

# Hidden Markov Models and Disease Mapping

Peter J. GREEN and Sylvia RICHARDSON

---

We present new methodology to extend hidden Markov models to the spatial domain, and use this class of models to analyze spatial heterogeneity of count data on a rare phenomenon. This situation occurs commonly in many domains of application, particularly in disease mapping. We assume that the counts follow a Poisson model at the lowest level of the hierarchy, and introduce a finite-mixture model for the Poisson rates at the next level. The novelty lies in the model for allocation to the mixture components, which follows a spatially correlated process, the Potts model, and in treating the number of components of the spatial mixture as unknown. Inference is performed in a Bayesian framework using reversible jump Markov chain Monte Carlo. The model introduced can be viewed as a Bayesian semiparametric approach to specifying flexible spatial distribution in hierarchical models. Performance of the model and comparison with an alternative well-known Markov random field specification for the Poisson rates are demonstrated on synthetic datasets. We show that our allocation model avoids the problem of oversmoothing in cases where the underlying rates exhibit discontinuities, while giving equally good results in cases of smooth gradient-like or highly autocorrelated rates. The methodology is illustrated on an epidemiologic application to data on a rare cancer in France.

**KEY WORDS:** Allocation; Bayesian hierarchical model; Disease mapping; Finite mixture distributions; Heterogeneity; Hidden Markov models; Markov chain Monte Carlo; Poisson mixtures; Potts model; Reversible jump algorithms; Semiparametric model; Spatial mixtures; Split/merge moves.

---

## 1. INTRODUCTION

### 1.1 Hidden Markov Random Fields

Hidden Markov models (HMMs) assume in general terms that the observations form a noisy realization of an underlying process that has a simple structure with Markovian dependence. The most studied class of HMMs has been temporal observations linked to an underlying Markov chain. This formulation originated in engineering and has since been used in many domains, ranging from finance to molecular biology. (A comprehensive review of mathematical properties, statistical treatment, and applications of hidden Markov chains can found in Künsch 2001.) Robert, Rydén, and Titterton (2000) reported a recent study that is particularly relevant to the work presented here because of the variable number of states in the hidden chain, the Bayesian treatment, and the use of reversible jump algorithms in the implementation.

When the data are spatially structured, a natural extension is to hidden Markov random fields, that is, Markov random fields degraded by (conditionally) independent noise. One context in which such models have been much used is image analysis, going back to Besag (1986) and beyond; another is disease mapping in epidemiology. By disease mapping, we mean studies aiming to uncover a potential spatial structure in disease risk when analyzing small numbers of observed health events in a predefined set of areas. In this case the noise is related to the rarity of the health event and the size of the population at risk, leading to the low disease counts per area ( $<10$ ) typically found in many studies (e.g., Elliott, Wakefield, Best, and Briggs 2000).

In this article we consider a class of hidden discrete-state Markov random field models related to an underlying

finite-mixture model that allows spatial dependence and does not predefine the number of “states” or components of the mixture. Our motivation is dual: We are interested in generalizing some of the useful features of hidden Markov chains with an unknown number of discrete states and also in proposing a flexible alternative to the current Markov random field models commonly used in disease mapping. This motivating context has both stimulated model development and driven some of the choices made in model implementation.

Although the extension from a hidden Markov chain to a hidden random field is an obvious one in modeling terms, it introduces disproportionate difficulties in implementation. In the case of a linear chain, there are fast methods such as the “forward-backwards” algorithm for computing likelihoods, and if a full Bayesian approach via Markov chain Monte Carlo (MCMC) is taken, then the normalizing constants of the joint prior distribution of the hidden chain are explicitly available. Neither of these is true of the general spatial case, and this has severely limited application of these methods (see Rydén and Titterton 1998 for a full discussion). One of the aims of this article is to demonstrate that implementation is perfectly practicable, at least for a moderate number of hidden variables, when using a model with just two parameters (the number of hidden states and the strength of interaction).

### 1.2 Models for Spatially Correlated Count Data

We consider the modeling of spatial heterogeneity for count data on a rare phenomenon, observed in a predefined set of areas. Throughout, our terminology refers to disease mapping, but we stress that our model is easily translatable to other contexts in which spatial heterogeneity is of interest, for example, in ecology or agricultural science.

There are many reasons for suspecting heterogeneity in an underlying disease event rate and wanting to characterize it. For example, the discovery of either local discontinuities

---

Peter J. Green is Professor of Statistics, Department of Mathematics, University of Bristol, Bristol BS8 1TW, U.K. (E-mail: [P.J.Green@bristol.ac.uk](mailto:P.J.Green@bristol.ac.uk)). Sylvia Richardson is Professor of Biostatistics, Department of Epidemiology and Public Health, Imperial College School of Medicine, Norfolk Place, London, W2 1PG (E-mail: [sylvia.richardson@ic.ac.uk](mailto:sylvia.richardson@ic.ac.uk)). The authors thank Julian Besag, Carmen Fernández, Alan Gelfand, Alex Lewin, Annie Mollié, and Christine Monfort for valuable interaction and comments. The work was partially supported by an EPSRC Visiting Research Fellowship, a travel grant from the ESF programme on Highly Structured Stochastic Systems, and the INSERM contract ITM 4TM05F. Most of this work was undertaken while Sylvia Richardson was at INSERM, Paris.

or smooth gradients can be exploited for further study or action. Indeed, the suspicion of a local excess in disease occurrence or the highlighting of geographic inequalities in health are important public health concerns that can be addressed through an analysis of spatial heterogeneity. Of course, the analysis must take into account all relevant risk factors that can be assessed at the area level. But it is hardly plausible that all the factors acting on the underlying disease risk can be identified or measured at the required geographic level. Thus there often remains residual heterogeneity in the disease event rate, which moreover is likely to have a spatial structure inherited from some of the unmeasured or undiscovered risk factors for the disease. Note that epidemiologic studies are observational by nature, and there is little or no control over the sources of variability. Furthermore, the delicate issue of ecological bias (Greenland and Robins 1994) must be kept in mind when interpreting sources of variability for disease outcomes analyzed at an aggregated level.

Modeling spatial heterogeneity of rare counts has usually been addressed in a hierarchic framework. We do the same in the development of our Bayesian approach, and specifically consider a Poisson model at the lowest level of the hierarchy,

$$y_i \sim \text{Poisson}(\lambda_i E_i) \quad \text{independently for } i = 1, 2, \dots, n, \quad (1)$$

where  $y_i$  denotes the observed count of disease incidences or deaths in area  $i$ ;  $E_i$  is the expected count based on population size, adjusted for, say, age and sex; and  $\lambda_i$  is an area-specific relative risk variable, the main object of our inference. We use the simple term “risk” to refer to the  $\{\lambda_i\}$  in the future.

This model may be extended straightforwardly to accommodate dependence on covariates  $\{x_{ij}\}$  measured in each area  $i$ , so that (1) is replaced by, for example,

$$y_i \sim \text{Poisson}(\lambda_i e^{\sum_j x_{ij} \gamma_j} E_i) \quad \text{independently for } i = 1, 2, \dots, n. \quad (2)$$

Illustrations of the use of the model both with and without covariates are given later in the article.

We now consider the choice of structure for the joint distribution of the  $\{\lambda_i, i = 1, 2, \dots, n\}$  at the next level of the hierarchy, a choice that can be influential on effective smoothing of the Poisson noise. In the seminal work of Besag, York, and Mollié (1991) and Clayton and Bernardinelli (1992) and subsequent work, log-linear Gaussian models for the  $\{\lambda_i\}$  were postulated using a conditional formulation that included a spatial autoregressive component based on contiguity in an undirected graph as well as a term modeling unstructured variability. This approach has been commonly adopted in recent work in disease mapping and has helped highlight many interesting features of the geographic distribution of some rare diseases. Alternative formulations of a multivariate Gaussian model for the  $\{\log \lambda_i\}$  that directly specify a spatially parameterized covariance matrix have also been discussed (Best, Arnold, Thomas, Waller, and Conlon 1999; Wakefield and Morris 1999). In both cases, the parameters characterizing the spatial dependence are constant across the entire study region, although models where the strength of spatial interaction is allowed to vary spatially have also been proposed in other contexts (Clifford 1986; Aykroyd and Zimeras 1999). When

using these parametric models, there is the potential risk of oversmoothing and masking of local discontinuities due to the global effect of the parameters. Concern about this is borne out by empirical studies, including the study reported in Section 4.

There have been several attempts to address this difficulty, which have in common the replacement of a continuously varying random field for  $\{\lambda_i\}$  by an allocation or partition model of the form

$$\lambda_i = \lambda_{z_i}, \quad (3)$$

where  $\{\lambda_j, j = 1, 2, \dots, k\}$  characterize  $k$  different components, and  $\{z_i, i = 1, 2, \dots, n\}$  are *allocation variables* taking values in  $\{1, 2, \dots, k\}$ . Moving the spatial dependence one level higher in the hierarchy to the discrete-valued process  $\{z_i\}$  has the potential of providing a greater degree of spatial adaptivity, again seen empirically. Note that discreteness in the prior is not imposed on posterior inference, in the sense that, marginalizing over the allocations, the posterior mean risk surface from any partition model can provide a smooth estimate of the risk surface. Models that can be described in this framework include the clustering or segmentation models of Knorr-Held and Raßer (2000) and Denison and Holmes (2001), which propose different spatial models for  $\{z_i\}$ . In the model investigated in this article we propose using a Potts model for  $\{z_i\}$ , with the number of states and strength of interaction unknown. In contrast to the partition models cited, we retain a Markovian structure for the  $\{z_i\}$ . Other models in this class are the spatial mixture models introduced by Fernández and Green (2002), in which the spatial dependence is pushed yet one level higher. The  $\{z_i\}$  are conditionally independent given weights  $w_{ij} = P(z_i = j)$  constructed from Gaussian random fields.

### 1.3 Mixtures and Other Related Models

Mixture models arise naturally whenever the existence of unknown subpopulations corresponding to different models for the quantity of interest can be hypothesized. They have found applications in many contexts, some of which are illustrated, along with comprehensive accounts of the theory, in the monographs by Titterton, Smith, and Makov (1985) and McLachlan and Peel (2000). It is not always possible or advisable to interpret the subgroups of areas identified by such models directly, so an important second perspective on the HMM adopted here is to view it as a semiparametric approach, following recent developments in both Bayesian and frequentist settings. Indeed, the question of the specification or the potential misspecification of the distribution of latent variables has been the subject of much attention (see e.g., Roeder, Carroll, and Lindsay 1996; Carroll, Roeder, and Wasserman 1999 for discussion of such issues in the measurement error context), and mixtures of distributions have been proposed as an alternative. Our proposed model is given extra flexibility by the treatment of the number of allocation classes and the strength of spatial interaction as unknowns, to be estimated together with the Poisson parameters.

Image analysis is another context in which hidden Markov random field models have been extensively used. Tjelmeland

and Besag (1998) provided a systematic study of Markov random fields with higher-order interactions, with the aim of producing well-calibrated posterior inference that goes beyond simple restoration. Johnson (1994) allowed a variable number of labels as we do, but generated rich geometric structures by constructing specific nonlocal potential functions. We stress that these extensions aim to recover high-level features of the image—an aim different from ours, of flexible analysis of spatial heterogeneity.

