



## A comparison of Bayesian methods for flexible modelling of spatial risk surfaces in disease mapping

Journal:	<i>Biometrical Journal</i>
Manuscript ID:	bimj.201200176
Wiley - Manuscript type:	Research Paper
Date Submitted by the Author:	06-Sep-2012
Complete List of Authors:	Sturtz, Sibylle; Institute for Quality and Efficiency in Health Care, Medical Biometry Ickstadt, Katja; TU Dortmund University, Statistics
Keywords:	Disease mapping, Poisson/gamma random field model, Spatial statistics

SCHOLARONE™  
Manuscripts

## A comparison of Bayesian methods for flexible modelling of spatial risk surfaces in disease mapping

Sibylle Sturtz<sup>\*1,2</sup> and Katja Ickstadt<sup>2</sup>

<sup>1</sup> Department of Medical Biometry, Institute for Quality and Efficiency in Health Care, Im Mediapark 8, 50670 Köln, Germany.

<sup>2</sup> Faculty of Statistics, TU Dortmund University, 44221 Dortmund, Germany

Received zzz, revised zzz, accepted zzz

Bayesian hierarchical models usually model the risk surface on the same arbitrary geographical units for all data sources. Poisson/gamma random field models overcome this restriction as the underlying risk surface can be specified independently to the resolution of the data. Moreover, covariates may be considered as either excess or relative risk factors. We compare the performance of the Poisson/gamma random field model to the Markov random field (MRF)-based ecologic regression model and the Bayesian Detection of Clusters and Discontinuities (BDCD) model, in both a simulation study and a real data example. We find the BDCD model to have advantages in situations dominated by abruptly changing risk while the Poisson/gamma random field model convinces by its flexibility in the estimation of random field structures and by its flexibility incorporating covariates. The MRF-based ecologic regression model is inferior. WinBUGS Code for Poisson/gamma random field models is provided.

**Key words:** Disease mapping, Poisson/gamma random field model, Spatial statistics

Supporting Information for this article is available from the author or on the WWW under <http://dx.doi.org/10.1022/bimj.XXXXXXX>

### 1 Introduction

Routinely available disease and population data used for disease mapping are typically reported as counts per administrative areas such as counties or electoral units. Bayesian hierarchical spatial models are widely used to model the underlying risk surface and produce disease maps. These typically model the risk surface on the same geographical units as the data are measured on. Examples are the Markov random field (MRF)-based model by Besag et al. (1991) further developed to an ecologic regression model by Clayton and Bernardinelli (1992) and the Bayesian Detection of Clusters and Discontinuities (BDCD) model (Knorr-Held and Raßer, 2000). However, administrative areas are arbitrary and there is no reason to believe that disease risk truly varies on this scale. Furthermore, the assumption of spatial dependence based on neighborhoods of administrative units not necessarily reflects environmental processes. The assumption of a smooth underlying risk surface from which area-specific risks are derived is a more reasonable one. For example, Kelsall and Wakefield (2002) model relative risks by a Gaussian random field model to overcome spatial aggregation and modeling dependencies based on neighboring administrative units. Within-area variability of covariates is incorporated into an ecologic regression model to reduce ecological bias caused by spatial aggregation by Haining et al. (2010).

\*Corresponding author: e-mail: sibylle.sturtz@iqwig.de, Phone: +00 49 221 35685456, Fax: +00 49 221 35685896

Wolpert and Ickstadt (1998) introduce a random field generalization of the Poisson/gamma hierarchical model by Clayton and Kaldor (1987) allowing for an underlying risk surface specified on a spatial resolution that is independent of that on which the data are measured. This model, generalized by Best et al. (2000) for an application in epidemiology, also allows for covariates to improve the estimation of the risk surface which may be modelled either as excess or relative risk factors (Breslow and Day, 1980) leading to different interpretations. The additive influence of excess risk factors offers alternative explanations of an event and is preferable for competing, non-interacting effects. The multiplicative influence of relative risk factors reflects different susceptibilities to a covariate.

Models such as the MRF-based ecologic regression model by Besag et al. (1991) and Clayton and Bernardinelli (1992) usually treat covariates as relative risk factors only, though more flexible functional forms are possible (Wakefield, 2007; Natario and Knorr-Held, 2003). The BDCD model only allows for categorical covariates introduced as relative risk factors (Giudici et al., 2000) to improve the estimation of the risk surface.

Often, such a restriction on categorical covariates is too limited to provide a good model fit for the given data. An example for such a data set is provided by Best et al. (2001) where childhood leukemia data given on a ward level are analyzed in dependence of a continuous variable, namely benzene emissions (in tonnes per year). This data set also forms the basis for this paper.

Best et al. (2005) compare five spatial models in the context of disease mapping within a Bayesian framework. Those represent three different classes which 1) either use correlated normal priors such as the MRF-based ecologic regression model, 2) are semi-parametric spatial models such as the BDCD model or 3) are spatial moving average models as the Poisson/gamma model. This Poisson/gamma model is a discrete version of the Poisson/gamma random field model considered here and includes covariates only additively. Due to the discretization they base their analysis on a fixed spatial partition. Best et al. (2005) restrict their analyses to only one specific spatial pattern including a covariate and two fixed hypothetical point sources.

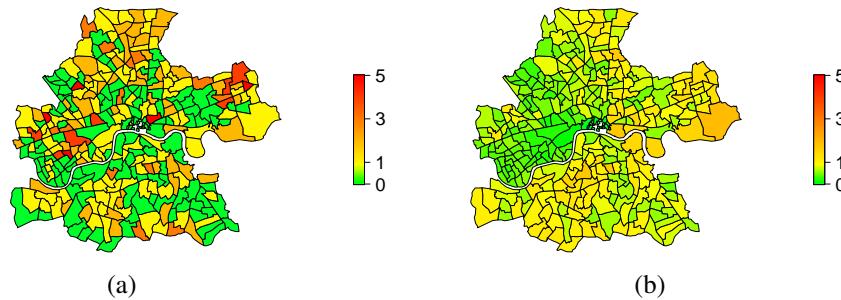
The aim of this paper is to compare the performance of the random field formulation of Poisson/gamma models to the MRF-based ecologic regression model and the BDCD model with regard to disease mapping (smoothing) both in a simulation study and in a real data example.

We overcome the restrictions of the work by Best et al. (2005) by the design of our study which is chosen to investigate the suitability of these models for different spatial patterns including smooth surfaces, sharp borders and clusters as well as for covariates modelled as either excess or relative risk factors. We also account for latent risks in both data generation and modelling to study the models' behavior with or without accounting for unobserved covariates. Characteristics of unobserved covariates such as their location and intensity are treated to be uncertain and estimated within our implementation.

Motivated by the example of Best et al. (2001), we simulate incidence rates of childhood leukemia using observed benzene data and population counts for Inner London. While benzene is aggregated to ward level, we model the latent covariate as a random field without choosing a specific spatial resolution in advance. Model comparison is based on the Deviance Information Criterion (Spiegelhalter et al., 2002, DIC) as well as on the Mean Square Error. We also analyze the real data set of leukemia rates in dependence to benzene emissions.

Furthermore, we focus on the implementation of Poisson/gamma random field models into the Bayesian software WinBUGS (Lunn et al., 2000), a user friendly software that is very popular for Bayesian applications. Selected spatial models such as the widely used MRF-based ecologic regression model are already implemented and ready to use which is one reason for their popularity. The implementation of Poisson/gamma random field models as presented in Appendix A will provide an alternative for an easily accessible modelling framework.

The outline of this paper is as follows: The underlying data set is described in the second section. The third section is dedicated to the Poisson/gamma random field models, the fourth one to the alternative models. In the fifth and sixth section the design and the results of the performed simulation study are



**Figure 1** Number of observed (a) and expected cases (b) in the period of 1985–1996.

presented, in the seventh section the results for the real data set are discussed. Section eight summarizes the results and discusses further developments.

## 2 Data

Incident cases of leukemia are registered by the Thames Cancer Registry for the period from 1985 until 1996. In the whole period 295 cases of cancer in children under 15 years old are registered in the area of Inner London. Figure 1(a) describes the spatial distribution of the observed leukemia cases  $O_i, i = 1, \dots, n$ , in the study region comprising  $n = 310$  electoral wards.

Population counts  $N_i, i = 1, \dots, n$ , are also given on ward level. They are available from 1981 and 1991 censuses and are stratified by age (0-4, 5-9 and 10-14 years) and sex. For intercensal years, counts are interpolated accounting for demographic changes. Based on the national leukemia rate  $r_{st}$  for age-sex stratum  $s$  and year  $t$  we calculate the expected numbers of cases  $E_i$  in ward  $i$  by

$$E_i = \sum_s \sum_t r_{st} N_{ist}, \quad (1)$$

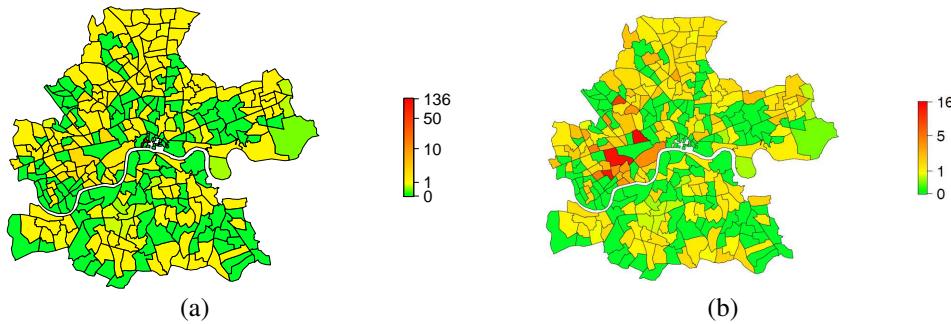
see Figure 1(b).

Given both, the number of observed cases  $O_i$  and expected cases  $E_i$ , we calculate the standardized mortality ratio (SMR) for region  $i, i = 1, \dots, n$ , by

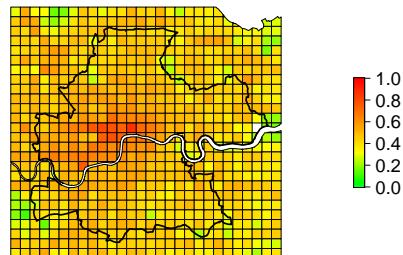
$$\text{SMR}_i = \frac{O_i}{E_i}.$$

The resulting spatial pattern is inhomogeneous, regions of low or high risk are difficult to characterize, cf. Figure 2. Furthermore, for one ward in the center we obtain a striking value of  $\text{SMR}_8 = 135.1$  due to a low population density resulting in 0.0074 expected cases but one registered case. One goal of spatial epidemiology which we also focus on in this paper is a model-based smoothing of this surface to allow for a reasonable interpretation as well as for the identification of regions with increased or lowered risk.

As discussed by Best et al. (2001) occupational exposure to benzene is accepted to be associated with an increased risk of various leukaemias. However, the effect of low level exposures to environmental benzene is unknown. Benzene emissions are therefore considered as a covariate. Emissions are available on 1 km ×



**Figure 2** Calculated SMRs for the leukaemia data set (a). In (b) the increased SMR of 135.1 in the central region 8 of (a) is set to be NA (white) for a better presentation of the risk surface.



