# Handbook of Spatial Epidemiology

Andrew Lawson, Sudipto Banerjee, Robert Haining, Lola Ugarte

# Chapter 1

# Estimating the health impact of environmental pollution fields

*Duncan Lee and Sujit K Sahu*

## 1.1   Introduction

The health impact of many facets of the natural and built environment have been well studied in recent years, including air pollution ([24]), green space ([37]) and water quality ([50]). This chapter focuses on quantifying the health impact of air pollution, although the environmental, epidemiological and statistical challenges discussed are applicable in the wider environmental context. Quantifying the impact of air pollution is an inherently spatial as well as a temporal problem, because air pollution concentrations vary at fine spatio-temporal scales. Furthermore, individuals move through this spatio-temporal pollution field, which makes quantifying both their exposure to air pollution and its resulting health impact a difficult modelling challenge. Nevertheless, this has been an active research topic since the 1990s, with one of the first studies quantifying the effect of short-term increases in exposure in London ([45]). Since then a truly voluminous literature has developed, which has collectively quantified the health effects resulting from exposure to air pollution in both the short and the long term. This literature has included both single site studies and large multi-city studies, the latter being advantageous because of the comparability of the results across multiple locations due to unified data and analysis protocols. Collectively, these studies have helped to drive and shape legislation limiting pollution concentrations around the world, with examples being the 1990 Clean Air Act in the USA, the 2007 Air Quality Strategy for England, Scotland, Wales and Northern Ireland, and the 2008 European Parliament directive on ambient air quality and cleaner air for Europe.

Three main study designs have been used to estimate the health impact of air pollution, namely time series studies, cohort studies and areal unit studies. Time series studies are used to estimate the health impact of short-term exposure to pollution, that is a few days of elevated concentrations, often termed an air pollution episode. The disease data used in such studies are population level summaries rather than individual disease cases, meaning that

this is an ecological association study and cannot be used to determine individual level cause and effect. However, due to the routine availability of population level disease summaries, time series studies are inexpensive, quick to implement and are the most common study design. Prominent examples include the large multi-city studies entitled Air Pollution and Health: A European Approach (APHEA-2, [43]) and the National Morbidity, Mortality and Air Pollution Study (NMMAPS, [10]). In contrast, cohort studies quantify the health effects resulting from long-term exposure to pollution, that is prolonged exposure over months or years. They utilise individual-level data, and as a result individual level cause and effect can be established. However they are costly to implement, due to the large amount of data collection required and the length of time required to conduct the study due to the need for a follow up period. Examples of cohort studies include the Six Cities Study ([9]), the Multi-Ethnic Study of Atherosclerosis (MESA, [22]) and the European Study of Cohorts for Air Pollution Effects (ESCAPE, [5]).

As a result of the high cost of cohort studies, areal unit study designs have also been used to quantify the long-term health impact of air pollution. These studies are the spatial analogue of time series studies, and estimate the effects of air pollution based on spatial contrasts in disease risk and pollution concentrations across a set of contiguous areal units. Like time series studies they utilise population-level rather than individual-level disease data, and cannot be used to quantify individual level cause and effect. However, the areal unit data required to implement such studies has become widely available in recent times, with examples being the Health and Social Care Information Centre (*https://indicators.ic.nhs.uk/*) in the UK and the Surveillance, Epidemiology and End Results Program (*http://seer.cancer.gov/*) in the USA. Therefore, areal unit studies are quick and inexpensive to implement, which means that they can contribute to and independently corroborate the evidence from cohort studies. These studies have been implemented using both spatial ([21] and [27]) and spatio-temporal ([16] and [25]) designs, and the latter has also been used to estimate the short-term impact of pollution (see e.g., [53], [13] and [7]).

This chapter provides a critique of the statistical and epidemiological challenges faced by researchers conducting areal unit studies, reviews the literature in this area to date, and provides a fully worked example to illustrate how to conduct such a study. We focus on the spatial modelling challenges that arise when conducting an areal unit study, although there are similar challenges to be encountered when conducting time series or cohort studies. For simplicity, we discuss these challenges in the context of a spatial rather than a spatio-temporal study, but we note that similar challenges exist in the latter design. The layout of the remainder of the chapter is as follows. The next two sections describe the study design and data used in areal unit studies, as well as the statistical models commonly used in the literature to analyse these data. This review is followed by three sections highlighting the main statistical challenges facing researchers in this area, focusing on modelling spatial autocorrelation, spatial misalignment of the data, and allowing for within area variability in pollution concentrations. These discussions are followed by an example that illustrates the issues discussed so far, and then the chapter ends with a section providing the main conclusions and a discussion of future work needed in this area.