The article is organized as follows. In Section 2 we present the spatially correlated allocation model. We describe our MCMC implementation, which requires variable-dimension moves, in Section 3. In Section 4, we analyze the performance of the model on a collection of synthetic datasets designed to test different features of the model in the context of disease-mapping data. We also present the results of a simulation study aimed at comparing some aspects of its performance with that of the Markov random field formulation of Besag et al. (1991). We discuss an epidemiologic application to French cancer data in Section 5, and conclude with a discussion of extensions and further work in Section 6.

## 2. POTTS MODELS WITH POISSON NOISE

The main novelty of our approach lies in the modeling of the allocation variables  $\{z_i\}$  in (3). First, the number of components  $k$  is treated as unknown, with prior distribution  $p(k)$ , typically either truncated Poisson in form or uniform on some range  $\{1, 2, \dots, k_{\max}\}$ . Then, given  $k$ ,  $\{z_i\}$  follows a spatially correlated process.

The formulation for this process is built on a prescribed undirected graph, which plays the role of the prior conditional independence graph of the hidden random field  $\{z_i\}$ . Two areas,  $i$  and  $i'$ , are said to be neighbors, written as  $i \sim i'$ ,  $i \in \partial i'$ , or  $i' \in \partial i$ , if they are adjacent with respect to this graph. Typically, areas are taken to be neighbors in this sense if and only if they are spatially contiguous. More sophisticated spatial relationships can be modeled with little difficulty; we give an example of this in Section 4.6. Apart from this specific development, we always use spatial contiguity as our graph structure.

In this article we concentrate on the Potts model, an allocation model often used in image processing applications and originating in statistical physics. In contrast to the hierarchical mixture model defined by Richardson and Green (1997), and indeed most mixture models, this formulation does not make use of explicit weights on components.

In the Potts model formulation, the  $z_i$  are modeled jointly,

$$p(z|\psi) = e^{\psi U(z) - \theta_k(\psi)}, \quad (4)$$

where

$$U(z) = \sum_{i \sim i'} I[z_i = z_{i'}] \quad \text{and} \quad (5)$$

$$\theta_k(\psi) = \log \left( \sum_{z \in \{1, 2, \dots, k\}^n} e^{\psi U(z)} \right)$$

are the *number of like-labeled neighbor pairs* in the configuration  $z$  and an additive normalizing constant. The interaction

parameter  $\psi$  is nonnegative;  $\psi = 0$  corresponds to independent allocations, uniformly on the labels  $\{1, 2, \dots, k\}$ . The degree of spatial dependence increases with  $\psi$ , whereas allocations remain *marginally* uniform on  $\{1, 2, \dots, k\}$ . It is clear that for positive  $\psi$ ,  $p(z|\psi)$  favors probabilistically those allocation patterns where like-labeled locations are neighbors.

To complete the model specification, we must define our prior models for  $\lambda$ ,  $\psi$ , and  $k$ . We place an independence prior on  $\{\lambda_j, j = 1, 2, \dots, k\}$ :

$$\lambda_j \sim \Gamma(\alpha, \beta) \quad \text{independently for } j = 1, 2, \dots, k.$$

Although we have occasionally considered also a hierarchical model in which  $\beta \sim \Gamma(b_1, b_2)$ , we usually take the hyperparameters  $\alpha$  and  $\beta$  as fixed. Our standard choice is  $\alpha = 1$ ,  $\beta = \sum_i E_i / \sum_i y_i$ . Usually in epidemiologic applications,  $\sum_i E_i = \sum_i y_i$ , and hence  $\beta = 1$ ; in other cases, this choice makes the analysis equivariant to multiplicative misspecification of the  $\{E_i\}$ . Although it is not strictly necessary, we prefer to ensure identifiability of the labeling of mixture components by indexing the  $\{\lambda_j\}$  in numerically increasing order:  $\lambda_1 < \lambda_2 < \dots < \lambda_k$ . Thus the joint prior for  $\lambda$  becomes

$$p(\lambda|k, \alpha, \beta) = k! I[\lambda_1 < \lambda_2 < \dots < \lambda_k] \prod_{j=1}^k \frac{\beta^\alpha \lambda_j^{\alpha-1} e^{-\beta \lambda_j}}{\Gamma(\alpha)}.$$

The whole issue of labeling and the impact of ordering on MCMC performance was comprehensively discussed in the discussion and rejoinder to Richardson and Green (1997). (Also see Stephens 2000 for an alternative approach).

We take  $p(\psi)$  to be a discrete distribution, uniform on the values  $\{0, 0.1, \dots, \psi_{\max}\}$ . The uniformity is arbitrary—as usual, other forms of distribution could be substituted after sampling, using importance reweighting. The discreteness, which we do not believe has a significant impact on our inference, is for the sake of computational convenience, because the normalizing constants  $\theta_k(\psi)$  can then be precomputed offline and stored in a table, with no approximation or interpolation necessary at run time. Finally, our prior on the number of components  $k$  is uniform on the values  $\{1, 2, \dots, k_{\max}\}$ .

Prior simulations are useful to inform the choice of  $\psi_{\max}$ . For the Potts model on the contiguity graph of the French départements used in our studies reported in Section 4, we found that the value  $\psi = 1.0$  is a very high level of interaction; the average prior probability that two neighboring regions have the same label when  $\psi = 1.0$  and  $k = 2$  is .96, declining only to .70 when  $k = 8$ . We thus choose default values of  $\psi = 1.0$  and  $k = 10$ , which seem sufficient to ensure flexibility for all practical purposes, although these could easily be increased if deemed unsuitable for a specific graph. In fact, exceptionally, in one of our examples, we extend  $\psi_{\max}$  up to 1.2.

All the foregoing specifications are expressed somewhat loosely, in the interest of economy of notation. Each of the model ingredients is actually a conditional distribution for the stated variable, conditional on both its immediate parameters and hyperparameters higher up in the directed acyclic graph (DAG) describing the model (Fig. 1). Thus the joint distribution of all variables corresponding to the Potts model formulation is

$$p(k)p(\psi)p(\lambda|k, \alpha, \beta)p(z|k, \psi)p(y|\lambda, z).$$

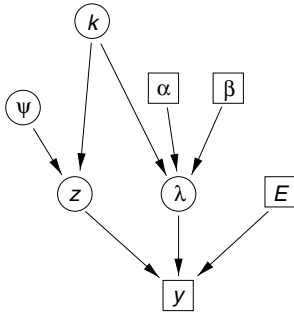


Figure 1. DAG for the Potts Spatial Mixture Model.

### 3. MARKOV CHAIN MONTE CARLO

Naturally, MCMC methods are needed to fit these spatial mixture models. Our sampler for the Potts mixture model uses four different moves, each move updating a subset of the variables, under detailed balance with respect to the posterior distribution, which forms an irreducible Markov chain when used together in a deterministic scan. In this general structure, our computational method follows other recent work, including that of Richardson and Green (1997). Three of the four are standard fixed-dimension moves; the fourth move proposes to change dimension by increasing or decreasing the number of components.

#### 3.1 Fixed-Dimension Moves

The three fixed-dimension moves update the allocations  $z$ , the spatial interaction parameter  $\psi$ , and the component parameters  $\{\lambda_j\}$ . The allocations  $z$  are updated by a Gibbs kernel. The full conditional from which an update for  $z_i$  is drawn is

$$p(z_i = j | \dots) \propto e^{-\lambda_j E_i} \lambda_j^{y_i} e^{\psi n_{ij}},$$

where  $n_{ij} = \#\{i' \in \partial i : z_{i'} = j\}$  is the number of neighbors of  $i$  currently assigned to component  $j$ . Note that in contrast to the simple random sample mixture case of Richardson and Green (1997), here the  $\{z_i\}$  are not conditionally independent given all other variables, so they may not be updated simultaneously.

The interaction parameter  $\psi$  has full conditional

$$p(\psi | \dots) \propto p(\psi) e^{\psi U(z) - \theta_k(\psi)},$$

a discrete distribution on a finite grid of values, like the prior. A random walk Metropolis kernel, proposing perturbations of  $\pm 1$  with equal probability, is convenient for updating  $\psi$ .

Finally, we need to update the parameters  $\{\lambda_j\}$ . An approach to the simultaneous update of these, maintaining the order restriction, exploits the following simple trick that we have not seen elsewhere. We propose simultaneous independent zero-mean normal increments to each  $\log \lambda_j$ ; the modified values of  $\lambda$  are then placed in increasing order to give say,  $\{\lambda'_j\}$ . Remarkably, the fact that the proposal density that we are using is actually a sum of  $k!$  rather complicated terms, due to the reordering, does not matter; the terms that appear in the sums in the numerator and denominator of the Metropolis–Hastings ratio are in constant proportion and so cancel out. The acceptance probability for the complete set of updates,

formed from the prior ratio, the likelihood ratio, and a Jacobian for the log transformation, reduces to

$$\min \left\{ 1, \prod_{j=1}^k \left[ \left( \frac{\lambda'_j}{\lambda_j} \right)^{\alpha + \sum_{i:z_i=j} y_i} \exp\{-(\lambda'_j - \lambda_j)(\beta + \sum_{i:z_i=j} E_i)\} \right] \right\}.$$

#### 3.2 Variable-Dimension Move

Changing the number of components under detailed balance with respect to the posterior requires a reversible jump move (Green 1995). We follow the general idea of a random choice between splitting an existing component into two components and merging two existing components into one component, as used by Richardson and Green; the probabilities of these alternatives are  $b_k$  and  $d_k$  when there are currently  $k$  components. Along with incrementing or decrementing  $k$ , the move also entails modifying the vector  $\lambda$  accordingly and reallocating observations into the new component(s) as necessary.

In contrast to Richardson and Green, we do the reallocation part of the proposal not independently for each observation, but rather in a way that approximately respects the spatial structure of the Potts model. This is with the usual aim of increasing the probability of the move's acceptance; exact detailed balance is, of course, ensured by correctly calculating the acceptance probability in terms of the model and proposal probabilities.

We describe the split move in some detail. First, a component to be split, say,  $j$ , is chosen uniformly at random from  $\{1, 2, \dots, k\}$ . This is replaced by two components that we label “–” and “+,” with  $\lambda$  values generated by

$$\lambda_- = \lambda_j u^c \quad \text{and} \quad \lambda_+ = \lambda_j u^{-c},$$

where  $u$  is generated from  $U(0,1)$  and  $c$  is a proposal spread parameter that we set at 0.1. If  $\lambda_- < \lambda_{j-1}$  or  $\lambda_+ > \lambda_{j+1}$  (with appropriate modifications if  $j = 1$  or  $k$ ), then the move is rejected, as the misordered vector has zero density under the ordered prior. Those observations  $i$  currently allocated to component  $j$  are then dynamically reallocated between – and +. We scan over such  $i$ , and the probability with which  $z_i$  is set to – rather than + is

$$\frac{e^{\psi n_- - \lambda_- E_i} \lambda_-^{y_i}}{e^{\psi n_- - \lambda_- E_i} \lambda_-^{y_i} + e^{\psi n_+ - \lambda_+ E_i} \lambda_+^{y_i}},$$

where  $n_-$  and  $n_+$  are the numbers of areas adjacent to  $i$  already proposed for allocation to – and + in this scan. This choice has the affect of mimicking the Potts model term in the target distribution to favor proposed allocations with spatial coherence. As the scan proceeds, the probability,  $P_{\text{alloc}}$ , of the allocation actually generated is accumulated. Denote the proposed new allocation vector by  $z'$ . Following logic very similar to that of Richardson and Green (1997, Eq. 11), the acceptance

probability for this complete proposal is  $\min\{1, R\}$ , where

$$\begin{aligned} R &= \prod_{i:z'_i=-} e^{-(\lambda_- - \lambda_j)E_i} \left(\frac{\lambda_-}{\lambda_j}\right)^{y_i} \prod_{i:z'_i=+} e^{-(\lambda_+ - \lambda_j)E_i} \left(\frac{\lambda_+}{\lambda_j}\right)^{y_i} \\ &\times \frac{\beta^\alpha}{\Gamma(\alpha)} \left(\frac{\lambda_- \lambda_+}{\lambda_j}\right)^{\alpha-1} e^{-\beta(\lambda_- + \lambda_+ - \lambda_j)} (k+1) \frac{p_{k+1}}{p_k} \\ &\times \exp\{\psi(U(z') - U(z)) + \theta_k(\psi) - \theta_{k+1}(\psi)\} \\ &\times \frac{d_{k+1}}{b_k P_{\text{alloc}}} \times \frac{2c\lambda_j}{u}. \end{aligned}$$

### 3.3 Approximating the Potts Model Partition Function

Although the MCMC moves for  $z$  and  $\lambda$  make no reference to  $\theta_k(\psi)$ , values of this normalizing constant, the logarithm of what is known as the partition function in statistical physics, are needed for the update for  $\psi$  and the split/merge move. Because in our model both  $k$  and  $\psi$  are discrete, we need to evaluate  $\theta_k(\psi)$  on a grid of  $(k, \psi)$  pairs. These are computed in offline simulations, specific to the assumed neighborhood graph, and provided in a look-up table for our MCMC sampler to use. We have found the following simple method for estimating  $\theta_k(\psi)$  easy to use and reliable on graphs of the size that we have encountered in real disease-mapping applications.