**Figure 3** Benzene exposure data on 1 km  $\times$  1 km grid cells.

1 km grid cells from the atmospheric emissions inventory for London (Buckingham et al., 1997). Benzene is not monitored directly but modelled as

$$\text{activity rate} \times \text{emission factor} = \text{emission rate}$$

where the emission factor depends on the source of pollutants which may be modelled traffic flows, petrol stations, and commercial, residential, or industrial combustion processes. The modelled atmospheric emissions for Inner London are presented in Figure 3.

### 3 Poisson/gamma random field model

Poisson/gamma random field models can be thought of a flexible extension of the conjugate Poisson/gamma models introduced by Clayton and Kaldor (1987). The random field formulation allows for location and observation specific covariates that may be modeled as excess or relative risk. More specifically, Ickstadt and Wolpert (1999) consider locations  $y \in \mathcal{Y}$  and attributes  $a \in \mathcal{A}$  which may be individual or environmental covariates, as marked points  $x = (y, a)$  on a set  $\mathcal{X} = \mathcal{Y} \times \mathcal{A}$  in Euclidean space. The number of

points  $N(dx) = N(dy da)$  is then modelled as a marked Poisson process on  $\mathcal{X}$  with mean  $\Lambda(x)w(dx)$ , i.e.,  $N(dy da) \sim \text{Pois}(\Lambda(y, a)w_Y(dy)w_A(da))$ , where  $w(dx)$  is the overall reference measure with reference measure  $w_Y(dy)$  of spatial covariates and  $w_A(da)$  of the attributes.

The reference weight prior  $w_A(da)$  of the attributes is usually chosen to be space independent, i.e.  $w_A(da|y_1) = w_A(da|y_2) = w_A(da)$  for  $y_1 \neq y_2$ , and set to  $w_A(da) \equiv 1$  but can also be location specific. Choosing  $w_A(da)$  to be location independent gives  $w_Y(dy)$  the role of a population reference measure in disease mapping (Ickstadt and Wolpert, 1999).

The mean  $\Lambda(y, a)$  depends on a set  $J_A$  of excess risk factors that are introduced additively and a multiplicatively modeled set  $J_M$  of relative risk factors.

To allow for non-measured covariates, we introduce a latent covariate  $X_*$  and corresponding correlation coefficient  $\theta_*$ .  $X_*$  is modeled as a kernel mixture of a random measure  $\Gamma(ds)$  on a space  $\mathcal{S}$  that may coincide with or may be larger than  $\mathcal{Y}$  representing possible locations of latent covariates within or perhaps outside the area of observed marked points. Here we choose  $\mathcal{S}$  to be the bounding box of the area of Inner London to account for an influence of covariates inside and outside the city limits. We use an inhomogeneous Gamma random field  $\Gamma(ds)$  with shape measure  $\alpha^\theta(ds)$  and inverse scale function  $\beta^\theta(ds) > 0$ , possibly depending on additional parameters  $\theta$ .  $\Gamma(ds)$  consists of infinitely many latent sources  $s, s \in \mathcal{S}$ , located at  $\mu_s = (\mu_1, \mu_2)$  with magnitudes  $\gamma_s$ . This gamma random field is combined with a bivariate Gaussian kernel  $k(y, s)$  with uncorrelated longitude and latitude, i.e.,

$$k(y, s) = \exp \left\{ -\frac{1}{2} \left( \frac{(y_1 - \mu_1)^2}{\sigma_1^2} + \frac{(y_2 - \mu_2)^2}{\sigma_2^2} \right) \right\}. \quad (2)$$

Thus, the influence of the latent covariate at location  $y$  is given by the kernel mixture  $X_* = \int_{\mathcal{S}} k(y, s)\Gamma(ds)$  and due to the discrete representation of the gamma random field can be rewritten as  $\sum_{s \in \mathcal{S}} k(y, s)\gamma_s$ .

The rate  $\Lambda(y, a)$  is then modelled in a Bayesian framework with  $\pi(\theta)$  denoting the prior distributions on all the parameters  $\theta$ , leading to the following three stage hierarchical Bayesian Poisson/gamma random field model:

Level 1: Marked points  $N(dx) \sim \text{Pois}(\lambda(x)w(dx)), x = (y, a) \in \mathcal{X} \equiv \mathcal{Y} \times \mathcal{A}$

$$\begin{aligned} \text{Intensity} \quad \Lambda(x) &= \left( \theta_0 + \sum_{j \in J_A} a_j \theta_j + \sum_{s \in \mathcal{S}} k(y, s) \gamma_s \theta_* \right) \\ &\quad \times \exp \left( \sum_{j \in J_M} a_j \theta_j \right) \end{aligned}$$

Level 2: Latent sources  $\Gamma(ds) \sim \text{Gamma}(\alpha^\theta(ds), \beta^\theta(ds))$

Level 3: Parameter  $\theta \sim \pi(\theta)d\theta$ .

For further details see Ickstadt and Wolpert (1999) and Best et al. (2000).

For an implementation in WinBUGS please note that the Poisson/gamma random field model is approximated by a finite sum  $\sum_{s \in \mathcal{S}} k(y, s)\gamma_s$  as follows. To ensure the full flexibility of the gamma random field we allow the center of each kernel  $\mu_s = (\mu_1, \mu_2)$  to be estimated by the model. We start with only one jump and corresponding kernel and successively add further jumps until the quality of the model fit as expressed by the DIC if accounted for MC error is not improved any more.

### 3.1 Prior settings

Prior distributions for all parameters need to be defined. The coefficients  $\theta_0$  and  $\theta_*$  are modelled as described above. The only other covariate is benzene with parameter  $\theta_{\text{benz}}$  which is treated to be either an excess risk factor in set  $J_A$  or a relative risk factor in set  $J_M$ . For the regression coefficients  $\theta_j$ ,  $j \in J = \{\theta_0, \theta_{\text{benz}}, \theta_*\}$ , we assume a Gamma distribution  $\text{Gamma}(\alpha, \tau)$  with density

$$f(\theta_j) = \begin{cases} \frac{\tau^\alpha \theta_j^{\alpha-1} \exp(-\tau\theta_j)}{\Gamma(\alpha)} & \text{if } \theta_j > 0 \\ 0 & \text{else} \end{cases}$$

and set the shape parameter equal to  $\alpha = 0.575$ . This corresponds to a ratio of the 90th/10th percentile of the prior distribution of 100 reflecting a prior probability of 80% of the number of cases associated to each factor to lie between a 1/10th and 10 times the prior mean. The prior scale parameter  $\tau$  is chosen to obtain a prior mean assuming an equal amount of association for each covariate. Since the intensity  $\Lambda$  depends on the ratio of (the number of observed cases)/(the number of expected cases) we use the prior mean  $(\sum_i N_i) / (|J| \sum_i E_i)$  for regression coefficient  $\theta_j$ .

For the latent magnitudes  $\gamma_m$  of kernel mixtures we use a Gamma distribution as well. Here we choose

$$\gamma_m \sim \text{Gamma}(\alpha_\gamma, \tau_\gamma)$$

with  $\alpha_\gamma$  and  $\tau_\gamma$  chosen as

$$\begin{aligned}\alpha_\gamma &= |S| \times \tau_\gamma \\ \tau_\gamma &= \frac{1}{m}\end{aligned}$$

proportional to the area  $|S|$  of the bounding box  $S$  of the modelled region with  $m$  being the number of latent risk sources.

Due to their structure, random field models are consistent under aggregation of areas. This is also considered in the choice of the priors. We set the prior mean of the magnitudes to be  $|S|/m$ . Hence,  $\alpha_\gamma$  changes in proportion to the number of involved latent risk sources. The location of each kernel is assumed to be uncertain. For each kernel, we draw  $\mu_1$  and  $\mu_2$  uniformly within the bounding box  $S$ .

For the kernel  $k(y, s)$  we assume a Gaussian kernel as in equation 2 with uncertain variance parameter  $\rho$ . The prior distribution of  $\sigma_1^2 = \sigma_2^2 = \rho$  is chosen according to a log-Normal distribution with mean 0; the precision varies between 1 and 3. This choice was motivated by a prior study on model adequacy assuming fixed variances for Gaussian kernels.

We use two Markov chains for MCMC simulation with 100 000 iterations after a burn-in of 50 000.

#### 4 Alternative models

Besides Poisson/gamma random field models we also apply the MRF-based ecologic regression model (Besag et al., 1991, Clayton and Bernardinelli, 1992). Observed cases are assumed to follow a Poisson distribution where the Poisson rate depends on the number of expected cases. Covariates as well as spatially unstructured and structured terms are considered as relative risk factors, the latter ones based on a Markov random field (Besag et al., 1991).

We also apply the BDCD model (Knorr-Held and Raßer, 2000), a spatial partition model. It also assumes a Poisson distribution for the observed cases. Here, the rate depends on the expected cases as well as on a relative risk which is assumed to be constant across neighboring regions referred to as a cluster. The partition of the region into  $k$  clusters with  $k$  unknown is estimated by Reversible Jump MCMC methods (Green, 1995). An extension of this model allows for categorical covariates (Giudici et al., 2000), but not for continuous ones as can be modeled by Poisson/gamma random fields.

For details of prior settings for these two models see Sturtz (2007).

#### 5 Simulation scenarios

To compare the performance of the Poisson/gamma random field models with MRF-based ecologic regression and the BDCD models we set up a simulation study. Simulated data sets correspond to the goals of our study, i.e. identification and modeling of clustered structures as well as smoothly changing risk surfaces by the estimation of the latent field. Furthermore, we analyze the necessity for different covariate interpretations by generating scenarios assuming the covariate, i.e. benzene, to be either an excess or a relative risk factor.

**Table 1** Structures used for data generation:  $\times$  marks structures not described in detail in this paper, capital letters correspond to structures involving benzene additively (A) or multiplicatively (M).

$\oplus$ benzene		$\otimes$ benzene		
low	high	low	high	
$\times$	$\times$	M <sub>1</sub>	$\times$	benzene only
A <sub>2</sub>	$\times$	M <sub>2</sub>	M <sub>2</sub> <sup>high</sup>	+ latent risk covariate
$\times$	$\times$	M <sub>3</sub>	$\times$	+ linear trend covariate
$\times$	$\times$	M <sub>4</sub>	$\times$	+ increased risk in the southern part
$\times$	$\times$	M <sub>5</sub>	$\times$	+ increased risk in 3 clusters
$\Downarrow$	$\Downarrow$	$\Downarrow$	$\Downarrow$	$\Downarrow$
330	770	330	770	+ 330 cases

Following the ideas of experimental design, we vary all relevant factors simultaneously. We assume two different scenarios for this covariate, a ‘low’ risk scenario in which we assume benzene to account for about 330 observations similar to our real data application and a ‘high’ risk scenarios in which we more than double this amount making benzene responsible for 770 generated cases. Furthermore, we combine benzenes’ influence with different spatial patterns. The scenarios of the simulation study (summarized in Table 1) therefore include: i) structures determined by the influence of benzene only (e.g., M<sub>1</sub>); ii) structures assuming influence of benzene in combination with a Gaussian kernel (e.g., A<sub>2</sub>, M<sub>2</sub>, M<sub>2</sub><sup>high</sup>); iii) structures combining benzenes’ influence with a linearly decreasing trend from North to South (e.g., M<sub>3</sub>); iv) structures combining benzene with a plateau of increased risk in all wards south of the Thames (e.g., M<sub>4</sub>); v) structures assuming an influence of benzene in combination with increased risk in three randomly located clusters of different size (e.g., M<sub>5</sub>). Locations for the three cluster centers at the border of the study region, in the Center of London and clusters north of the Thames in scenario v) were selected at random provided they fulfilled the following criteria: a) the number of comprised wards are different; b) one cluster is divided by the river Thames; c) there are clusters at the border of Inner London and clusters completely surrounded by low risk regions; and d) distances between the clusters are different. For scenario ii) we choose the location and the variance such that the kernel has an influence in the south-western part of the study region only.