## 1.2 Areal unit studies

The study region $\mathcal{A}$ is a large geographical region such as a city, state or country, and is partitioned into $n$ areal units $\mathcal{A} = \{\mathcal{A}_1, \ldots, \mathcal{A}_n\}$ such as local authorities or census tracts. The areal units are typically defined by administrative boundaries, and the populations living in each one will be of different sizes and demographic structures. The disease data are denoted by $\mathbf{y} = (y_1, \ldots, y_n)$, and are counts of the total numbers of disease cases observed for each areal unit during an extended period of time such as a year. Hospitalisations and mortalities due to numerous causes have been associated with air pollution by existing studies, including cardiovascular disease ([51]), cerebrovascular disease ([35]) and respiratory disease ([27]). The differences in the population sizes and demographics between areal units are accounted for by computing the expected numbers of disease cases based on national disease rates, which are denoted here by $\mathbf{e} = (e_1, \ldots, e_n)$. For this calculation the population in each areal unit are split into a total of $R$ strata based on their age, sex and possibly ethnicity, so let $N_{ik}$ denote the number of people from areal unit $i$ in strata $k$. Letting $r_k$ denote the strata specific disease risk for the entire population, then $e_i$ is computed as $e_i = \sum_{k=1}^{R} N_{ik} r_k$. Based on the pair $(\mathbf{y}, \mathbf{e})$, we define the Standardised Morbidity/Mortality Ratio (SMR) as the ratio $\hat{\theta}_i = y_i / e_i$, which is a simple estimate of disease risk in areal unit $i$. A SMR value of one represents an average risk, while a SMR value of 1.2 means an area has a 20% increased risk of disease.

A vector of representative pollution concentrations for the $n$ areal units is denoted by $\mathbf{x} = (x_1, \ldots, x_n)$, and is typically measured in micrograms per cubic metre ($\mu g m^{-3}$). For simplicity of exposition, we work with one particular pollutant, which can also be taken as a continuous index of air quality, but the methodology can be easily extended to include multiple pollutants in which case each $x_i$ will be a vector of pollution concentrations. The health impact of numerous different pollutants have been investigated in areal unit studies, including carbon monoxide ([35]), nitrogen dioxide ([18]), ozone ([51]), particulate matter ([21]) and sulphur dioxide ([12]). The most common of these associations is with airborne particulate matter, which are small solid and/or liquid particles in the air. Particulate matter is classified by the maximum size of these particles, and particles having diameters less than 10 microns ($PM_{10}$, see [26]) and 2.5 microns ($PM_{2.5}$, see [20]) have been associated with ill health. However, estimating $\mathbf{x}$ is a challenging task, and two different data types can be used. The first of these are data from a pollution monitoring network, with examples being the State and Local Air Monitoring Stations (SLAMS) network run by the United States Environmental Protection Agency (USEPA), and the Automatic Urban and Rural Network (AURN) maintained by the Department for Environment, Food and Rural Affairs (DEFRA) of the UK government (http://www.gov.uk/defra). The second type of data are modelled concentrations from an air pollution dispersion model, with examples being the Community Multi-scale Air Quality Model (CMAQ), and the Air Quality in the Unified Model (AQUM, see [44]) developed by the UK Met Office. Monitoring networks measure air pollution concentrations with little error, but they do not have good spatial coverage and in particular some the $n$ areal units may not have any air pollution monitoring site at all. In contrast, computer dispersion models estimate pollution concentrations on a regular grid, and give complete spatial coverage of the study region without any missing observation. However, modelled concentrations are known to contain errors and biases, and are less accurate than the monitored values.