*3.3.1 The Thermodynamic Integration Approach.* This method has a long history; according to Gelman and Meng (1998), it has been used in statistical physics since the 1970s, and Ogata and Tanemura (1984) are responsible for its first use in spatial stochastic processes. Consider the Potts model with  $k$  labels on a graph with  $n$  vertices, defined in (4) and (5). Differentiating  $\theta_k(\psi)$ , we obtain

$$\begin{aligned} \frac{\partial}{\partial \psi} \theta_k(\psi) &= \frac{\partial}{\partial \psi} \log \sum_{z \in Z} e^{\psi U(z)} \\ &= \sum_{z \in Z} U(z) p(z|\psi) \\ &= E(U(z)|\psi, k), \end{aligned} \quad (6)$$

the expectation of  $U(z)$  when  $z$  is distributed according to the assumed Potts model. Here  $Z = \{1, 2, \dots, k\}^n$  is the set of possible labelings of the graph. But  $\theta_k(0) = \log \sum_{z \in Z} 1 = n \log k$ , so

$$\theta_k(\psi) = n \log k + \int_0^\psi E(U|\psi', k) d\psi'. \quad (7)$$

In particular, note that when  $k = 1$ ,  $U(z)$  is the constant  $n_E$ , the number of edges in the graph, so that  $\theta_1(\psi) = \psi n_E$ .

Equation (7) is the basis of a simple method for estimating  $\theta_k(\psi)$ ; the expectation is replaced by a sample average in a MCMC simulation of the Potts model for specific  $(k, \psi')$ . In our implementation, the integral is computed by numerical integration of a cubic spline smooth of the simulated averages with respect to  $\psi'$ .

*3.3.2 Improving Partition Function Estimates Using the Potts Model Mixture Sampler.* Suppose that we place a prior  $p(k, \psi)$  on  $(k, \psi)$  and conduct an MCMC simulation of the distribution of  $(k, \psi, z)$  assuming an approximate trial value  $\tilde{\theta}_k(\psi)$  for the log partition function. This can be accomplished by the posterior sampler for the Potts model mixture model derived in Sections 3.1 and 3.2, with the data suppressed and likelihood terms omitted from the model. We are simulating from the joint distribution

$$\propto p(k, \psi) \exp(\psi U(z) - \tilde{\theta}_k(\psi)), \quad (8)$$

so the marginal distribution for  $(k, \psi)$  is

$$\propto p(k, \psi) \exp(\theta_k(\psi) - \tilde{\theta}_k(\psi)). \quad (9)$$

This observation can be used in two ways. First, comparing the prior  $p(k, \psi)$  to the observed frequencies  $\hat{p}(k, \psi)$ , say, provides a check on the departure of  $\tilde{\theta}_k(\psi)$  from  $\theta_k(\psi)$ . Second, if we instead equate (9) to  $\hat{p}(k, \psi)$ , then we can solve to give improved estimates of  $\theta_k(\psi)$ , namely

$$\theta_k(\psi) = \tilde{\theta}_k(\psi) + \log(\hat{p}(k, \psi)/p(k, \psi)), \quad (10)$$

up to an additive constant. This device takes advantage of the often improved mixing offered by variable-dimension samplers. It cannot cope with very poor initial estimates, but works well as a supplement to thermodynamic integration or any other method for deriving the normalizing constant. In the numerical experiments that follow, we use thermodynamic integration, based on runs of length 50,000 for each  $(k, \psi)$  combination, followed by the improvement just described.

A full analysis of the problem of estimating normalizing constants has been given by Gelman and Meng (1998). They discussed several methods that are more sophisticated, but also more cumbersome; these will handle more challenging problems and might be needed to adapt our methods to bigger graphs or a wider range of  $(k, \psi)$  values.

### 3.4 Implementation of the BYM Model

As mentioned in Section 1, we make use of the method of Besag et al. (1991) as a standard for comparison in our experiments. This is based on a hierarchical model; given the values of the variance parameters, the logarithms of the risks have a certain multivariate normal distribution, a priori. In contrast, in our model, conditional on the allocations and hyperparameters, the risks are gamma distributed. To eliminate the impact on our comparisons of this basic difference, the model that we implement, and refer to as the BYM model, is a minor reformulation of the model of Besag et al.

We suppose that the risk in area  $i$  is  $\lambda_0 e^{u_i + v_i}$ , where, conditional on  $\alpha, \beta, \tau_u$ , and  $\tau_v$ , the terms  $\lambda_0, u = (u_i)$  and  $v = (v_i)$  are independent, with

$$\lambda_0 \sim \Gamma(\alpha, \beta),$$

$$p(u|\tau_u) \propto \exp(-\tau_u \sum_{i \sim i'} (u_i - u_{i'})^2/2),$$

and

$$p(v|\tau_v) \propto \exp(-\tau_v \sum_{i=1}^n v_i^2/2),$$

where both  $\{u_i\}$  and  $\{v_i\}$  are constrained to sum to 0. (In the case of a disconnected graph, we would apply these constraints separately in each connected component, but in fact we do not need this in our examples herein.) In this formulation, if  $\tau_u$  and  $\tau_v$  go to  $\infty$ , then we obtain the standard nonspatial conjugate model, which also arises if  $k$  is fixed at 1 in our mixture model. Allowing  $k > 1$  or  $\tau_u, \tau_v < \infty$  are alternative approaches to fitting spatial heterogeneity. We assume that  $\tau_u \sim \Gamma(\alpha_u, \beta_u)$  and  $\tau_v \sim \Gamma(\alpha_v, \beta_v)$  a priori, with  $\alpha_u = \beta_u = \alpha_v = \beta_v = .1$ .

Constructing a sampler for the BYM model requires some care both because of the strong spatial interaction among the  $u$  variables and because of the high correlation a posteriori between  $u$  and  $\tau_u$  and between  $v$  and  $\tau_v$ . In addition, the sum-to-zero constraints that we impose in our form of the model pose extra difficulties. Recent work by Knorr-Held and Rue (2002) has focused on the first of these problems. Our sampler, which seems to perform quite adequately, explicitly addresses the second and third problems. The variables in question are updated in two blocks,  $(u, \tau_u)$  and  $(v, \tau_v)$ , using moves of a similar design. Taking the  $(u, \tau_u)$  block as an example, we use a Metropolis–Hastings proposal on the whole vector  $u$ , making simultaneous independent Gaussian perturbations constrained to sum to 0. The proposal is accepted or rejected by reference to the target distribution in the usual way, but with  $\tau_u$  integrated out; that is, the acceptance ratio for the update from  $u$  to  $u'$  is

$$\frac{(\beta_u + (1/2) \sum_{i \sim j} (u'_i - u'_j)^2)^{\alpha_u + n/2}}{(\beta_u + (1/2) \sum_{i \sim j} (u_i - u_j)^2)^{\alpha_u + n/2}} \times \exp \left\{ \sum_{i=1}^n y_i (u'_i - u_i) - \lambda_0 (e^{u'_i + v_i} - e^{u_i + v_i}) \right\}.$$

Following this, whether or not the update for  $u$  is accepted,  $\tau_u$  is updated by a Gibbs move, drawing a new value from its full conditional  $\Gamma(\alpha_u + n/2, \beta_u + (1/2) \sum_{i \sim j} (u_i - u_j)^2)$ . We adjust the spread of the Gaussian perturbations on  $u$  and  $v$  in pilot runs to achieve reasonable acceptance rates. On larger graphs, it might prove necessary to apply such perturbations to subsets of the areas instead of the whole graph. Finally,  $\lambda_0$  is updated by a Gibbs kernel.

## 4. MODEL PERFORMANCE AND COMPARISON

### 4.1 The Simulated Datasets

In this section we investigate the distinguishing features of the model and its performance on simulated datasets. Throughout, we use the spatial layout of the 94 mainland French départements. To test different characteristics of the model, we generated three datasets corresponding to contrasting geographic features of the underlying simulated risks. Specifically, the “Block4” case refers to a situation of a background value of  $\lambda$  equal to .7, with four groups of areas having values of  $\lambda$  equal to 1.5. These four groups consist of either a single well-populated area or a group of five rural départements or are on the border (Fig. 2). “North-South” simulates a simple north/south divide, with  $\lambda$  equal to .8 in the north and 1.2 in the south (Fig. 3), whereas “Gradient NS” corresponds to risks smoothly (linearly) decreasing from north to south (Fig. 4). For each dataset, the observed number of events were simulated as

$$y_i \sim \text{Poisson}(\lambda_i E_i) \quad \text{independently for } i = 1, 2, \dots, n,$$

where the expected numbers of events were chosen on the basis of the French population structure and correspond to real data on two types of cancer. For Block4 and “North-South,” these numbers correspond to the expected number of deaths for laryngeal cancer in females for the period 1986–1993 and range from 2 to 58, whereas for Gradient NS they are about three times larger and correspond to those of the gall bladder dataset analyzed by Mollié (1996), for which risks were found to have a gradient-like structure. Figures 2(b), 3(b), and 4(b) display the maximum likelihood estimate of  $\lambda_i$ ,  $y_i/E_i$ , commonly referred to as the standardized mortality ratio (SMR) in the epidemiologic literature.