Parameters for additionally involved covariates are always chosen to account for another 330 cases. This leads to 20 different simulation scenarios which are summarized in Table 1. We focus our detailed presentation mainly on data sets assuming a multiplicative influence of benzene (M<sub>i</sub>,  $i = 1, \dots, 5$ ). This corresponds to the most frequently used scenario in spatial epidemiology as applied models usually support relative risk factors. Furthermore, we discuss the ‘low’ case scenario in more detail as this corresponds more to real world examples. In order to allow for a comparison with data sets assuming an additive influence of the covariate (such as A<sub>2</sub>) or a high multiplicative influence (e.g., M<sub>2</sub><sup>high</sup>) we exemplarily present the results for those models in the scenario of additionally generated risk represented by a Gaussian kernel, see Table 1. For a detailed discussion of the results of other structures marked by an  $\times$  see Sturtz (2007).

The number of cases  $O_i$  in region  $i$ ,  $i = 1, \dots, n$ , is generated from a Poisson distribution, i.e.,

$$O_i \sim \text{Pois}(\Lambda_i E_i)$$

with  $\Lambda_i$  calculated according to the simulation scenarios. We used the expected number of cases  $E_i$  adjusted for different age and sex distributions in each area as given in Best et al. (2001).

The Poisson/gamma random field models applied to the generated data sets include benzene either as an excess (model a) or a relative risk factor (model m) in order to determine whether we can identify the true

**Table 2** Summary of the DICs achieved for different structures, capital letters correspond to the underlying simulated structures, lower case letters to the applied Poisson/gamma random field model (benzene included as excess risk factor: model a; benzene included as relative risk factor: model m; benzene not included: model o), the MRF-based ecologic regression model (MRF) and the BDCD model (BDCD), bold numbers indicate best-fitting model.

Structure		model a	model m	model o	MRF	BDCD
A <sub>2</sub>	(+ 1 cluster)	<b>346.2</b>	+ 28.5	+ 38.2	+ 41.4	+ 61.8
M <sub>1</sub>	(benz only)	<b>321.9</b>	+ 0.9	+ 4.0	+ 0.5	+ 6.1
M <sub>2</sub>	(+ 1 cluster)	<b>373.2</b>	+ 0.3	+ 2.0	+ 13.7	+ 24.8
M <sub>2</sub> <sup>high</sup>	(+ 1 cluster)	+ 83.9	<b>341.1</b>	+137.3	+ 59.8	+ 141.1
M <sub>3</sub>	(+ smooth cov)	+ 1.4	<b>338.3</b>	+ 4.4	+ 14.7	+ 7.2
M <sub>4</sub>	(+ step)	+ 25.2	+ 25.5	+ 23.1	+ 28.7	<b>347.8</b>
M <sub>5</sub>	(+ 3 clusters)	+ 34.2	+ 34.2	+ 58.7	+ 68.1	<b>358.9</b>

interpretation. Additionally, we use a Poisson/gamma random field model that ignores information available from the covariate (model o). The gamma random field is modelled as discussed above. Additionally, we apply the BDCD algorithm where wards separated by the river Thames are assumed to be neighbors (BDCD model) and a MRF-based ecologic regression model with the same neighborhood and benzene as an relative risk factor.

Model comparison is done by the DIC and the Mean Square Error (MSE). We present the resulting DICs only as the MSE shows a corresponding behavior, see Sturtz (2007).

## 6 Results of simulations

### 6.1 Main conclusions

To summarize the results of our simulations study, we report the value of the best model, in combination with the differences between this and the DIC of every other model in Table 2. Our main conclusions are as follows:

Only Poisson/gamma random field models allow for additive modelling. Here, we find additive structures usually favoring an additive model while multiplicative structures lead to a multiplicative model as best fitting model. The amount of increase in the DIC compared to the best model varies. As we can see in Table 2 for multiplicative structures differences can be negligible (see structures M<sub>1</sub> and M<sub>2</sub>, for example), but can also be substantial. This holds especially for structures assuming a “high” influence of the covariate leading to an increase of 83.9 points in the DIC for structure M<sub>2</sub><sup>high</sup> when modelled by an additive Poisson/gamma random field model.

Covariates should be included in the model; if they are missed the DIC is increased. One example is again structure M<sub>2</sub><sup>high</sup> combined with Poisson/gamma random field model o where we observe an increase in the DIC of 137.3 points. The BDCD model without covariates has a DIC increased by 141.1 points compared to the best model. This holds for almost all structures analyzed with exceptions only for clustered structures for other reasons. In general, the increase in DIC values of Poisson/gamma random field models tends to be smaller compared to the BDCD model.

The estimation of the size of the covariate effects is of different quality. Additive structures in the low risk scenario ( $\theta_{\text{benz}} = 3$  in data generation) results in an estimate of  $\hat{\theta}_{\text{benz}}$  with a minimum of 2.202 and a maximum of 3.979 whereas in the high risk scenario ( $\theta_{\text{benz}} = 7$  in data generation) we estimate  $\min(\hat{\theta}_{\text{benz}}) = 5.834$ ,  $\max(\hat{\theta}_{\text{benz}}) = 6.706$ . The underestimation also holds for the high risk scenario for multiplicative models when latent terms are included in data generation. When only benzene or a latent

risk covariate is considered, the effect is reasonably well reproduced. As the complexity of the latent risk covariate increases, the accuracy of the estimate of the benzene covariate decreases. Additive structures require a higher number of flexible kernels which also account for the influence of benzene.

Typically, the inclusion of only five to ten latent risk sources is sufficient for a good model fit of the Poisson/gamma random field models. Again, the model corresponding to the underlying simulated structure is usually characterized by a lower number of kernels. Models not including benzene require a higher number of kernels to achieve a comparable model fit.

In our comparison, MRF models assume spatial dependencies most similar to multiplicative Poisson/gamma random field models. Nevertheless, their results are inferior. Most similar results are obtained for structure  $M_1$  determined by benzene only, compare Table 2. Deviations from the best fitting Poisson/gamma random field model tend to be highest for the additive structure  $A_2$  and the one with an increased number of cases  $M_2^{\text{high}}$ . Within the whole simulation study, the MRF model never represents the best fitting model. The size of the true covariate effect cannot be reliably reproduced by the model. For most structures, it is of comparable size as the covariate effect of multiplicative Poisson/gamma random field models.

Clustered structures  $M_4$  and  $M_5$  are best modelled by a cluster model, i.e. the BDCD model in our study. Here, the DIC is 25 to 35 points higher for Poisson/gamma random field models, compare Table 2. Details for structure  $M_5$  are discussed below. For structures assuming latent risks generated by a Gaussian kernel ( $A_2$ ,  $M_2$ ,  $M_2^{\text{high}}$ ) the deviations of the DIC from the best model are rather high.

## 6.2 Results of structure $M_5$

In the following, we discuss the results of structure  $M_5$  in more detail. As we use Gaussian kernels, sharp risk changes are more challenging for our model than smoothly varying risk surfaces in contrast to the BDCD model. Furthermore, this structure can be seen as an extension of structure  $M_4$  which has only one cluster of increased risk. For other important structures, more detailed results are presented in Appendix A2 - A7. Those structures are  $A_2$ ,  $M_2^{\text{high}}$  and  $M_1$  to  $M_4$  as presented in Table 1.

To evaluate the results, we present a scatterplot matrix where we plot estimated  $\hat{\Lambda}_i E_i$  of the applied models, see Figure 4;  $E_i$  is a population adjusted for different age and sex distributions in each area calculated as in Equation 1 and  $\hat{\Lambda}_i$  is the estimated rate. Additionally we give the estimated risk surface  $\hat{\Lambda}_i$  on a log-scale for a Poisson/gamma random field model neglecting the influence of benzene (model o). The notation for Poisson/gamma random field models is completed by a number referring to the number of employed kernels, e.g.  $M_5a15$  corresponds to a data set generated according to structure  $M_5$  (assuming an increased risk in 3 clusters and a low multiplicative influence of benzene) that is modelled by a Poisson/gamma random field model assuming an additive influence of benzene (model a) including 15 kernels.

For Poisson/gamma random field models including benzene the process of inclusion of latent covariates can be stopped at 15 or less kernels since no substantial improvement can be seen for higher numbers, see Table 3. As indicated by the DIC, there is no preference for an additive or multiplicative interpretation of benzene, both having a minimal DIC of 393.1 with 15 latent covariates. Without benzene, the inclusion of seven latent covariates leads to a minimal DIC of 417.6 which is about 25 points higher. The risk surface estimated by model  $M_5o7$  (Figure 4) is nevertheless a good representation as indicated by the scatterplots comparing the Poisson/gamma random field models. Table 3 also gives the estimates of  $\theta_{\text{benz}}$  where appropriate. The value used in the underlying model was 0.9, i.e., about twice the estimated value of the corresponding multiplicative Poisson/gamma random field model. The intercept  $\theta_0$  that was set 1 in data generation is estimated to be 1.071 (standard error 0.126) in  $M_5m15$ . Estimates of the MRF model are even lower.

**Table 3** DIC values for Poisson/gamma random field models with different number of kernels applied to simulated structure M<sub>5</sub> as well as for the BDCD and MRF model, estimates for  $\theta_{\text{benz}}$  (standard error) are given in *italic*.