There are a number of confounding factors that must be adjusted for when estimating the health impact of air pollution using a population level spatial design, the most prominent of which is the differential rates of smoking across the $n$ areal units. However, reliable smoking data can be hard to obtain, so many studies have used measures of socio-economic deprivation as a proxy for smoking, due to the likely high correlation between smoking and deprivation ([23]). Socio-economic deprivation is multi-dimensional, and variables that have been used to account for it include individual measures of income ([25]), unemployment ([28]) and house price ([21]), as well as the Carstairs ([12]) and Townsend ([18]) deprivation indices. Let $\mathbf{U} = (\boldsymbol{u}_1^{\mathrm{T}}, \ldots, \boldsymbol{u}_n^{\mathrm{T}})$ denote the matrix of $p$ confounders, where the values relating to areal unit $\mathcal{A}_i$ are denoted by $\boldsymbol{u}_i^{\mathrm{T}} = (u_{i1}, \ldots, u_{ip})$.

## 1.3   Modelling

Poisson log-linear models are typically used to estimate the health impact of air pollution, and both classical (e.g., [21] and [16]) and Bayesian approaches have been used for inference (e.g., [13] and [27]). For the latter, Markov Chain Monte Carlo (MCMC, for details see [39]) simulation and Integrated Nested Laplace Approximations (INLA, for details see [40]) have both been used, although the latter is rare with one of the few example studies being [28]. A Bayesian approach is the most popular inferential framework in these studies, because the models used are typically hierarchical in nature and include spatial autocorrelation and different levels of variation. The first stage of a general Bayesian hierarchical model for these data is given by

$$
\begin{aligned}
y_i &\sim \text{Poisson}(e_i\theta_i) \quad \text{for} \ \ i = 1, \ldots, n, \qquad (1.1) \\
\ln(\theta_i) &= \beta_0 + x_i\beta_x + \boldsymbol{u}_i^{\mathrm{T}}\boldsymbol{\beta}_u + \phi_i, \\
\boldsymbol{\beta} = (\beta_0, \beta_x, \boldsymbol{\beta}_u) &\sim \text{N}(\boldsymbol{\mu}_\beta, \Sigma_\beta).
\end{aligned}
$$

Here the expected value of the disease count $y_i$ is the product $e_i\theta_i$, where $\theta_i$ is the risk of disease in areal unit $i$. Here a value of $\theta_i$ greater (less) than one indicates that areal unit $\mathcal{A}_i$ has a higher (lower) than average disease risk, and $\theta_i = 1.15$ corresponds to a 15% increased risk of disease. The log risk is modelled as a linear combination of an overall intercept term $\beta_0$, air pollution concentrations $(x_i\beta_x)$, confounding factors $(\boldsymbol{u}_i^{\mathrm{T}}\boldsymbol{\beta}_u)$ and a vector of random effects $\boldsymbol{\phi} = (\phi_1, \ldots, \phi_n)$. The latter accounts for the spatial autocorrelation remaining in the data after the covariate effects have been removed, as well as any overdispersion resulting from the restrictive Poisson assumption that $\text{Var}[y_i] = \mathbb{E}[y_i]$. One possible cause of this is unmeasured confounding, which occurs when an important spatially correlated covariate is either unmeasured or unknown. The spatial structure in this covariate induces spatial autocorrelation into the response, which cannot be accounted for in a regression model. Other possible causes of residual spatial autocorrelation are neighbourhood effects, where, in general, a subject's behaviour is influenced by that of neighbouring subjects, and grouping effects, where subjects choose to be close to similar subjects.

The regression parameter $\beta_x$ quantifies the relationship between air pollution and disease risk on the log scale, and is transformed to a relative risk for the purposes of interpreta-

tion. The relative risk for a $\nu$ (say) unit increase in pollution concentrations measures the proportional increase in risk from increasing pollution by $\nu$, and is calculated as

$$\mathrm{RR}(\beta_x, \nu) \;=\; \frac{e_i \exp(\beta_0 + (x_i + \nu)\beta_x + \boldsymbol{u}_i^{\mathrm{T}}\boldsymbol{\beta}_u + \phi_i)}{e_i \exp(\beta_0 + x_i\beta_x + \boldsymbol{u}_i^{\mathrm{T}}\boldsymbol{\beta}_u + \phi_i)} = \exp(\nu\beta_x). \tag{1.2}$$

Hence a relative risk of 1.05 means a 5% increase in disease risk when the pollution level increases by $\nu$ $\mu gm^{-3}$. The posterior distribution and hence 95% credible intervals for the relative risk can be computed by applying the transformation given by (1.2) to the posterior distribution for $\beta_x$, and there is substantial evidence of a relationship if the 95% credible interval does not include 1.