### 4.2 Output Analysis and Criteria

All of the results displayed correspond to runs of 500,000 sweeps of the algorithm after a burn in of 20,000 sweeps. The mixing performance of the split and merge moves was satisfactory, with acceptance rates generally around 10%, except in cases where the data support a low number of components, as in the North-South example, where the acceptance rate

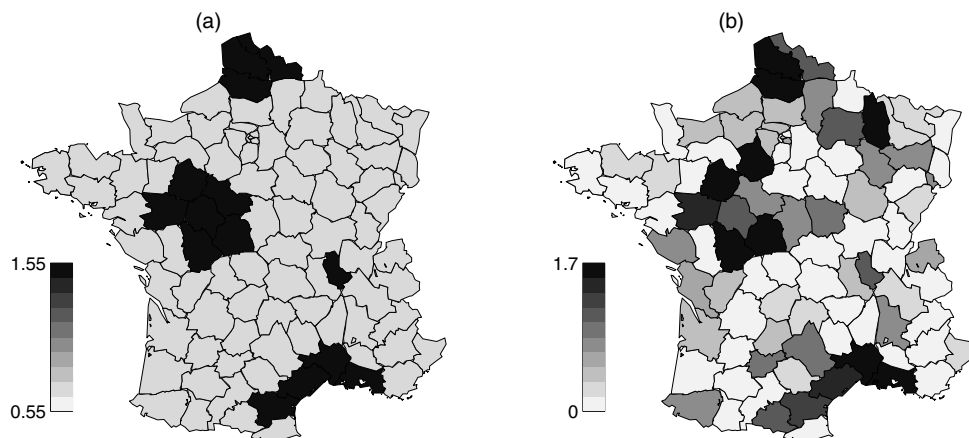


Figure 2. True Risks (a) and Observed SMRs (b) for the Block4 Dataset.

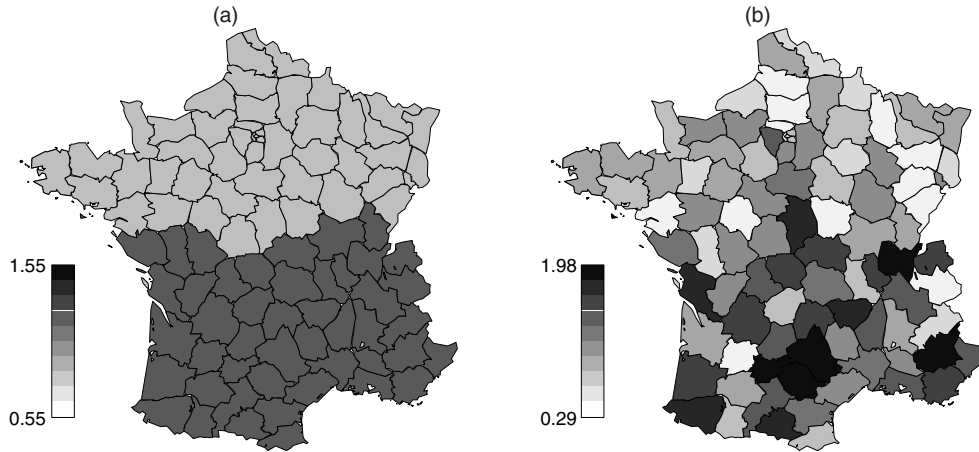


Figure 3. True Risks (a) and Observed SMRs (b) for the North-South Dataset.

drops to 5%. From the samples, different summaries of the posterior distribution can be computed. Our main interest is in the spatial variation of the  $\{\lambda_{z_i}\}$  and associated posterior probabilities.

The section concludes by reporting some simulation results that aim to compare some aspects of the performance of our spatial mixture model to that of the BYM model. Our criteria for comparison are fairly straightforward. For simulated data, we know the “true” underlying risks, and these will be denoted by  $\{\lambda_i^t\}$ . We then calculate for each area  $i$ ,  $MSE_i = E((\lambda_i^t - \lambda_i)^2 | y)$ , the (posterior) mean squared error (MSE), where  $\lambda_i$  corresponds to  $\lambda_{z_i}$  for the mixture model and to  $\lambda_0 \exp(u_i + v_i)$  for the BYM model. To summarize the performance over the whole map, we compute  $RAMSE = (\sum_i MSE_i / n)^{1/2}$ , the root averaged MSE. Because this criterion has the disadvantage of penalizing multiplicative overestimation of a risk more than underestimation, we also compute a corresponding criterion on the log scale—that is, replacing above  $\lambda_i^t$  and  $\lambda_i$  by  $\log \lambda_i^t$  and  $\log \lambda_i$ . The corresponding summary over the whole map, denoted by  $RAMSEL$ , now treats symmetrically a risk that is, say, doubled or halved.

Turning to a measure applicable to real data, where the true risk map is not available and that aims to balance fit

and complexity, we have computed the deviance information criterion (DIC) proposed by Spiegelhalter, Best, Carlin, and van der Linde (2002). The DIC is the sum of two terms:  $E(D|y)$ , the posterior expected deviance, and  $p_D$ , a penalty term.  $E(D|y)$  is evaluated from the MCMC output in a standard way; at each sweep, values of the parameters are produced from which can be calculated the Poisson deviance,  $D = 2 \sum_i (y_i \log(y_i / \mu_i) - y_i + \mu_i)$ , where  $\mu_i = \lambda_{z_i} E_i$ . The penalty  $p_D$  is the difference between the posterior expectation of this deviance and the deviance at the posterior mean of the parameters. Here we have used the posterior mean of the  $\{\lambda_{z_i}\}$  for the mixture model and that of  $\lambda_0 \exp(u_i + v_i)$  for the BYM model. For each dataset and each model, we report the value of DIC, of  $E(D|y)$ , and of  $p_D$ , which can be interpreted as a measure of model complexity. (See also Best et al. 1999 for a discussion of using DIC in comparisons of spatial models.)

### 4.3 Posterior Inference on $k$ and $\psi$

Figure 5 displays the joint posterior distribution of  $k$  and  $\psi$  for the three datasets. There are clear differences between the “images.” There is support for low values of  $k$  in the first two datasets, whereas larger  $k$  and  $\psi$  are necessary to fit the gradient-like structure. The posterior for  $\psi$  concentrates

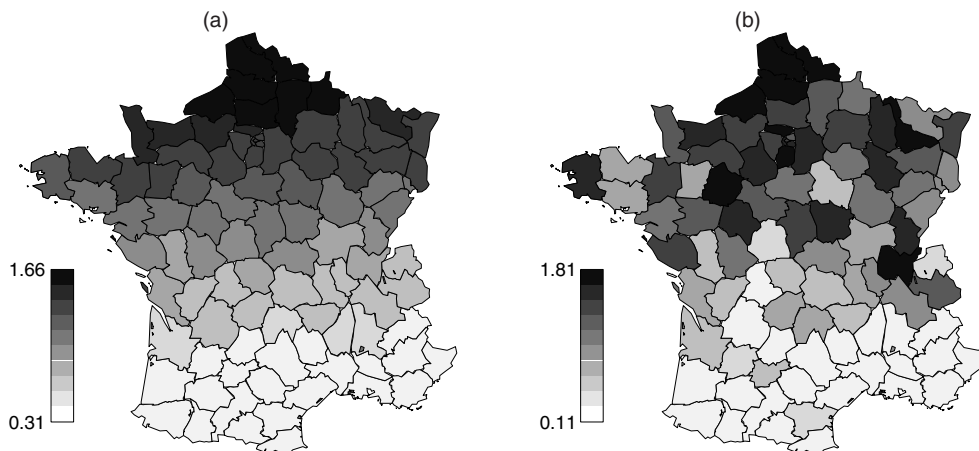


Figure 4. True Risks (a) and Observed SMRs (b) for the Gradient NS Dataset.

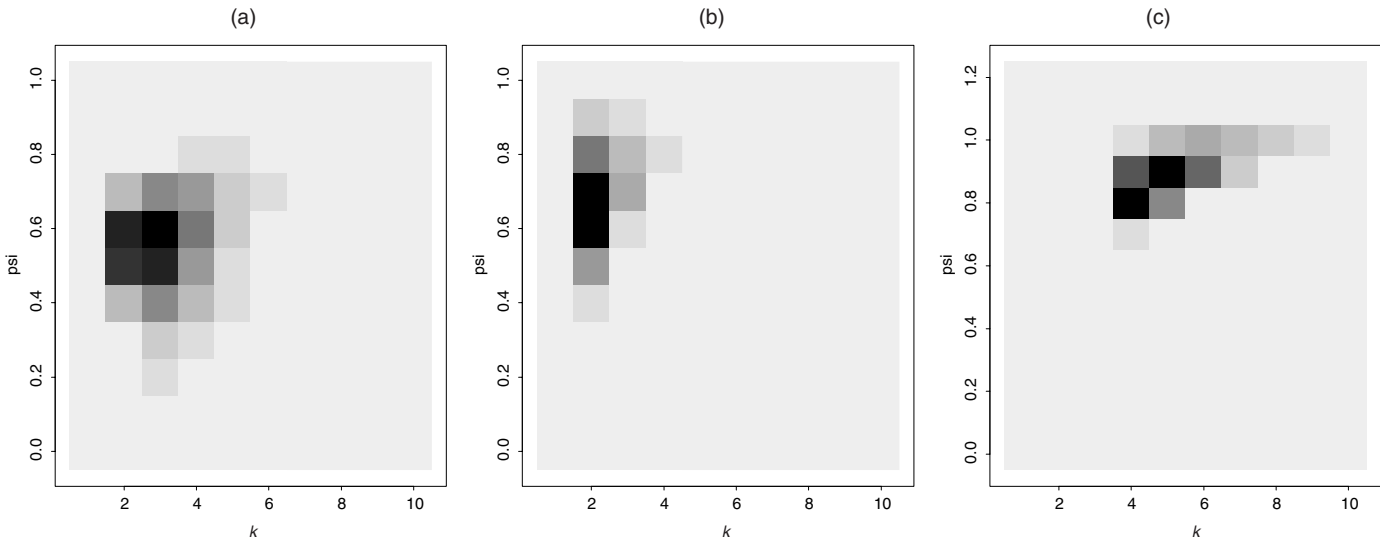


Figure 5. Joint Posterior Distribution of  $k$  and  $\psi$  for the Three Simulated Datasets: (a) Block4, (b) North-South, and (c) Gradient NS.

around higher values when relatively large “clusters” (i.e., neighboring areas with the same  $\lambda$ ) exist in the true setup. The different types of correlation between  $k$  and  $\psi$  show the adaptivity of the mixture to different geographic patterns. For North-South, there is a clear peak of  $p(k|y)$  for  $k = 2$ , the true number of components, whereas for Block4, the mode of  $p(k|y)$  is at  $k = 3$ , with the mixture model preferring to fit more than one component to model the variability of the large number of areas having the background risk.

#### 4.4 Posterior Estimates of the $\{\lambda_{z_i}\}$

Figures 6(a) and 7 show the posterior mean of  $\lambda_{z_i}$  for the three simulated datasets. Visually comparing these to the true simulated risks shows an excellent match. Because of model averaging, the posterior means of the  $\lambda_{z_i}$  are smoothly varying over the space and are not steplike. The flexibility of the mixture to adapt to very different patterns of risks is apparent. Figure 6(b) displays the map of posterior standard deviations of the  $\lambda_{z_i}$  for North-South. Note that variability is a little higher on the border areas between the contrasting zones; of course, this is also modulated by the size of the expected

counts. Because mean estimates can be a misleading summary in cases of high variability or skewness, Figure 8(a) displays for Block4 a representation for each area of the posterior distribution of the  $\{\lambda_{z_i}\}$  as a five-bin histogram with break points at .7, .9, 1.1, and 1.3. We see that the histograms corresponding to areas of simulated elevated risk in the four blocks are clearly right-skewed in comparison to the prevalent left-skewed histograms corresponding to the background areas. It is also interesting to display maps of posterior probabilities that the risk in each area exceeds certain thresholds. These can be easily computed from the output. Figure 8(b) shows that  $\Pr(\text{RR} > 1)$  in the North-South example provides a clear indication that the southern areas have more elevated risks than the northern areas.