# latent factors	0	1	2	3	4	5	6	7
model a	1194.5 <i>3.501</i> (0.436)	913.4 <i>4.898</i> (0.289)	527.3 <i>0.591</i> (0.416)	446.0 <i>0.974</i> (0.469)	435.0 <i>1.015</i> (0.469)	434.3 <i>0.834</i> (0.468)	426.1 <i>0.792</i> (0.463)	418.6 <i>0.763</i> (0.456)
	8	9	10	11	12	13	14	15
	413.2 <i>0.743</i> (0.456)	404.9 <i>0.738</i> (0.456)	403.8 <i>0.752</i> (0.463)	399.3 <i>0.754</i> (0.465)	397.1 <i>0.764</i> (0.468)	396.4 <i>0.785</i> (0.474)	394.7 <i>0.796</i> (0.474)	393.1 <i>0.808</i> (0.476)
	1186.7 <i>1.351</i> (0.155)	1071.5 <i>1.089</i> (1.073)	529.4 <i>0.0711</i> (0.0794)	445.1 <i>0.544</i> (0.247)	444.5 <i>0.490</i> (0.274)	436.6 <i>0.455</i> (0.279)	430.5 <i>0.366</i> (0.259)	421.3 <i>0.370</i> (0.263)
	8	9	10	11	12	13	14	15
	416.2 <i>0.360</i> (0.259)	408.9 <i>0.365</i> (0.256)	403.2 <i>0.375</i> (0.258)	401.2 <i>0.390</i> (0.265)	397.3 <i>0.405</i> (0.266)	395.9 <i>0.423</i> (0.272)	393.3 <i>0.441</i> (0.274)	393.1 <i>0.456</i> (0.280)
model o	—	1210.0	528.5	450.0	438.0	433.8	426.3	417.6
	8	9	10	11	12	13	14	15
	419.0	—	—	—	—	—	—	—
BDCD model	358.9							
MRF model	427.4, 0.041 (0.505)							

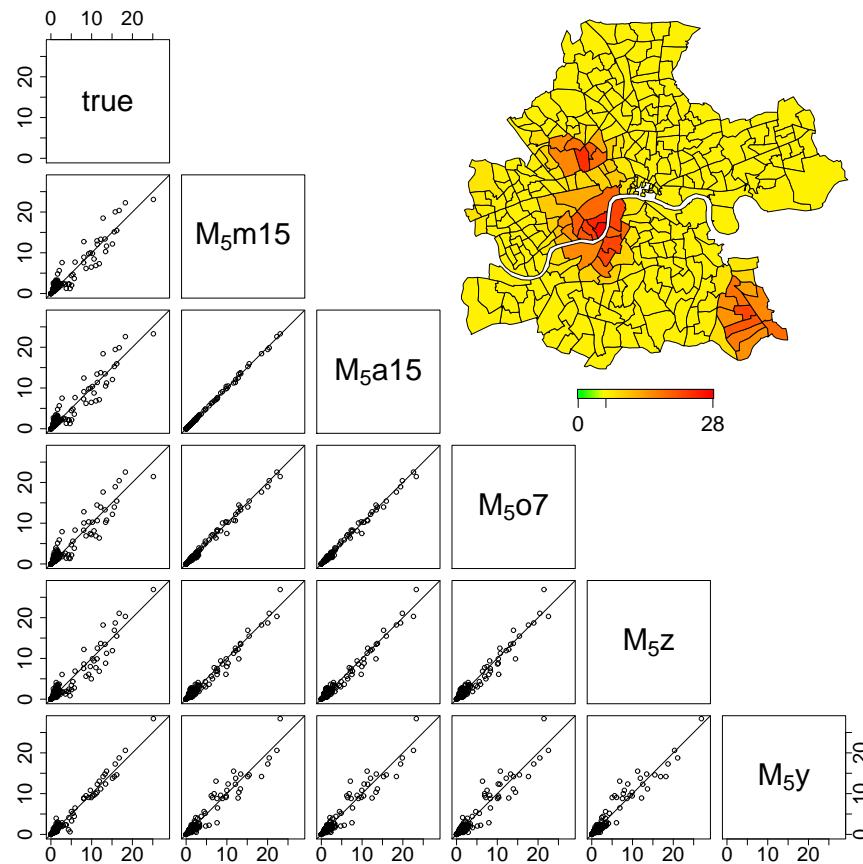
### 6.3 Identification of clusters

Poisson/gamma random field models are able to model clusters as well as sharp borders. They convince by the identification of location and size of clusters by allowing for flexible combinations of Gaussian kernels. The latent field also adopts adequately to the sharp risk decrease through estimating high precisions for Gaussian kernels.

It is possible to map the latent random field such as in Figure 5. Here, we always choose the model that allows for a multiplicative influence of benzene for a fair comparison. The plots give estimated locations of the kernels for structures where latent risk is involved in data generation. The frequency is represented by the color of each cell, cells not visited during any MCMC iterations except the burn-in period remain uncolored. On the axes we present density estimates of the location of the Gaussian kernels for longitude and latitude. For all plots the same scale is used that ranges between 0 (yellow) and 0.2 (red) corresponding to the posterior probability of being located in the corresponding cell being about 100 m<sup>2</sup> in size. These examples show different realizations of Gamma random fields and demonstrate their flexibility to identify relevant (unobserved) risk factors regardless of the underlying structure.

For the MRF-based ecologic regression model we recognize difficulties to model the sharp decrease as all wards are assumed to be neighbors leading to the highest DIC values within this structure. Concerning the DIC, best results are achieved by the BDCD model. Nevertheless, the spatial risk surface is not identified correctly as the model underestimates the risk in the very southern part substantially.

In summary, our detailed simulation study shows that Poisson/gamma random field models provide a useful and flexible tool for the identification of high risk regions of different spatial patterns. Model



**Figure 4** Scatterplot matrix of generated  $\Lambda_i E_i$  for structure M<sub>5</sub> and estimated  $\hat{\Lambda}_i E_i$  of models M<sub>5</sub>m15, M<sub>5</sub>a15, M<sub>5</sub>o7, M<sub>5</sub>Z (MRF model), and M<sub>5</sub>y (BDCD model), and plot of rates  $\hat{\Lambda}_i$  for model M<sub>5</sub>o7.

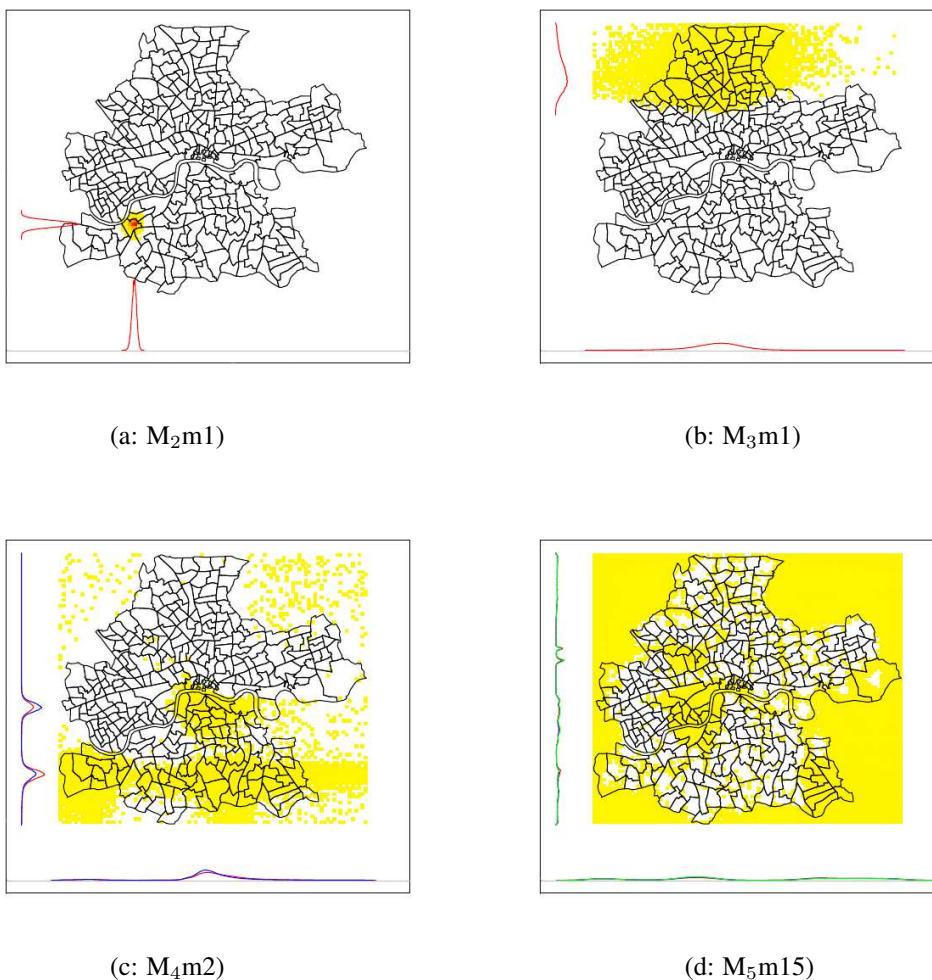
performance as well as flexibility within this model class are superior compared to those provided by MRF-based random field models. They also outperform the BCDC algorithm except for clustered structures not dominated by excess risk factors.



## 7 Applications to the real data set

We now apply Poisson/gamma random field models with different covariate interpretations as well as the MRF and BDCD model to the data set of childhood leukemia cases. For the MRF model we ignore covariate information in one scenario and use benzene as a relative risk factor in a second scenario.

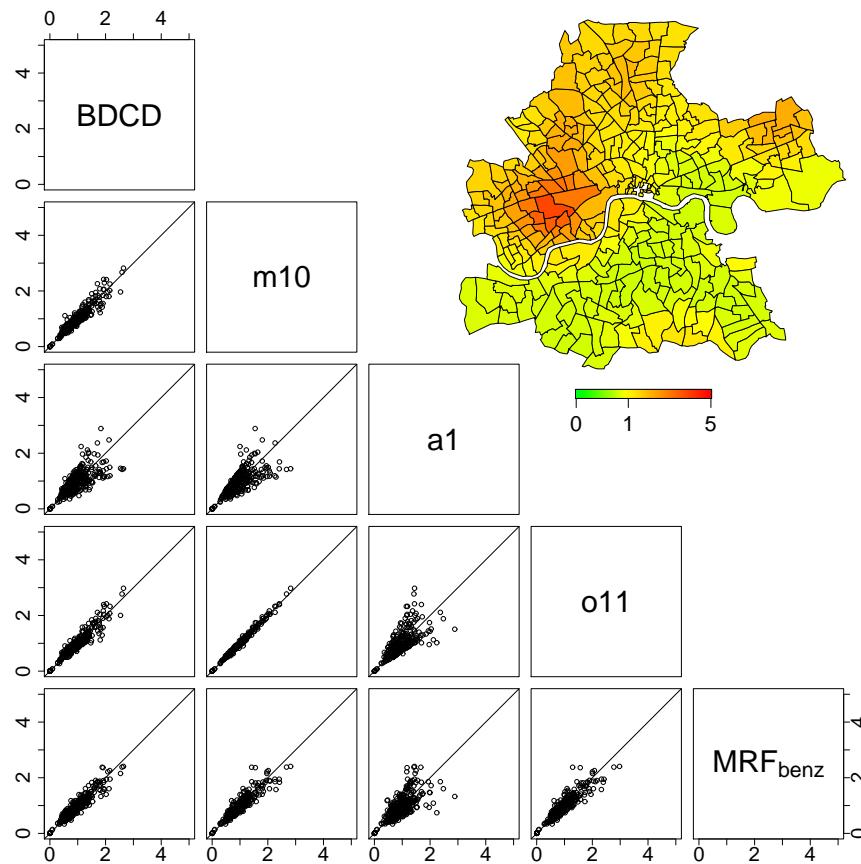
From the DIC values presented in Table 4 we can draw the following conclusions. Poisson/gamma random field models assuming benzene to be a covariate of any type require at least one latent covariate to drop the DIC substantially, compare the DIC values of models m0/a0 and m1/a1 in Table 4. When one latent risk source is involved, the DIC value is similar in the situation with (model m1/a1) or without benzene (model o1). Therefore, unobserved spatially varying risk needs to be considered in our real data example.



