1, 3–16.
- Hamilton, J. D. (1989), “A New Approach to the Economic Analysis of Nonstationary Time Series and the Business Cycle,” *Econometrica*, 57, 357–384.
- Hurvich, C. M., Simonoff, J. S., and Tsai, C.-L. (1998), “Smoothing Parameter Selection in Nonparametric Regression Using an Improved Akaike Information Criterion,” *Journal of the Royal Statistical Society, Ser. B*, 60, 271–294.
- Lenk, P. (1999), “Bayesian Inference for Semiparametric Regression Using a Fourier Representation,” *Journal of the Royal Statistical Society, Ser. B*, 61, 863–879.
- Kiefer, N. (1980), “A Note on Switching Regression and Logistic Discrimination,” *Econometrica*, 48, 1065–1069.
- Quandt, R. E., and Ramsey, J. B. (1978), “Estimating Mixtures of Normal Distributions and Switching Regressions” (with discussion), *Journal of the American Statistical Association*, 73, 730–752.

- Richardson, S., and Green, P. J. (1997), "On Bayesian Analysis of Mixtures With an Unknown Number of Components" (with discussion), *Journal of the Royal Statistical Society, Ser. B*, 59, 731–792.
- Robert, C. P. (1996), "Inference in Mixture Models," in *Markov Chain Monte Carlo in Practice*, W. R. Gilks, S. Richardson, and D. J. Spiegelhalter (eds.), London: Chapman and Hall, pp. 441–464.
- (1997), Discussion of "On Bayesian Analysis of Mixtures with an unknown Number of Components" by S. Richardson and P. Green, *Journal of the Royal Statistical Society, Ser. B*, 59, 758–764.
- Robert, C. P., and Casella, G. (1999), *Monte Carlo Statistical Methods*, New York: Springer-Verlag.
- Rue, H. (1995), "New Loss Functions in Bayesian Imaging," *Journal of the American Statistical Association*, 90, 900–908.
- Shumway, R. H., and Stoffer, D. S. (1991), "Dynamic Linear Models with Switching," *Journal of the American Statistical Association*, 81, 763–769.
- Stephens, M. (2000a), "Bayesian Methods for Mixtures of Normal Distributions—An Alternative to Reversible Jump Methods," *The Annals of Statistics*, 28, 40–74.
- (2000b), "Dealing With Label Switching in Mixture Models," *Journal of the Royal Statistical Society, Ser. B*, 62, 795–810.
- Viallefont, V., Richardson, S., and Green, P. J. (2002), "Bayesian Estimation of Poisson Mixtures," *Journal of Nonparametric Statistics*, 14, 181–202.