Another posterior summary that is easily obtainable from our mixture model is the posterior distribution of the allocations  $z_i$  between different components, conditional on values of  $k$ . Such allocation graphs are simple and visually effective in isolating areas of particularly high or low risk. But their interpretation is conditioned by the separation between the components, and in the examples that we looked at, we

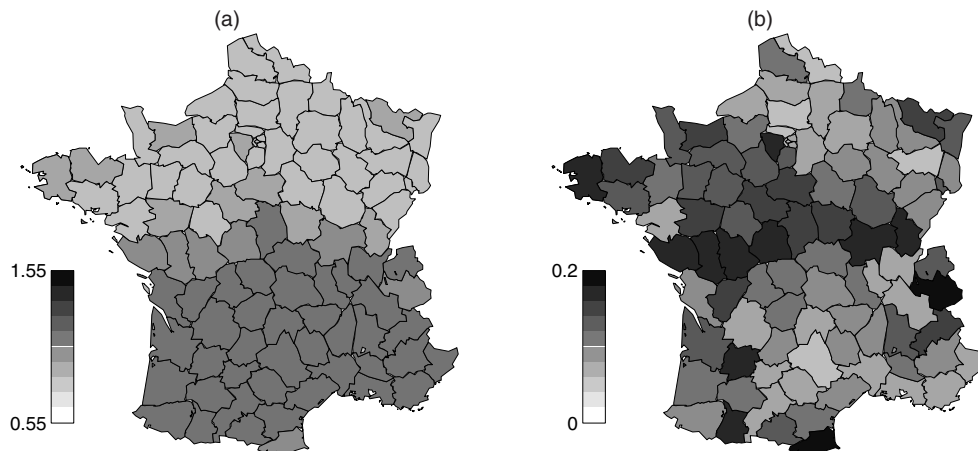


Figure 6. Mixture Estimates of the Risks for North-South, Posterior Means (a) and Posterior Standard Deviations (b).



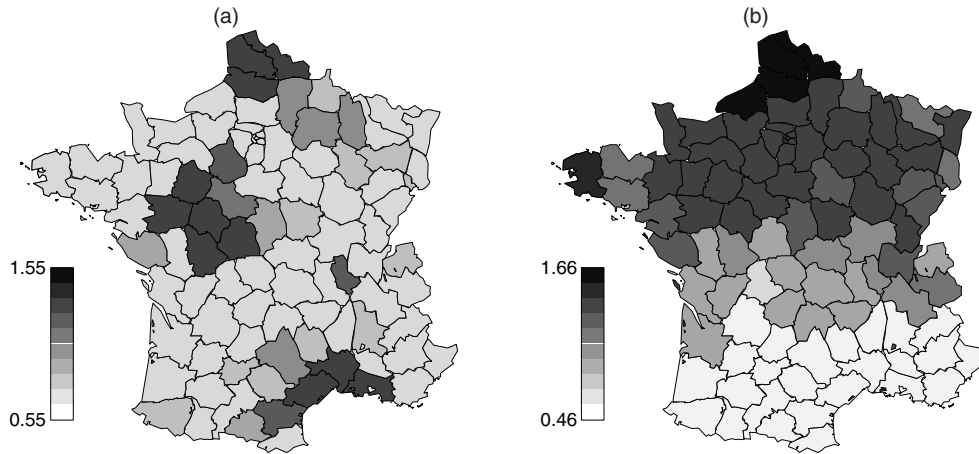


Figure 7. Mixture Estimates of the Risks for Block4 (a) and Gradient NS (b): Posterior Means.

did not find that such graphs uncover new features not visually apparent in the histogram of the  $\{\lambda_{z_i}\}$ . One important feature to point out is the wide posterior variability of these allocations. Indeed, a priori, our model assumes that for each  $k$ , the different labels are equally probable. This prior assumption, akin to the assumption of uniformly distributed weights usually made in mixture models, nevertheless allows the posterior allocations to be far from uniform when there is information in the data. This can be seen in Figure 9, which displays the boxplots of the modal allocation probabilities for the Block4 and the Gradient NS datasets and a selected range of values of  $k$ . To be precise, for each area  $i$  and each  $k$ , we determine  $\max_j\{P(z_i = j|k, y)\}$  and form a boxplot of these probabilities over  $i$ .

For Block4, most of the modal allocation probabilities are above .90 when  $k = 2$ , reflecting the real contrast in the data. When  $k \geq 3$ , the areas with background risk are split between several components with close values of  $\lambda$ , and their allocation probabilities are much closer to their prior mean of  $1/k$  (indicated by a dot in Fig. 9) as could be expected, because there is little information in the data about these further com-

ponents. For Gradient NS, there is substantial structure in the data, and the modal allocation probabilities reflect this, being well above  $1/k$  for all values of  $k$ .

### 4.5 Clusters

We define a *cluster* as a set of like-labeled areas connected by paths from neighbor to neighbor. To be precise, areas  $i$  and  $j$ , say, are in the same cluster if there is a path  $i = l_0 \sim l_1 \sim \dots \sim l_r = j$  such that  $z_{l_p}$  is the same for all  $0 \leq p \leq r$ . Note that in the segmentation approach of Knorr-Held and Raßer (2000), each cluster is labeled differently, whereas in our model, disconnected areas can have the same label. Thus it is interesting to compare the prior and posterior distributions of the number of clusters  $m$ . This can be done conditionally on a fixed value of  $k$ , and we do this in our synthetic examples for the “true” and the modal  $k$ , or integrating over  $k$  using all of the output.

Figure 10 displays the prior distribution of  $m$  (integrating over the uniform priors for  $k$  and  $\psi$ ), as well as the posterior

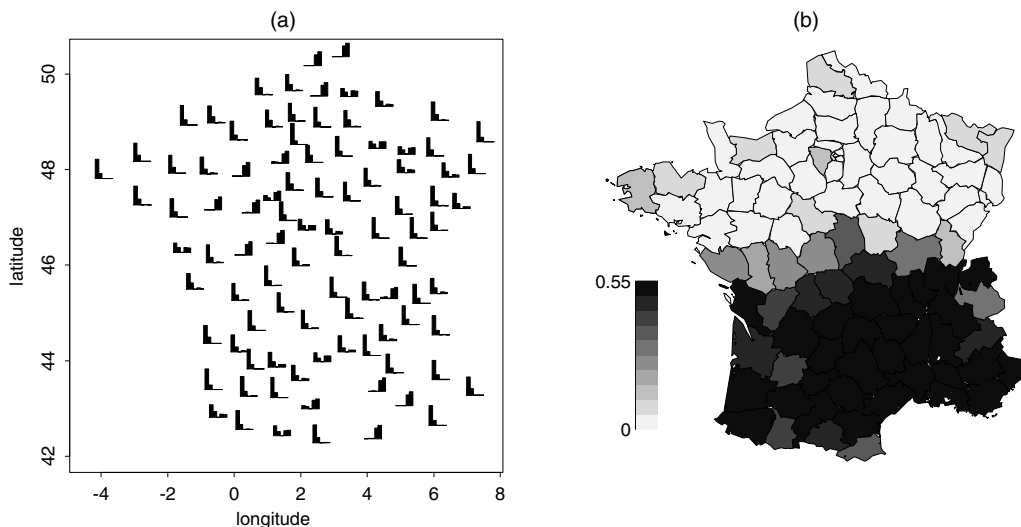


Figure 8. Examples of Posterior Summaries Obtained Using the Mixture Model: Histograms of Posterior Risks for Block4 (a) and Posterior Probability of Risk > 1 for North-South (b).

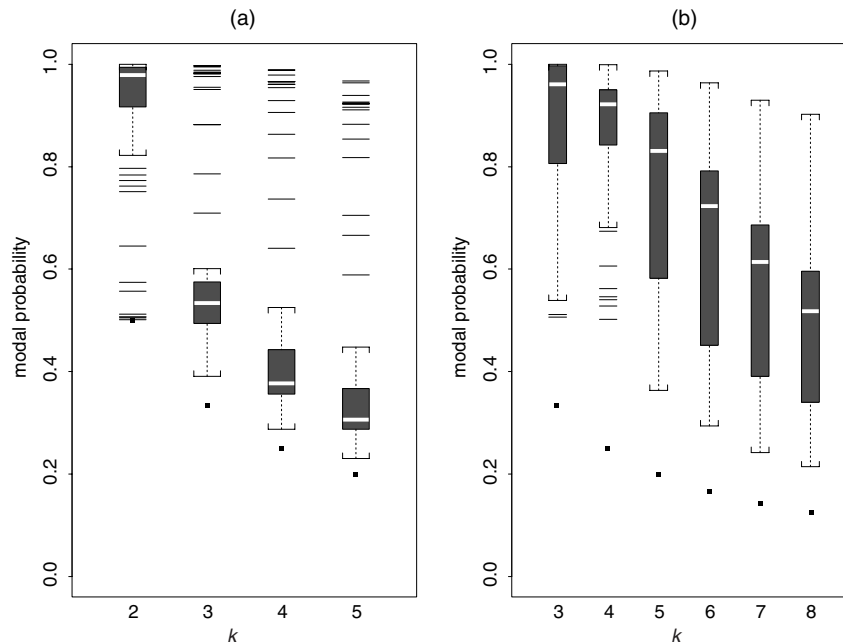


Figure 9. Boxplots of the Allocation Probabilities for the Block4 (a) and Gradient NS (b) Datasets. The dots represent the prior means of  $1/k$ .

distribution of  $m$  for the three datasets, integrating over  $k$  (solid line), conditional on the value of  $k$  corresponding to the “true  $k$ ” when it exists (dotted line), and conditional on the mode of  $p(k|y)$  (dashed line). We see peaked patterns for the three datasets, contrasting with the flat shape and extended tail of the distribution of  $m$  for the prior model. This shows that the spatial pattern in the data resulted in concentration of the posterior distribution of  $m$  on smaller values.

A simple structure for  $p(m|y)$  is apparent for North-South and Gradient NS, with hardly any shift between the conditional and the overall cluster distribution. On the other hand, for Block4, the bimodality of  $p(m|y)$  indicates hesitation between a simple and a more complex clustering pattern arising when the areas with the background risks are split between more than one component. We see the flexibility of our mixture model in generating different cluster patterns. Clearly,

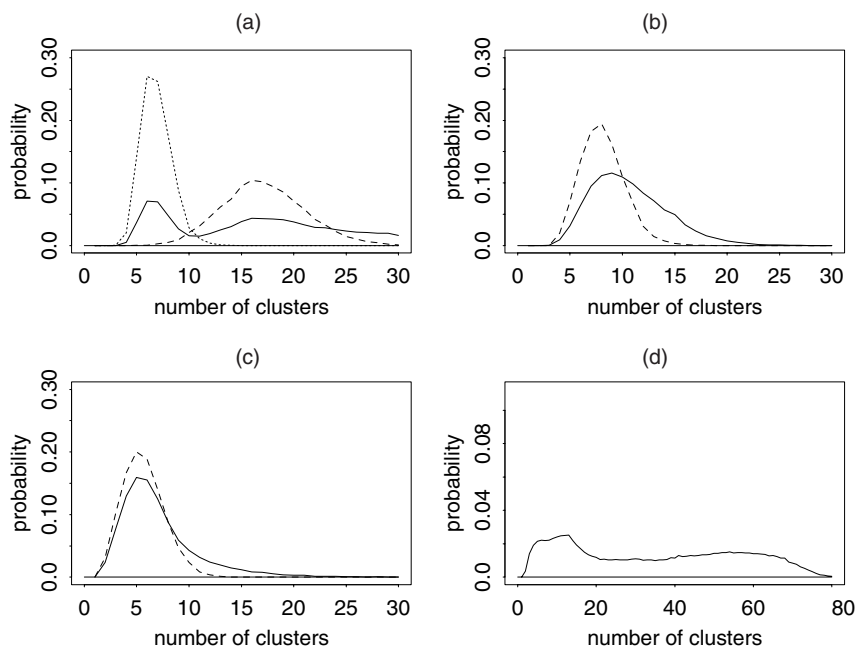


Figure 10. Posterior Distributions of the Number  $m$  of Clusters Obtained Using the Mixture Model, Integrated Over  $k$  (solid line), Conditional on “True”  $k$  (dotted line), and Conditional on the Posterior Mode of  $k$  (dashed line). (a) Block 4; (b) Gradient NS; (c) North-South; (d) Prior. In the Gradient NS case, there is of course no true  $k$ ; in the North-South case, the true and posterior modal  $k$  coincide. For the Prior case, we plot only the integrated distribution, conditional on  $k > 1$ .

using the Potts model in conjunction with variable  $k$  has not “frozen” the cluster pattern, as might have been anticipated from the experience of Tjelmeland and Besag (1998) using the Potts model with a fixed  $k$  in a related context in image analysis.