**Figure 5** Latent random field for Poisson/gamma random field models M<sub>2</sub>m1 (one latent Gaussian kernel in data generation), M<sub>3</sub>m1 (smoothly decreasing risk assumed), M<sub>4</sub>m2 (increased risk in southern wards assumed) and M<sub>5</sub>m15 (increased risk in three clusters assumed).

**Table 4** DIC values for Poisson/gamma random field models with different number of kernels as well as for the BDCD and MRF model applied to the leukemia data set.

# latent factors	0	1	2	3	4	5	6	7	8	9	10	11
model a	417.0	390.8	394.6	417.1	—	—	—	—	—	—	—	—
model m	415.5	391.0	397.3	391.8	388.9	386.9	385.2	384.1	383.4	382.8	382.5	382.5
model o	—	390.8	396.6	393.3	390.5	388.2	386.4	384.8	383.8	383.2	383.0	382.9
BDCD model							375.6					
MRF model with benzene								390.7				
MRF model without benzene									391.0			



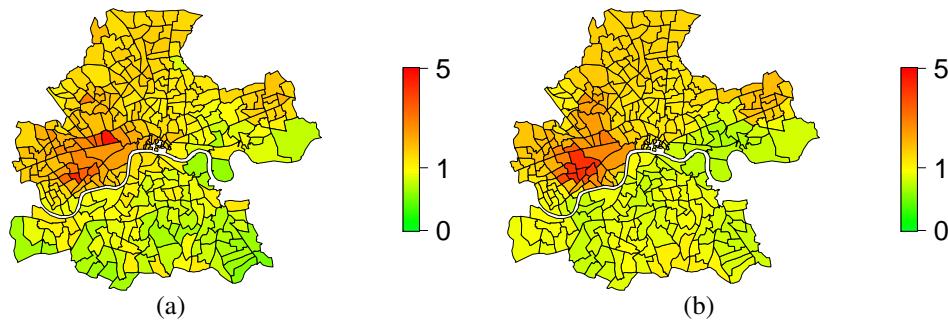
**Figure 6** Scatterplot matrix of estimated  $\hat{\Lambda}_i E_i$  of models applied to the leukemia data set: BDCD model with the lowest overall DIC, Poisson/gamma random field model a1, model m10, model o11 and MRF model including benzene, and plot of rates  $\hat{\Lambda}_i$  for model o11.

In the additive setting (model a) the model including one latent risk source (model a1) leads to the lowest DIC value of 390.8. Inclusion of further latent sources is not required. When benzene is considered to be a relative risk factor (model m), the DIC continuously decreases with inclusion of latent risk sources until the final model incorporating 10 additional Gaussian kernels (model m10, DIC=382.5). Inclusion of any further kernels does not decrease the DIC value. For Poisson/gamma random field models not including benzene (model o) we observe a similar behavior. Here, the final model requires 11 latent kernels corresponding to a DIC value of 382.9 (model o11).

For the MRF-based ecologic regression model we hardly see any benefit when benzene is included into the model. In terms of DIC, the model fit is slightly worse in general compared to Poisson/gamma random field models.

The BDCD clustering model leads to the lowest DIC value of 375.6 which reflects the fact that the covariate has hardly any influence.

Similar to the simulated scenarios we give a scatterplot matrix of the estimated values of preferred models in Figure 6. We observe a high agreement between the Poisson/gamma random field model including benzene multiplicatively and 10 latent risk sources (m10) and the one with eleven kernels only (o11) due



**Figure 7** Spatial risk surfaces of the rates  $\hat{\Lambda}_i$  estimated by the MRF model including benzene (a) and BDCD (b).

to a similar model formulation. The spatial risk surface of the model m10 is presented in Figure 8 (a), for model o11 it is given in Figure 6. The Poisson/gamma random field model assuming benzene as an additive risk factor has highest deviations to all other models as already indicated by the highest DIC value. We observe minor deviations for the MRF model and the BDCD model from all other models but the additive Poisson/gamma random field model indicating similar model fit.

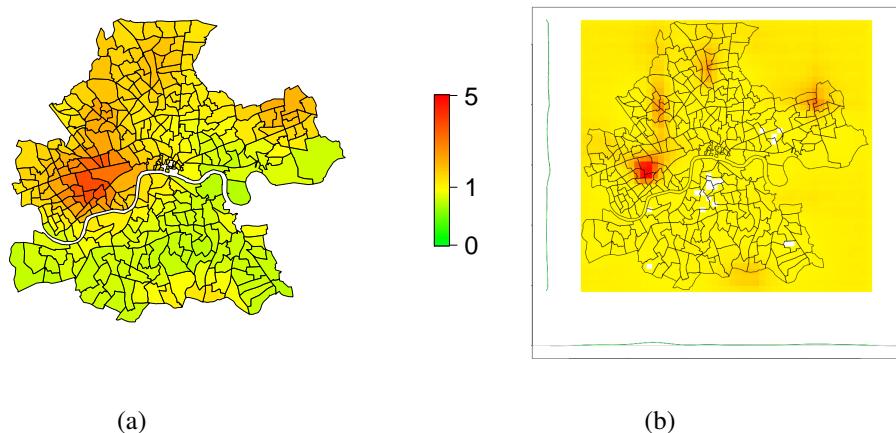
Spatial risk surfaces of the BDCD model and the MRF model including benzene are given in Figure 7. Though differences can be observed, conclusions from all three models are identical. For models that allow to incorporate covariate information we find that the covariate benzene should be included in the model, and if there is a choice — as for Poisson/gamma random field models — it is preferably modeled as a relative risk factor. Here the relative risk (RR) due to benzene is not significantly increased (RR 1.476 [95% CI 0.880; 2.071]). Using the MRF model, it is significantly higher (RR 2.478 [95% CI 1.734; 3.239]). A possible explanation is the coverage by latent risk.

The risk of leukemia is increased in the central western part of Inner London, while all wards south of the Thames are characterized by low risk. The BDCD model estimates the risk in high risk regions higher than the other two models while the MRF model has a less smoothed risk surface. The Poisson/gamma random field model without benzene tends to compromise between the two models. The estimated latent random field allows for an identification of potentially relevant covariates, compare Figure 8 (b). This figure is similar to Figure 5 but with a different color scale for the locations (red indicates a posterior density of about 0.04) for a better interpretation.

## 8 Discussion

Our paper focuses on the performance of Poisson/gamma random field models compared to MRF-based ecologic regression models and the BDCD model for disease mapping (smoothing) and cluster detection.

Both the simulation study and the real data analysis show an appropriate estimation of the risk surface by all three model classes. Nevertheless, for certain data structures certain classes are preferable. The Poisson/gamma random field model convinces by its flexibility in the estimation of random field structures. Only in situations with abruptly dropping risk (i.e. clusters) the BDCD model is preferable. As no continuous covariates can be considered by the BDCD model this holds especially for situations where covariates have hardly any influence; for clustered structures including excess risk factors the model is inferior. For MRF-based ecologic regression models inclusion of relative risk factors is possible, Poisson/gamma random field models additionally allow for excess risk factors. In contrast to the MRF model,



**Figure 8** Poisson/gamma random field model including benzene as a relative risk factor and 10 Gaussian kernels: (a) spatial risk surface of the rates  $\hat{\Lambda}_i$ , (b) latent random field.

the latent structure of the latter one is not restricted to an MRF with a prespecified neighborhood structure but can be adapted flexibly. This leads to improved results for Poisson/gamma random field models for additive and multiplicative structures.

The leukemia data set turns out to be dominated by a clustered structure where benzene has hardly any influence leading to comparable DIC values and conclusions for all model classes. Hence, the lowest DIC indicating the best fit is given by the BD<sub>CD</sub> model. The MRF-based ecologic regression model has a 15 points higher DIC, whereas the Poisson/gamma random field including benzene as a relative risk factor compromises between both models in terms of estimated risk and DIC which is increased by 7 points compared to the BD<sub>CD</sub> model.

Best et al. (2005) performed a similar simulation study using five replicates of a single data set showing different results. In contrast to our model formulation, they allow for excess risk factors only and they consider only a discrete version of Poisson/gamma models. Furthermore, a fixed number of prespecified kernel functions is assumed leading to a lack in overall fit.

To put the findings of our simulation study in context to the results of Best et al. (2005) we implemented a discrete version of Poisson/gamma random field models in WinBUGS for comparison. We used a fixed number of Gaussian kernels at evenly spaced fixed locations with a common unknown variance. The number of kernels ranged between 9 (distance 15 km) and 36 (distance 5 km). Models were employed for selected structures from the simulation study. When latent risk is involved in data generation, the discrete implementation turns out to be too restrictive. The underlying structure cannot be identified correctly, especially if the position of the kernel is not in the vicinity of the kernel used for data generation. For structure M<sub>5</sub> the DIC of the discrete version of the model is 1177.0 for the additive model and 1154.2 for the multiplicative model when 9 kernels at a fixed position are considered. Using 36 latent risk factors, it drops to 821.5 (additive model) and 820.5 (multiplicative model). These DIC values are more than twice as large as those achieved by the Poisson/gamma random field models. The problem cannot be solved by increasing the number of latent risk factors as this results in overestimating low risk regions caused by additional covariates as well as inflating computational time. Thus allowing for random locations of latent risk sources as provided by Poisson/gamma random field models leads to superior model fits.

Further improvements for Poisson/gamma random field models are possible. In our paper, kernel functions were restricted to Gaussian distributions with uncorrelated longitude and latitude. Different choices such as Gaussian distributions with a correlation structure, alternative kernel functions, or mixtures, using, e.g., uniformly distributed kernels are possible and easy to implement in the provided WinBUGS code (see Appendix A.1). Those may improve model performance concerning clustered structures.

We implemented Poisson/gamma random field models in WinBUGS using MCMC methods. As latent risk covariates have to be added successively, the process can be time-consuming. An alternative is to treat the most suitable number of latent sources as uncertain and estimate it via Reversible Jump MCMC methods (Green, 1995). Such methods are available in a beta-release from the WinBUGS development site (Lunn et al., 2008).

An alternative to the time-consuming MCMC approach is to compute the posterior distributions directly. Rue et al. (2009) developed Integrated Nested Laplace Approximation (INLA) to perform Bayesian inference for Gaussian latent fields. As Schrödle and Held (2011) show for spatial and spatio-temporal disease mapping models, INLA provides accurate results within short computational time compared to MCMC methods. However, the Poisson/gamma random field model relies on a mixture of Gamma distributions for which INLA cannot easily be applied.

In summary, Poisson/gamma random field models are a very flexible model class leading to a benefit in risk surface estimation as well as in the corresponding interpretation. They can be easily implemented in WinBUGS. The price to pay is the increased expenditure of computation time.

**Acknowledgements** The authors would like to thank Nicky Best for her collaboration on the project and helpful comments on a draft version of this paper. This work was supported by the German Academic Exchange Service [D/05/47880 to S.S.]