#### 4.6 Assumptions and Sensitivity

Our hierarchical formulation involves different levels of assumptions: distributional, quantitative (in relation to hyperparameter specification), and structural. The distributional assumption for the observed counts must be adapted to the data. For disease mapping, it is standard to use the Poisson assumption, but in other cases where spatial HMMs are used (e.g., in ecologic applications), this component of the model could be replaced by alternative distributions, such as the binomial or negative binomial.

As in any mixture-like problem, one might anticipate that some aspects of the model are sensitive to the choice of prior distribution for the component parameters. We have chosen to use gamma distributions for the  $\lambda$ s with  $\alpha = 1$  and  $\beta = \sum_i E_i / \sum_i y_i$ . We conducted a sensitivity study of this choice by letting  $\alpha$  take the alternative values of .4 or 2.5, with  $\beta$  adjusted correspondingly so that  $\alpha/\beta = \sum_i y_i / \sum_i E_i$ ; we also allowed an additional level in the hierarchy and treated  $\beta$  as random, with  $\beta \sim \Gamma(b_1, b_2)$ . The posterior distributions of the area-specific risks were highly stable under these various choices. For the three datasets, neither the posterior means nor the posterior standard deviations vary by more than .03 from their values under the standard choice  $\alpha = 1$  and  $\beta = \sum_i E_i / \sum_i y_i$ . This is a welcome feature of the model. On the other hand, as anticipated from other mixture studies that we have conducted (Richardson and Green 1997), inference on  $k$  and  $z$  is less stable, with a tendency for the model to fit more components as the variability of the gamma distribution is decreased. This sensitivity is the reason that one must be careful to not overinterpret the posterior on  $k$ , and to use the mixture simply for exploring interesting features of the heterogeneity.

The single-parameter Potts model on the nearest-neighbor graph with variable number of labels is a particular formulation of a spatial HMM that allows computational tractabil-

ity and, we believe, sufficient flexibility. In some applications, allowing higher-order neighbors might be called for, and this can be done without changing the computational strategy. However, it becomes more cumbersome to compute the look-up tables for the log partition function  $\theta_k(\psi)$  if there are additional parameters. We have investigated the effect of replacing the 0–1 contiguity coefficients by a piecewise linear function,

$$w_{i,i'} = \begin{cases} 1, & d_{i,i'} \leq 60 \\ (120 - d_{i,i'})/60, & 60 < d_{i,i'} \leq 120 \\ 0, & d_{i,i'} > 120, \end{cases} \quad (11)$$

of the distance  $d_{i,i'}$  (in km) between the administrative centers of each area and redefining the prior potential function  $U(z)$  as  $\sum_{i,i'} w_{i,i'} I[z_i = z_{i'}]$ , giving a modified Potts model with smoother spatial dependence and effectively larger neighborhoods.

We found that the posterior means for the area-specific risks are quite robust overall to the modification of the Potts model, as can be seen by comparing Figures 11 and 7. Nonetheless, differences can be seen for some areas, and it is interesting to explore these further. For Block4, the largest difference of estimated risks between the two maps was .2, occurring for an area in the middle of the central elevated block. This area had a low SMR of .96, even though its true risk was 2. With the standard Potts model, the posterior mean risk for this area increases to 1.12, because of the influence of the contiguous neighbors with elevated risks. On the other hand, with the modified model, the influence of the further-distant areas with a low true risk of .7 predominates, and the posterior mean risk for this area drops to .92. A similar story is seen for Gradient NS. The largest difference of risk between the two maps, .21, occurs for an area on the northern border, the estimated risk for that area being pulled down when using the modified model. For both datasets, the DIC was lower and the RAMSE was slightly smaller for the standard model. As expected, changing the model influences the posterior distribution of  $k$ . For example, mixtures with fewer components were fitted with the modified model in the Gradient NS case, reflecting the flexible interaction between the type of spatial structure of the data, the chosen spatial model, and the estimated mixture characteristics. This modest investigation has thus indicated that

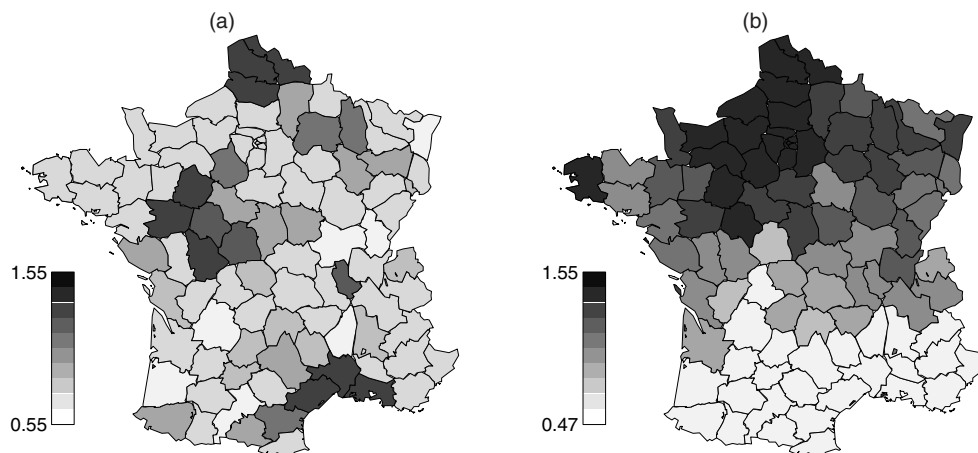


Figure 11. Mixture Estimate of the Risks for Block4 (a) and Gradient NS (b) Computed Under the Modified Potts Model Based on Eq. (11).

the model shows good adaptivity. A full investigation of the choice of spatial model for the allocations is beyond the scope of this article.

#### 4.7 Comparison With BYM

We first comment briefly on the results of a small simulation study conducted on the basis of the three underlying risk patterns described previously, together with three other setups. In the first one, “CAR”, the (log) risks follow an intrinsic autoregressive Gaussian model with zero mean based on the contiguity matrix  $W$ , where  $w_{ij} = 1$  if  $i \sim j$ . To be precise, the joint distribution of  $\{\log(\lambda_i)\}$  is simulated with covariance matrix  $.16$  times the generalized inverse of the matrix  $\text{diag}(v_1, \dots, v_n) - W$ , where  $v_i$  is the number of neighbors of area  $i$ , and the constant  $.16$  is used to scale appropriately the variability of the risk across the map. The corresponding maps of true risk and SMR are displayed in Figure 12. This setup was chosen to correspond to the spatial model underlying the BYM analysis. The last two examples have no spatial pattern: “Flat” refers to risks displaying no trend or spatial correlation, the “null hypothesis” for an epidemiologist, taken here as drawn at random from a uniform  $(.9, 1.1)$  distribution, and “Over” simulates overdispersed risks using a gamma mixture of Poissons chosen so that  $\text{var}(y) = 1.5 E(y)$ . For those last three datasets, the expected number of deaths are again those of the dataset on larynx cancer for women.

The spatial mixture and the BYM are compared on the basis of RAMSE, RAMSEL, DIC,  $E(D)$ , and  $p_D$ . Each line of Table 1 corresponds to the criterion averaged over five independent Poisson replications of the data pattern (the first replication for four of the datasets having been displayed previously). Note that out of the six risk patterns considered, only the first two have discontinuities. The comparisons are thus designed to investigate the versatility of the two models in recovering risk patterns for which they are not necessarily well adapted.

The table shows that the spatial mixture model gives more faithful estimation of the underlying risks, with smaller RAMSE in most cases. The RAMSEL criterion gives a similar picture except for the Over dataset, where it is smaller

for BYM. As could be expected, the difference is accentuated when there are contrasting zones, as in Block4 or North-South, or when the pattern is fairly uniform, as in Flat. When the risks are smoothly varying, as in Gradient NS or CAR, the two models give similar results; it is perhaps surprising that the mixture model performs competitively. Moreover, the spatial structure of the  $\{\lambda_i\}$  induced by the chosen allocation process leads to a more parsimonious model with consistently lower DIC and  $p_D$ .

If one accepts the DIC principle at face value, then the BYM model overfits the data in five of the six cases. With the exception of Block4, the posterior deviance under the BYM model is substantially smaller, but the  $p_D$  is so much greater that the DIC is at least as large. In the balance between recovering the true underlying scene and fitting the data, the mixture model is clearly less influenced by the noise in the data than the BYM model and is able to effect some spatially adaptive smoothing in a variety of situations.

Complementing the results of Table 1, we display the map of posterior means of  $\{\lambda_i\}$  estimated by the BYM model for two datasets (Block4 and North-South) in Figures 13(a) and 14(a). Comparing these with those corresponding to our mixture model (Figs. 6 and 7), against the maps of the simulated risks (Figs. 2 and 3) reveals evidence of remaining noise, unsmoothed by the analysis. This illustrates what was quantified by the two MSE criteria in Table 1, that recovery of the unblurred picture is less effective for the BYM model. It is also interesting to see that the posterior variability of the risks in both datasets is quite different from that of the mixture model, whether one compares the posterior standard deviations for the North-South (Figs. 6 and 14) or the histograms for Block4 (Figs. 8 and 13). For the BYM model, variability is not increased for areas along discontinuities. Further, the variability of the risks is higher overall and is closely linked to the size of the risks. This phenomenon is also true for smoothly varying risks, as in the Gradient NS example (results not shown).

Finally, Figure 15 displays for each area the root MSE between the simulated and estimated  $\{\lambda_i\}$  corresponding to the mixture model (a) or the BYM model (b) for the North-South dataset. For the mixture model, the highest errors are

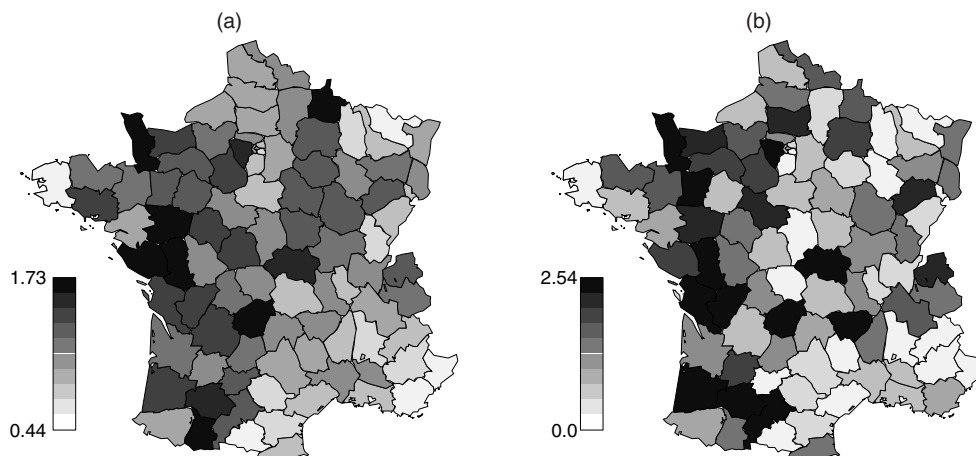


Figure 12. True Risks (a) and Observed SMRs (b) for the CAR Dataset.

Table 1. Simulation Results Comparing Spatial Mixture and BYM Models

Datasets	RAMSE		RAMSEL		DIC		E(D y)		p <sub>D</sub>	
	MIX	BYM	MIX	BYM	MIX	BYM	MIX	BYM	MIX	BYM
Block4	.22	.27	.22	.30	118.2	138.8	89.5	91.1	28.7	47.7
North-South	.15	.22	.15	.22	116.2	124.1	97.8	87.0	18.4	37.1
Gradient NS	.22	.24	.27	.26	125.4	129.7	94.6	86.6	30.8	43.1
CAR	.27	.28	.27	.27	133.8	136.7	102.1	95.1	31.7	41.6
Flat	.09	.19	.09	.19	93.4	108.6	89.4	77.2	4.0	31.4
Over	.23	.26	.24	.21	127.5	128.3	112.2	92.1	15.3	36.2

along the discontinuity, whereas for the BYM model, this pattern is less clear, and higher errors can be seen over all of the southern areas.

### 5. EPIDEMIOLOGIC APPLICATION TO DISEASE DATA

The performance of the model is illustrated on data concerning larynx cancer mortality in France at the level of the 94 mainland French départements reported by Rezvani, Molié, Doyon, and Sancho-Garnier (1997) for the period 1986–1993. For this dataset, we also illustrate how the introduction of area-level covariates in the model, as set out in (2), reduces the spatial heterogeneity of the risks.

The update of the regression parameters  $\gamma_j$  in (2) was performed using random walk Metropolis, other updates being as described earlier, with the necessary adjustment to the likelihood. Acceptance rates were around 49%, using normally distributed perturbations with standard deviations 1.5 times the reciprocal of the range of the corresponding covariate.

Laryngeal cancer is rare in women, with the observed number of deaths per area in this dataset ranging from 0 to 148 and SMR ranging from 0 to 2.1. The epidemiology of this cancer site has been studied mainly in men, in whom such risk factors as smoking, alcohol consumption, dietary factors, and specific occupational exposure have been brought to light in case-control studies (Austin and Reynolds 1996). Here we

include two covariates in our model: the per capita sales of cigarettes in 1975 (an available proxy for smoking) and an indicator of the urbanization of the area as recorded in the 1975 census. These variables are both time-lagged to allow for a delay between putative exposure and disease. Under the Potts mixture model adjusted without covariates, the posterior distribution of  $k$  peaks at  $k = 2$ , with  $p(k|y) = 0, .49, .23, .12$  for  $k = 1, 2, 3, 4$ .

Regions of higher risk are apparent in the north and south east. Even though the posterior means of the risks are fairly similar for the spatial mixture and the BYM model (Fig. 16), there is markedly higher variability for the BYM model (Fig. 17), a phenomenon described previously.

When the two covariates are introduced in the Potts mixture model, the posterior distribution of  $k$  shifts markedly toward the left, with  $p(k|y) = .31, .39, .18, .07$  for  $k = 1, 2, 3, 4$ . As expected, inclusion of the two covariates has partially explained the heterogeneity. In fact, the range of the posterior mean of the residual risks [ $\lambda$ 's in (2)] is substantially reduced to approximately one-third of the range of the risks displayed in Figure 16.

This example supports our view that by combining information from the mixture structure and the display of the posterior distribution of the risk estimates, one can characterize the heterogeneity of the risks and investigate how this heterogeneity is affected by the introduction of covariates. The

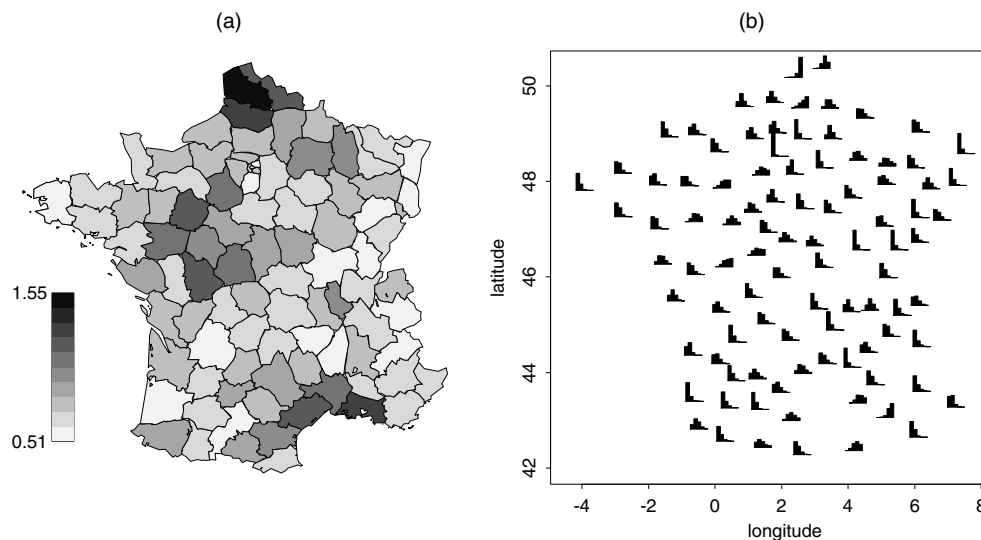


Figure 13. BYM Estimates of the Risks for Block4: Posterior Means (a) and Histogram of the Risks (b).

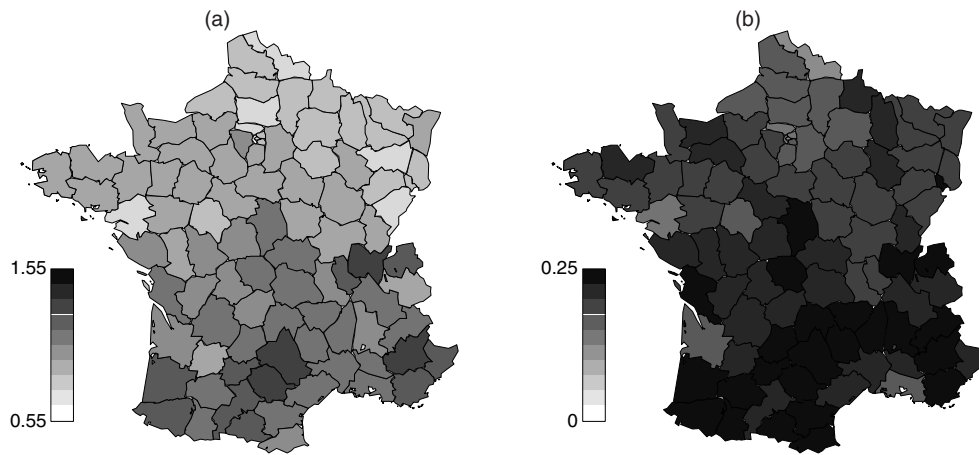


Figure 14. BYM Estimates of the Risks for North-South: Posterior Means (a), Posterior Standard Deviations (b).

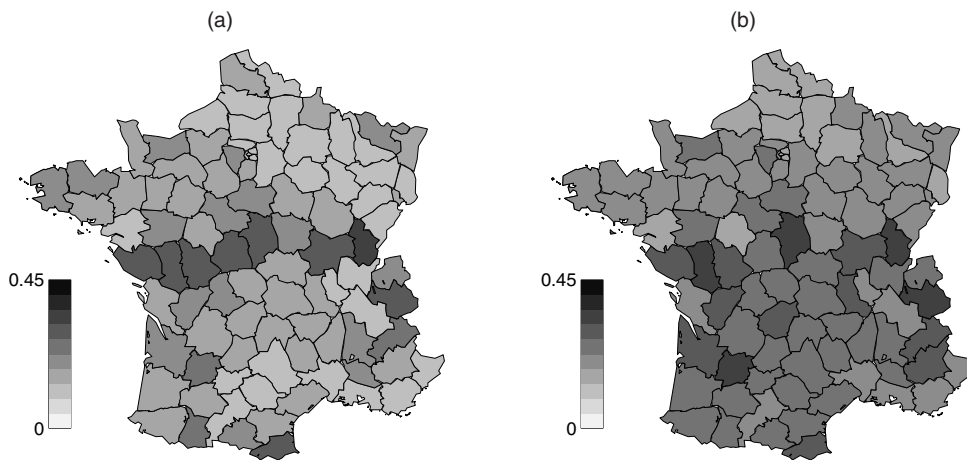


Figure 15. Root Mean Squared Errors for the North-South Dataset: Comparison Between the Mixture Model (a) and the BYM Model (b).

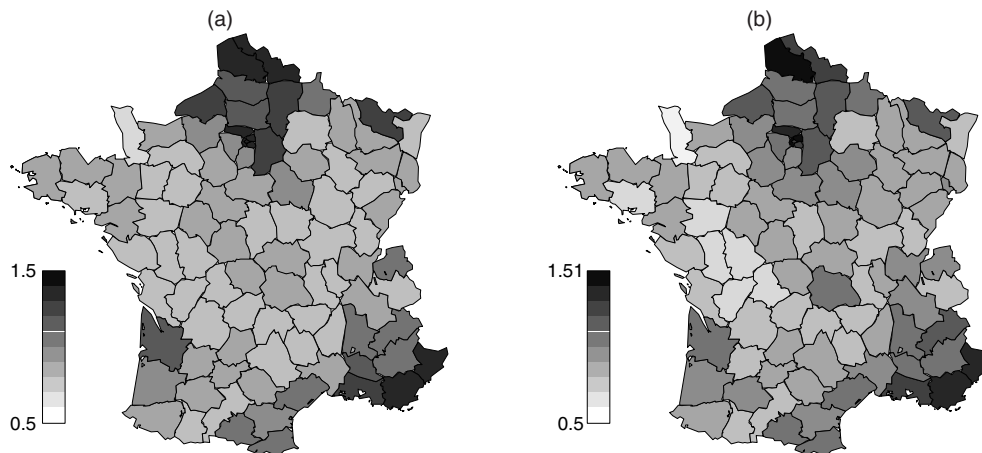


Figure 16. Larynx Cancer Mortality: Comparison of Posterior Means for the Risks Obtained Using the Mixture Model (a) and the BYM Model (b).