**Conflict of Interest** *The authors have declared no conflict of interest.*

## Appendix

### A.1. WinBUGS-code

```
# multiplicative model, random location of m latent kernels

model
{
# required constants
=====
expect <- mean(E[])

# no additive risk factors
=====

# multiplicative risk factors
=====
theta.0      ~ dgamma(a.0      , tau.0)      # prior for intercept
theta.benz   ~ dgamma(a.benz   , tau.benz)   # prior for benzene
theta.latent ~ dgamma(a.latent, tau.latent)  # prior for latent coefficient
a.0          <- 0.575                         # shape parameter for intercept
tau.0        <- a.0 * 3 * expect              # scale for intercept
a.benz       <- 0.575                         # shape for benzene
tau.benz    <- a.benz * 3 * expect             # scale for benzene
```

```

a.latent      <- 0.575                                # shape for latent coefficient
tau.latent   <- a.benz * 3 * expect                  # scale for latent coefficient

#Priors for gamma[m]'s
#=====
for (s in 1:m)
{
  gamma[s] ~ dgamma(a.gamma, tau.gamma)
}
a.gamma <- area * tau.gamma
tau.gamma <- 1/m

#Location of the latent kernels
#=====
for(i in 1:m)
{
  muX[i] ~ dunif(dist1[i], dist2[i])
  muY[i] ~ dunif(dist3[i], dist4[i])
  Sx.kernelMove[i] <- Sx.kernel[i] + muX[i]
  Sy.kernelMove[i] <- Sy.kernel[i] + muY[i]
}

#Kernel: kernel with uncertain variance
#=====
for (sx in 1:m)           # loop over latent kernels
{
  for (i in 1:I)          # loop over areas
  {
    distanceX[sx, i] <- abs(wardXcenter[i] - Sx.kernelMove[sx])
    distanceY[sx, i] <- abs(wardYcenter[i] - Sy.kernelMove[sx])
    k[sx, i]           <- exp(-(pow(distanceX[sx,i]/(2*sigmaX[sx]), 2) +
                                 pow(distanceY[sx,i]/(2*sigmaY[sx]), 2)) )
  }
  logsigmaY[sx] ~ dnorm(0, precision)
  logsigmaX[sx] ~ dnorm(0, precision)
  sigmaY[sx] <- exp(logsigmaY[sx])
  sigmaX[sx] <- exp(logsigmaX[sx])
}

#Intensities
#=====
for (i in 1:I)
{
  count[i] ~ dpois(lambda[i])
  lambda[i] <- p[i] * E[i]
  benz.term[i] <- theta.benz * benzene[i] # covariate is mean adjusted
  latent[i] <- inprod2(gamma[], k[,i])
  latent.term[i] <- theta.latent * latent[i]
  p[i] <- (theta.0 + latent.term[i]) * exp(benz.term[i])
}
}

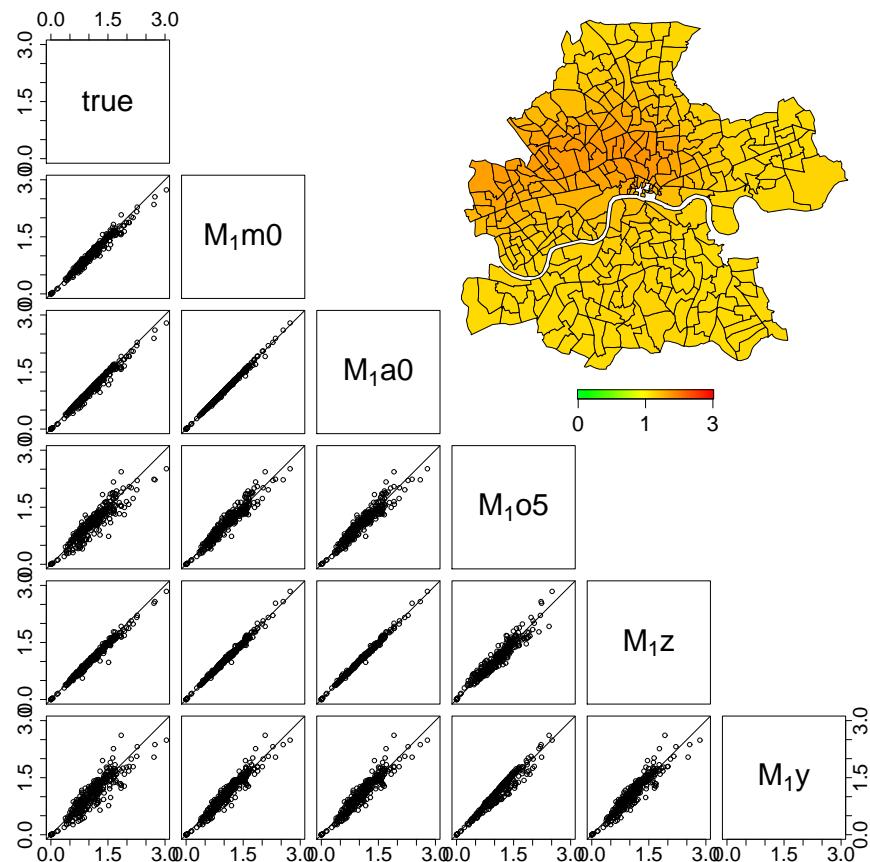
```

**Table 5** DIC values for Poisson/gamma random field models, BDCD and MRF models applied to simulated structure M<sub>1</sub>.

# latent factors	0	1	2	3	4	5	6
model a	321.9	323.6	—	—	—	—	—
model m	322.8	324.3	—	—	—	—	—
model o	—	329.7	327.8	327.2	326.0	325.9	386.4
BDCD model				328.0			
MRF model				322.4			

### A.2 Structure M<sub>1</sub>

Structure M<sub>1</sub> is characterized by a multiplicative influence of benzene accounting for about 330 expected cases, i.e.,  $\theta_{\text{benz}} = 0.9$ . When applying Poisson/gamma random field models on this structure, interpretations of the covariate lead to similar results, although  $\theta_{\text{benz}}$  is estimated to be lower, 0.929 (SE 0.373) by the additive Poisson/gamma random field model M<sub>1a0</sub> and 0.576 (SE 0.247) by the multiplicative model M<sub>1m0</sub>. Inclusion of any latent covariates is not required which corresponds to the underlying structure. Without benzene, the Poisson/gamma random field model requires five latent kernels leading to a comparable model fit. The MRF-based ecologic regression model leads to a similar risk surface with  $\theta_{\text{benz}}$  estimated as 0.706 (SE 0.255). The BDCD model produces a slightly higher DIC. Table 5 gives the calculated DIC values for all models, Figure 9 compares the estimated values in a scatterplot matrix and gives the surface estimated by model M<sub>1o5</sub>.



**Figure 9** Scatterplot matrix of generated  $\Lambda_i E_i$  for structure  $M_1$  and estimated  $\hat{\Lambda}_i E_i$  of models  $M_{1m0}$ ,  $M_{1a0}$ ,  $M_{1o5}$ ,  $M_{1z}$  (MRF model), and  $M_{1y}$  (BDCD model), and plot of rates  $\hat{\Lambda}_i$  for model  $M_{1o5}$ .

**Table 6** DIC values for Poisson/gamma random field models, BDCD and MRF models applied to simulated structure  $M_2$ .

# latent factors	0	1	2	3	4
model a	803.8	373.2	375.5	377.3	378.5
model m	821.4	373.5	373.5	377.6	373.9
model o	—	375.2	378.2	379.7	380.9
BDCD model			398.0		
MRF model			386.9		

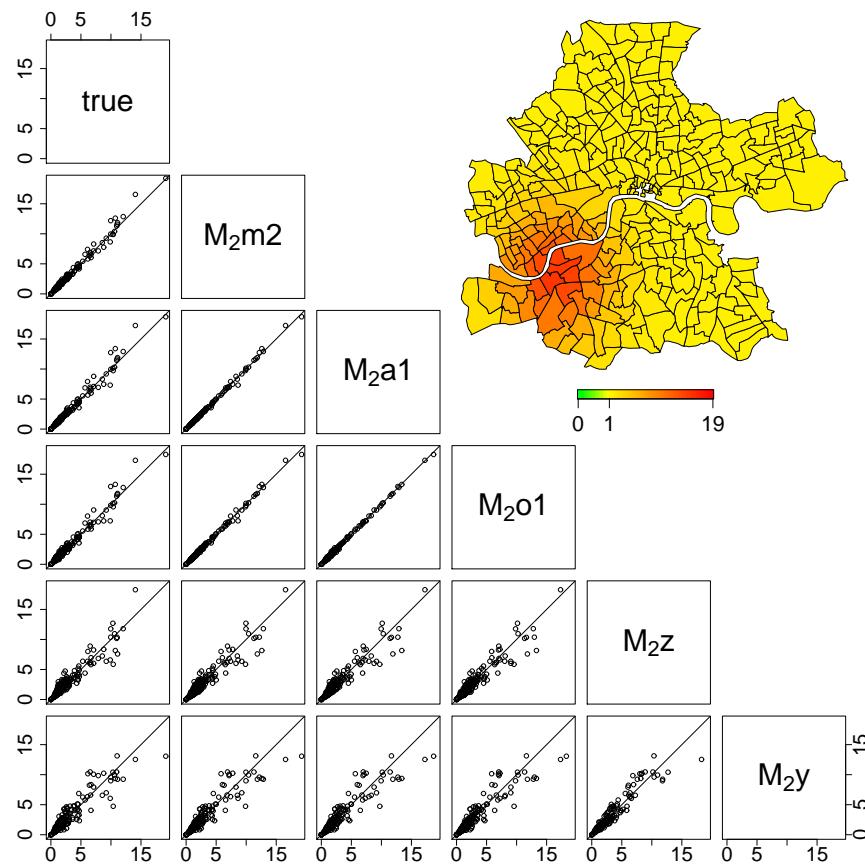
### A.3 Structure $M_2$

Structure  $M_2$  is characterized by a multiplicative influence of benzene accounting for about 330 cases, i.e.,  $\theta_{\text{benz}} = 0.9$ , combined with a latent risk covariate with an additional amount of 330 cases.

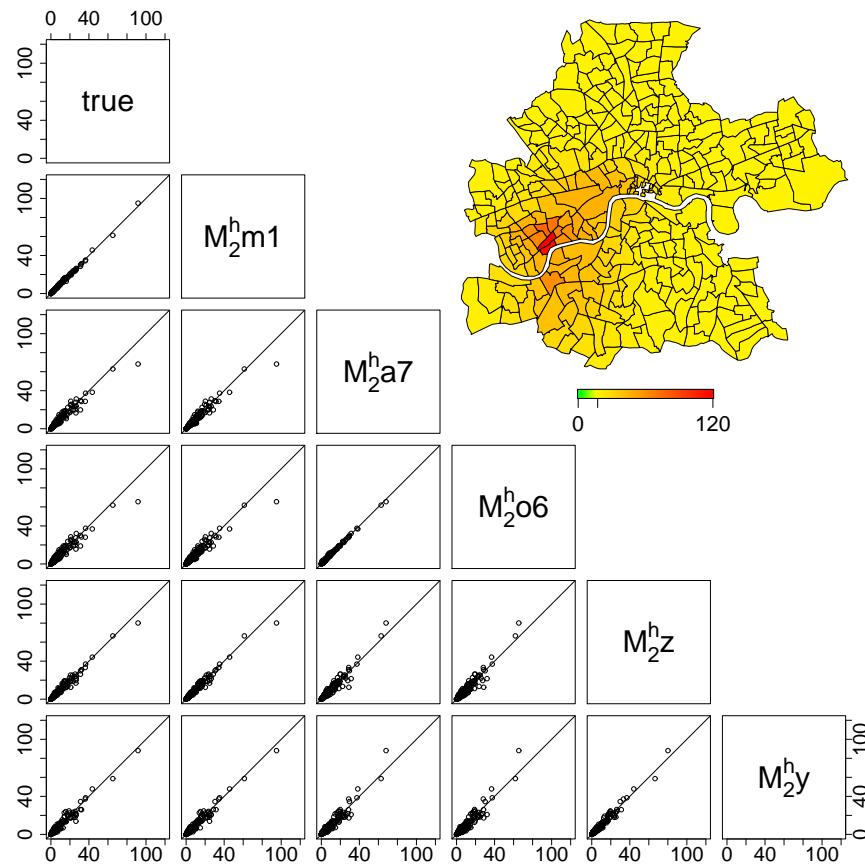
For this structure, the assumption of an either additive or multiplicative influence of benzene in Poisson/gamma random field models lead to almost identical results. The DIC drops substantially when latent covariates are considered, but as in data generation the inclusion of one Gaussian kernel is sufficient. The parameter  $\theta_{\text{benz}}$  is estimated to be lower compared to data generation, namely 0.820 (SE 0.521) by the additive Poisson/gamma random field model  $M_2a1$  and 0.375 (SE 0.230) by the multiplicative model  $M_2m2$ . Non-consideration of benzene increases the DIC slightly, here one Gaussian kernel is sufficient again for the latent field.

The MRF based ecologic regression model and the BDCD model lead to an inferior model fit. For the MRF model,  $\theta_{\text{benz}}$  is estimated as 0.868 (SE 0.359).

DIC values are summarized in Table 6, scatterplots of the results are presented in Figure 10.



**Figure 10** Scatterplot matrix of generated  $\Lambda_i E_i$  for structure  $M_2$  and estimated  $\hat{\Lambda}_i E_i$  of models  $M_{2m2}$ ,  $M_{2a1}$ ,  $M_{2o1}$ ,  $M_{2z}$  (MRF model), and  $M_{2y}$  (BDCD model), and plot of rates  $\hat{\Lambda}_i$  for model  $M_{2o2}$ .



**Figure 11** Scatterplot matrix of generated  $\Lambda_i E_i$  for structure  $M_2^{\text{high}}$  and estimated  $\hat{\Lambda}_i E_i$  of models  $M_2^{\text{high}} m1$ ,  $M_2^{\text{high}} a7$ ,  $M_2^{\text{high}} o6$ ,  $M_2^{\text{high}} z$  (MRF model), and  $M_2^{\text{high}} y$  (BDCD model), and plot of rates  $\hat{\Lambda}_i$  for model  $M_2^{\text{high}} o6$ .

#### A.4 Structure $M_2^{\text{high}}$

In contrast to structure  $M_2$ , benzene accounts for about 770 cases, i.e.,  $\theta_{\text{benz}} = 2.7$ , leading to a total of 1100 expected observations. Corresponding to the simulation scenario, the Poisson/gamma random field model including benzene as a relative risk factor leads to the best overall fit; the influence of latent kernels is similar to structure  $M_2$ , see Table 7. Benzene as an excess risk factor requires inclusion up to five to seven latent covariates without substantial decrease in DIC, which remains to be about 84 points higher compared to the lowest overall DIC. The parameter  $\theta_{\text{benz}}$  is estimated to be 2.753 (SE 0.123) by the multiplicative Poisson/gamma random field model  $M_2^{\text{high}} m1$  which is almost identical to the value used in data generation. It is estimated as 4.990 (SE 0.567) by the additive model  $M_2^{\text{high}} a7$ .

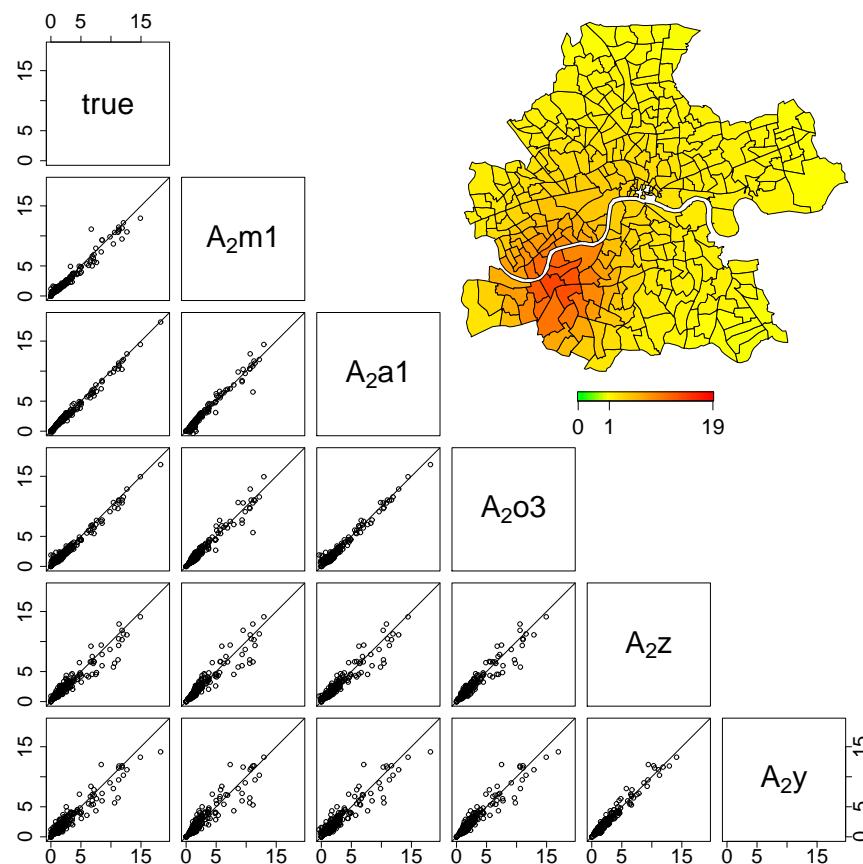
When including only latent risk sources, a similar amount of kernels is required but leads to an inferior model fit. Therefore, the inclusion of the covariate used in data generation corresponding to the interpretation of a relative risk factor is necessary. The BDCD model shows a similar performance as the Poisson/gamma random field model without any covariates. In comparison, the MRF-based ecologic regression model reduces the DIC as benzene is included, nevertheless the DIC remains to be substantially higher than for the best overall model which is the Poisson/gamma random field model considering benzene as a relative risk factor. The covariate's influence  $\theta_{\text{benz}}$  is estimated as 3.100 (SE 0.262).

**Table 7** DIC values for Poisson/gamma random field models, BDCD and MRF models applied to simulated structure  $M_2^{\text{high}}$ .

# latent factors	0	1	2	3	4	5	6	7
model a	2055.6	664.8	506.3	494.6	446.4	425.5	425.3	425.0
model m	1785.8	341.1	342.6	344.1	—	—	—	—
model o	—	768.9	627.7	524.5	551.0	479.0	478.4	478.8
BDCD model				482.2				
MRF model				400.9				

**Table 8** DIC values for Poisson/gamma random field models, BDCD and MRF models applied to simulated structure A<sub>2</sub>.

# latent factors	0	1	2	3	4
model a	732.7	346.2	346.7	352.1	353.1
model m	766.8	374.7	375.0	—	—
model o	—	410.1	408.0	384.4	385.2
BDCD model			408.0		
MRF model			387.6		



**Figure 12** Scatterplot matrix of generated  $\Lambda_i E_i$  for structure A<sub>2</sub> and estimated  $\hat{\Lambda}_i E_i$  of models A<sub>2m1</sub>, A<sub>2a1</sub>, A<sub>2o3</sub>, A<sub>2Z</sub> (MRF model), and A<sub>2y</sub> (BDCD model), and plot of rates  $\hat{\Lambda}_i$  for model A<sub>2o3</sub>.

### A.5 Structure A<sub>2</sub>

This structure is characterized by a latent risk determined by a Gaussian kernel accounting for 330 cases, i.e.,  $\theta_{\text{benz}} = 3$ . Here, the influence of benzene is assumed to be additive in data generation. We give the resulting tables for this section in Table 8, plots are presented in Figure 12.

Lowest DIC values are achieved by the additive Poisson/gamma random field model with one latent kernel corresponding to the underlying risk surface.  $\theta_{\text{benz}}$  is estimated as 3.979 (SE 0.447) in model A<sub>2a1</sub>. The multiplicative model manages to identify the number of kernels correctly, but the different covariate interpretation leads to an increase in DIC values. Here,  $\theta_{\text{benz}}$  is estimated as 1.413 (SE 0.213) in model A<sub>2m1</sub>. When benzene is not included three latent kernels are required to achieve similar results as the multiplicative model.

The MRF-based ecologic regression model also leads to a similar DIC. The covariates' influence  $\theta_{\text{benz}}$  is estimated as 1.495 (SE 0.344). In contrast, the BDCD model is not able to model the gradual descent of the data leading to an increased DIC.

**Table 9** DIC values for Poisson/gamma random field models, BDCD and MRF models applied to simulated structure  $M_3$ .

# latent factors	0	1	2	3	4
model a	517.9	339.7	340.0	343.2	390.5
model m	517.9	338.3	340.2	341.1	—
model o	—	343.5	346.0	342.7	343.4
BDCD model			345.5		
MRF model			353.0		

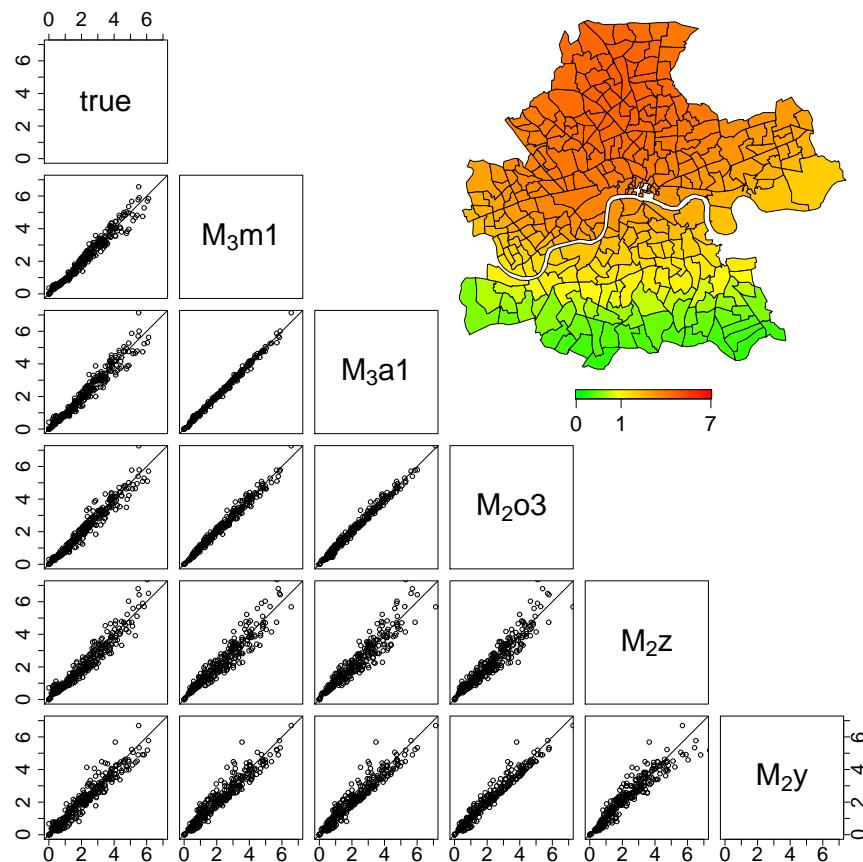
### A.6 Structure $M_3$

In contrast to structure  $M_5$  where the risk changes abruptly, it smoothly decreases for this structure. This is represented by a linear trend covariate decreasing from north to south. The influence of benzene is multiplicative with  $\theta_{\text{benz}} = 0.9$ .