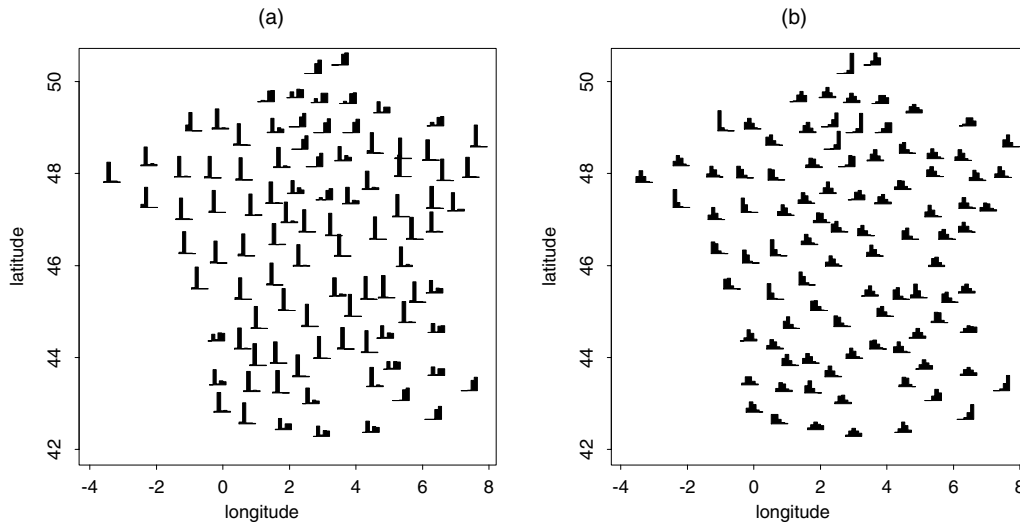


Figure 17. *Larynx Cancer Mortality: Comparison of the Posterior Distribution of the Risks Obtained Using the Mixture Model (a) and the BYM Model (b).*

ultimate goal, from a public health standpoint, is to uncover sets of covariates that can account for most of the geographic variability of the risks.

## 6. CONCLUDING REMARKS

Spatially structured heterogeneity is a common phenomenon that can be tackled in a hierarchical framework using a hidden Markov random field approach, through direct multivariate specification of the underlying field, or via partition models. Here, focused on rare outcomes and hierarchical Poisson models, we have proposed a new model within the hidden Markov random field framework. Whether the features that this model offers are useful depends on the purpose of the modeling exercise. In our study of the performance of the allocation model, we have been strongly influenced by the specific epidemiologic context. Recovery of the underlying Poisson rates is often the prime object of inference, supplemented by quantification of the extent of the underlying variability and probability statements about areas of high risk.

In terms of recovery of the “hidden state” or “true image,” we have shown that the posterior mean risks estimated under the allocation model give a faithful representation of the true risks in a wide variety of situations, encompassing both smooth and discontinuous cases. This flexibility is important, because there is usually little prior information on the underlying risks. Further, we have seen that information from the posterior estimate of the mixture structure can also be helpful in characterizing the spatial variability. We have explored some other features of the joint distribution of the risks, such as the number of clusters, but will leave a more in-depth study to further work.

Although these do not receive much space in this article, we have conducted numerous checks on the correctness of our MCMC samplers, and on the adequacy of their performance. In particular, we have verified that the prior distribution is recovered if we implement our computations without likelihood or data. We are satisfied that the range chosen for the interaction parameter  $\psi$  allows substantial spatial dependence

in the allocations, so that higher values of  $\psi$ , for which mixing could be slower, are not necessary. Overlong runs have been used on purpose in our examples, the algorithm being quite fast; our sampler, coded in Fortran, makes about 900 sweeps per second on a 300-MHz PC.

Further work on model comparison is certainly needed; we regard our comparisons with the BYM model as quite limited, in terms of both the scope and the criteria chosen for comparison. It would be interesting to extend comparisons to include the models proposed by Knorr-Held and Raßer and by Fernández and Green, and also to include other non-Gaussian Markov random field approaches. Indeed, when discontinuities are expected, it is advisable to replace the quadratic potential, leading to a Gaussian prior for the spatially structured random effects,  $u_i$ , by an absolute value difference potential. This has been discussed by Besag et al. (1991) and used by Best et al. (1999), and is related to smoothing using medians instead of means. Regarding the relevant criteria for comparison and choice, Bayesian model comparison is an active area of research, and there is much debate on the most appropriate approach.

As we mentioned at the end of Section 2, the model could be extended in several directions. One interesting avenue is the inclusion of covariates with heterogeneous effect, that is, where covariate effects and allocations interact. We have implemented the model outlined in (2), where the effect of the covariates is homogeneous. A Poisson mixture model with interaction between allocations and covariates but independent, nonspatial allocations has been discussed by Viallefont, Richardson, and Green (2002), and there should be no obstacle to extending this to spatially correlated allocations. Another extension that could be considered is to combine the spatial mixture with a BYM model, in effect replacing  $\lambda_0$  by  $\lambda_{z_i}$  in Section 3.4. It remains to be seen whether this could usefully combine the best features of both models, or whether gross overparameterization will result. A model containing a mixture of Gaussian and non-Gaussian (median-based) conditional autoregressive components was recently proposed by Lawson

and Clark (2001). Finally, there is great scope for extensions to spatiotemporal modelling, both for epidemiologic applications and more generally.

[Received July 2001. Revised April 2002.]

## REFERENCES

- Austin, D. F., and Reynolds, P. (1996), "Laryngeal Cancer," in *Cancer Epidemiology and Prevention* (2nd ed.), eds. D. Schottenfeld and J. F. Fraumeni, Oxford, UK: Oxford University Press, pp. 619–636.
- Aykroyd, R. G., and Zimeras, S. (1999), "Inhomogeneous Prior Models for Image Reconstruction," *Journal of the American Statistical Association*, 94, 934–946.
- Besag, J. (1986), "On the Statistical Analysis of Dirty Pictures" (with discussion), *Journal of the Royal Statistical Society, Ser. B*, 48, 259–302.
- Besag, J., York, J., and Mollié, A. (1991), "Bayesian Image Restoration With Applications in Spatial Statistics" (with discussion), *Annals of the Institute of Mathematical Statistics*, 43, 1–59.
- Best, N. G., Arnold, R. A., Thomas, A., Waller, L. A., and Conlon, E. M. (1999), "Bayesian Models for Spatially Correlated Disease and Exposure Data," in *Bayesian Statistics 6*, eds. J. M. Bernardo, J. O. Berger, A. P. Dawid, and A. F. M. Smith, Oxford, UK: Oxford University Press, pp. 131–156.
- Carroll, R. J., Roeder, K., and Wasserman, L. (1999), "Flexible Parametric Measurement Error Models," *Biometrics*, 55, 44–54.
- Clayton, D., and Bernardinelli, L. (1992), "Bayesian Methods for Mapping Disease Risk," in *Geographical and Environment Epidemiology: Methods for Small Area Studies*, eds. P. Elliott, J. Cuzick, D. English, and R. Stern, Oxford, UK: Oxford University Press, pp. 205–220.
- Clifford, P. (1986), Discussion of "On the Statistical Analysis of Dirty Pictures" by J. Besag, *Journal of the Royal Statistical Society, Ser. B*, 48, 284.
- Denison, D. G. T., and Holmes, C. C. (2001), "Bayesian Partitioning for Estimating Disease Risk," *Biometrics*, 57, 143–149.
- Elliott, P., Wakefield, J. C., Best, N. G., and Briggs, D. J. (2000), *Spatial Epidemiology: Methods and Applications*, Oxford, UK: Oxford University Press.
- Fernández, C., and Green, P. J. (2002), "Modelling Spatially Correlated Data via Mixtures: A Bayesian Approach," *Journal of the Royal Statistical Society, Ser. B*, 64, to appear (part 4).
- Gelman, A., and Meng, X.-L. (1998), "Simulating Normalizing Constants: From Importance Sampling to Bridge Sampling to Path Sampling," *Statistical Science*, 13, 163–185.
- Green, P. J. (1995), "Reversible Jump Markov Chain Monte Carlo Computation and Bayesian Model Determination," *Biometrika*, 82, 711–732.
- Greenland, S., and Robins, J. (1994), "Ecological Studies: Biases, Misconceptions, and Counterexamples," *American Journal of Epidemiology*, 139, 747–760.
- Johnson, V. E. (1994), "A Model for Segmentation and Analysis of Noisy Images," *Journal of the American Statistical Association*, 89, 230–241.
- Knorr-Held, L., and Raßer, G. (2000), "Bayesian Detection of Clusters and Discontinuities in Disease Maps," *Biometrics*, 56, 13–21.
- Knorr-Held, L., and Rue, H. (2002), "On Block Updating in Markov Random Field Models for Disease Mapping," *Scandinavian Journal of Statistics*, 29, in press.
- Künsch, H. R. (2001), "State Space and Hidden Markov Models," in *Complex Stochastic Systems*, eds. O. E. Barndorff-Nielsen, D. R. Cox, and C. Klüppelberg, London: Chapman and Hall/CRC, pp. 109–173.
- Lawson, A. B., and Clark, A. (2002), "Spatial Mixture Relative Risk Models Applied to Disease Mapping," *Statistics in Medicine*, 21, 359–370.
- McLachlan, G., and Peel, D. (2000), *Finite Mixture Models*, Chichester, UK: Wiley.
- Mollié, A. (1996), "Bayesian Mapping of Disease," in *Markov Chain Monte Carlo in Practice*, eds. W. Gilks, S. Richardson, and D. J. Spiegelhalter, London: Chapman and Hall, pp. 359–379.
- Ogata, Y., and Tanemura, M. (1984), "Likelihood Analysis of Spatial Point Patterns," *Journal of the Royal Statistical Society, Ser. B*, 46, 496–518.
- Rezvani, A., Mollié, A., Doyon, F., and Sancho-Garnier, H. (1997), *Atlas de la Mortalité par Cancer en France, Période 1986–1993*, Paris: Editions INSERM.
- 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., Rydén, T., and Titterton, D. M. (2000), "Bayesian Inference in Hidden Markov Models Through the Reversible Jump Markov Chain Monte Carlo Method," *Journal of the Royal Statistical Society, Ser. B*, 62, 57–75.
- Roeder, K., Carroll, R. J., and Lindsay, B. G. (1996), "A Semiparametric Mixture Approach to Case-Control Studies With Errors in Covariables," *Journal of the American Statistical Association*, 91, 722–732.
- Rydén, T., and Titterton, D. M. (1998), "Computational Bayesian Analysis of Hidden Markov Models," *Journal of Computational and Graphical Statistics*, 7, 194–211.
- Stephens, M. (2000), "Bayesian Analysis of Mixture Models With an Unknown Number of Components—An Alternative to Reversible Jump Methods," *Annals of Statistics*, 28, 40–74.
- Spiegelhalter, D. J., Best, N. G., Carlin, B. P., and van der Linde, A. (2002), "Bayesian Measures of Model Complexity and Fit," *Journal of the Royal Statistical Society, Ser. B*, 64, to appear.
- Titterton, D. M., Smith, A. F. M., and Makov, U. E. (1985), *Statistical Analysis of Finite Mixture Distributions*, Chichester, UK: Wiley.
- Tjelmeland, H., and Besag, J. (1998), "Markov Random Fields With Higher-Order Interactions," *Scandinavian Journal of Statistics*, 25, 415–433.
- Viallefont, V., Richardson, S., and Green, P. J. (2002), "Bayesian Analysis of Poisson Mixtures," *Nonparametric Statistics*, 14, 181–202.
- Wakefield, J., and Morris, S. (1999), "Spatial Dependence and Errors-in-Variables in Environmental Epidemiology," in *Bayesian Statistics 6*, eds. J. M. Bernardo, J. O. Berger, A. P. Dawid, and A. F. M. Smith, Oxford, UK: Oxford University Press, pp. 657–684.