Similarly for both structures, the performance of Poisson/gamma random field models including benzene as a covariate is improved when one latent covariate is considered. The multiplicative model performs best with small deviations to the additive model only.  $\theta_{\text{benz}}$  is estimated as 1.039 (SE 0.0583) in the additive Poisson/gamma random field model  $M_3a1$  and 0.466 (SE 0.240) in the multiplicative model  $M_3m1$ .

Non-consideration of benzene leads to an increase in the DIC of about five points. Here, the Poisson/gamma random field model considering three Gaussian kernels is most preferable in this group as the latent field adapts to the decreasing risk surface. The resulting risk surface is given in Figure 13.

The BCDC model shows a slightly inferior model fit letting us conclude that cluster configurations are less suitable to model a smoothly decreasing trend. The MRF based ecologic regression model leads to an increase in DIC by about 15 points compared to the best model, compare Table 9.  $\theta_{\text{benz}}$  is estimated as 1.018 (SE 0.293).



**Figure 13** Scatterplot matrix of generated  $\Lambda_i E_i$  for structure M<sub>3</sub> and estimated  $\hat{\Lambda}_i E_i$  of models M<sub>3m1</sub>, M<sub>3a1</sub>, M<sub>3o3</sub>, M<sub>2z</sub> (MRF model), and M<sub>2y</sub> (BDCD model), and plot of rates  $\hat{\Lambda}_i$  for model M<sub>3o3</sub>.

**Table 10** DIC values for Poisson/gamma random field models, BDCD and MRF models applied to simulated structure M<sub>4</sub>.

# latent factors	0	1	2	3
model a	596.3	383.6	373.0	373.1
model m	596.4	383.7	373.3	374.0
model o	—	383.5	370.9	373.4
BDCD model	335.7			
MRF model	347.8			

### A.7 Structure M<sub>4</sub>

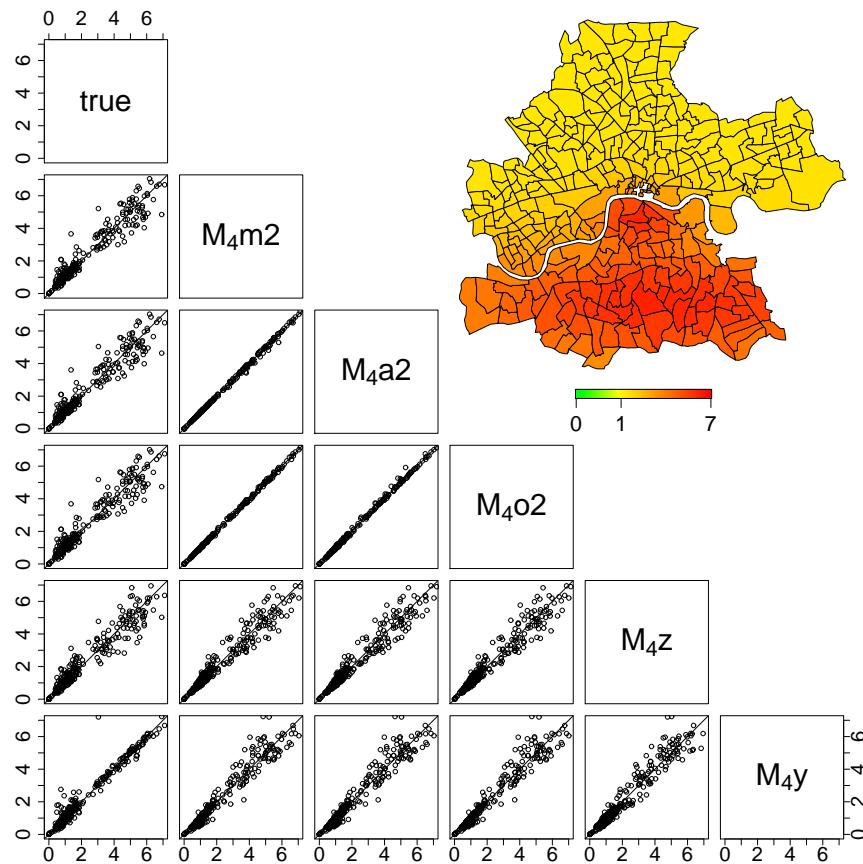
This model differs from structure M<sub>5</sub> in that one cluster of increased risk south of the Thames is generated instead of three separate clusters. We generated less additional cases per ward as the total number of cases remains constant while a higher number of wards builds the cluster. The influence of benzene is as in structure M<sub>5</sub>, namely multiplicative with  $\theta_{\text{benz}} = 0.9$ .

Results are similar as for structure M<sub>5</sub>. Different settings in Poisson/gamma random field models according to the covariate lead to similar DICs, the number of required kernels tends to be smaller as for structure M<sub>5</sub>. Inclusion of one latent kernel leads to an immense drop in DIC, a second kernel is required to achieve the minimal DIC for this model class. Due to independently estimated longitudinal and latitudinal variances of the Gaussian kernels the sharp risk decrease which is mainly in east-west direction is modeled adequately. Nevertheless, improvement is possible by alternative kernels. This is also indicated by the alternative models. While the MRF based ecologic regression model with a comparable assumption to Gaussian kernels leads to similar DICs, we notice a immense decrease for the BDCD model. Even if the sharp decrease in risk is not reproduced exactly, the model formulation allows for an abrupt change in neighboring regions leading to the lowest DIC.

Concerning the modeled influence of benzene, we estimate  $\theta_{\text{benz}}$  as 0.237 (SE 0.318) in the MRF model, as 0.346 (SE 0.320) by the additive Poisson/gamma random field model M<sub>4a2</sub>, and as 0.160 (SE 0.159) by the multiplicative Poisson/gamma random field model M<sub>4m2</sub>.

## References

- Besag, J., York, J., and Mollié, A. (1991). Bayesian image restoration, with two applications in spatial statistics, (with discussion). *Annals of the Institute of Statistical Mathematics* **43**, 1–59.
- Best, N. G., Cockings, S., Bennett, J., Wakefield, J., and Elliott, P. (2001). Ecological regression analysis of environmental benzene exposure and childhood leukaemia: Sensitivity to data inaccuracies, geographical scale and ecological bias. *Journal of the Royal Statistical Society/A* **164**, 155–174.
- Best, N. G., Ickstadt, K., and Wolpert, R. L. (2000). Spatial Poisson regression for health and exposure data measured at disparate resolutions. *Journal of the American Statistical Association* **95**, 1076–1088.
- Best, N. G., Richardson, S., and Thomson, A. (2005). A comparison of Bayesian spatial models for disease mapping. *Statistical Methods in Medical Research* **14**, 35–99.
- Breslow, N. and Day, N. (1980). *Statistical Methods in Cancer Research, Vol. 1: The Analysis of Case-Control Studies*. International Agency for Research on Cancer, Lyon.
- Buckingham, C., Clewley, L., Hutchinson, D., Sadler, L., and Shah, S. (1997). London atmospheric emissions inventory. *Technical report*, London Research Centre, London.
- Clayton, D. and Bernardinelli, L. (1992). Bayesian methods for mapping disease risk. In J. Cuzick et al. (eds), *Geographical and Environmental Epidemiology. Methods for Small Area Studies*, Oxford University Press, 205–220.
- Clayton, D. and Kaldor, J. (1987). Empirical Bayes estimates of age-standardized relative risks for use in disease mapping. *Biometrics* **43**, 671–681.
- Giudici, P., Knorr-Held, L., and Raßer, G. (2000). Modelling categorical covariates in Bayesian disease mapping by partition models. *Statistics in Medicine* **19**, 2579–2593.



**Figure 14** Scatterplot matrix of generated  $\Lambda_i E_i$  for structure  $M_4$  and estimated  $\hat{\Lambda}_i E_i$  of models  $M_{4m2}$ ,  $M_{4a2}$ ,  $M_{4o2}$ ,  $M_{4z}$  (MRF model), and  $M_{4y}$  (BDCD model), and plot of rates  $\hat{\Lambda}_i$  for model  $M_{4o2}$ .

- Green, P. J. (1995). Reversible jump Markov chain Monte Carlo computation and Bayesian model determination. *Biometrika* **82**, 711–732.
- Haining, R., Li, G. Q., Maheswaran, R., Blangiardo, M., Law, J., Best, N., and Richardson, S. (2010). Inference from ecological models: estimating the relative risk of stroke from air pollution exposure using small area data. *Spatial and Spatio-temporal Epidemiology* **1**, 123–131.
- Ickstadt, K. and Wolpert, R. L. (1999). Spatial regression for marked point processes. *Bayesian Statistics* **6**, 323–341.
- Kelsall, J. and Wakefield, J. (2002). Modeling spatial variation in disease risk: A geostatistical approach. *Journal of the American Statistical Association* **97**, 692–701.
- Knorr-Held, L. and Raßer, G. (2000). Bayesian detection of clusters and discontinuities in disease maps. *Biometrics* **56**, 13–21.
- Lunn, D. J., Best, N., and Whittaker, J. (2008). Generic reversible jump MCMC using graphical models. *Statistics and Computing* **19**, 395–408.
- Lunn, D. J., Thomas, A., Best, N., and Spiegelhalter, D. (2000). WinBUGS – a Bayesian modelling framework: concepts, structure, and extensibility. *Statistics and Computing* **10**, 325–337.
- Natario, I. and Knorr-Held, L. (2003). Non-parametric ecological regression and spatial variation. *Biometrical Journal* **45**, 670–688.

- Rue, H., Martino, S., and Chopin, N. (2009). Approximate Bayesian inference for latent Gaussian models using integrated nested Laplace approximations (with discussion). *Journal of the Royal Statistical Society/B* **71**, 319–392.
- Schrödle, B. and Held, L. (2011): A primer on disease mapping and ecologic regression using INLA. *Computational Statistics* **26**, 241–258.
- Spiegelhalter, D., Best, N., Carlin, B., and van der Linde, A. (2002). Bayesian measures of complexity and fit. *Journal of the Royal Statistical Society/B* **64**, 583–639.
- Sturtz, S. (2007). *Comparing models for variables given on disparate spatial scales: An epidemiological example*. Ph.D. thesis, Faculty of Statistics, TU Dortmund University., URL: <http://hdl.handle.net/2003/24952>
- Wakefield, J. (2007). Disease mapping and spatial regression with count data. *Biostatistics* **8**, 158–183.
- Wolpert, R. and Ickstadt, K. (1998). Poisson/gamma random field models for spatial statistics. *Biometrika* **85**, 251–267.