A number of approaches have been proposed for modelling the spatial autocorrelation and overdispersion in the data unaccounted for by the covariates, including geographically weighted regression ([51]), geostatistical models ([12]) and simultaneous autoregressive models ([21]). However, the most common approach is to model the vector of random effects $\boldsymbol{\phi}$ with a conditional autoregressive (CAR) prior (see [35], [13], [27], [25] and [29]), which is a special case of a Gaussian Markov Random Field (GMRF). This prior can be written as $\boldsymbol{\phi} \sim \mathrm{N}(\boldsymbol{0}, \tau^2 \mathbf{Q}(\mathbf{W}))$, where $\mathbf{Q}(\mathbf{W})_{n \times n}$ is a, potentially singular, precision matrix, $\tau^2$ is a variance parameter and $\boldsymbol{0}$ is an $n \times 1$ mean vector of zeros. Spatial autocorrelation is induced into this joint distribution via a binary $n \times n$ neighbourhood matrix $\mathbf{W}$, which determines the spatial adjacency structure of the $n$ areal units. If element $w_{ij} = 1$ then $(\mathcal{A}_i, \mathcal{A}_j)$ are spatial neighbours and share a common border (denoted $i \sim j$), while if $w_{ij} = 0$ (denoted $i \nsim j$) they do not.

The intrinsic model (ICAR, [3]) is the simplest CAR prior, and has a singular precision matrix (its row sums equal zero) given by $\mathbf{Q}(\mathbf{W}) = \mathrm{diag}(\mathbf{W1}) - \mathbf{W}$, where $\mathrm{diag}(\mathbf{W1})$ is an $n \times n$ diagonal matrix containing the row sums of $\mathbf{W}$. The spatial correlation structure implied by this prior is more easily observed from its full conditional form, that is as $f(\phi_i|\boldsymbol{\phi}_{-i})$ for $i = 1, \ldots, n$, where $\boldsymbol{\phi}_{-i} = (\phi_1, \ldots, \phi_{i-1}, \phi_{i+1}, \ldots, \phi_n)$. The Markov nature of this model means that the conditioning is in fact only on the random effects in geographically adjacent areal units, which induces spatial autocorrelation into $\boldsymbol{\phi}$. Using standard multivariate Gaussian theory the full conditional distribution $f(\phi_i|\boldsymbol{\phi}_{-i})$ is given by

$$\phi_i|\boldsymbol{\phi}_{-i}, \tau^2, \mathbf{W} \;\sim\; \mathrm{N}\left( \frac{\sum_{j=1}^{n} w_{ij}\phi_j}{\sum_{j=1}^{n} w_{ij}}, \; \frac{\tau^2}{\sum_{j=1}^{n} w_{ij}} \right). \tag{1.3}$$

The Bayesian model specification in (1.3) is completed by assuming $\tau^2 \sim$ inverse gamma$(a, b)$, where the hyperparameters $(a, b)$ are typically chosen to make the prior proper but weakly informative such as $(a = 2, b = 1)$ ([15]), to avoid controversy regarding the use of non-informative priors which may lead to improper posterior distributions. The ICAR prior is a natural model for strong spatial autocorrelation, because the conditional expectation is the mean of the random effects in neighbouring areas, while the conditional variance is inversely proportional to the number of neighbours. The rationale for the latter is that the more neighbours an area has the more information there is about the value of its random effect, hence its variance is smaller. However, the ICAR model can only capture strong spatial

autocorrelation, because it does not have a spatial autocorrelation parameter. Note that, if $\boldsymbol{\phi}$ is multiplied by 10 then the spatial autocorrelation structure will remain unchanged but $\tau^2$ will increase.

Therefore a number of different approaches have been proposed for allowing for varying levels of spatial autocorrelation in $\boldsymbol{\phi}$, the most popular of which is the *convolution* or *BYM* model proposed by [3]. This model augments the linear predictor in (1.1) with a second set of random effects, say $\boldsymbol{\gamma} = (\gamma_1, \ldots, \gamma_n)$, which are modelled independently as $\gamma_i \sim \mathrm{N}(0, \sigma^2)$. However, [47] discuss the identifiability problems that arise in this model due to having $n$ data points and $2n$ random effects. Therefore [31] proposed a CAR prior with a spatial autocorrelation parameter $\rho$, which has full conditional distributions given by

$$\phi_i | \boldsymbol{\phi}_{-i}, \tau^2, \rho, \mathbf{W} \;\sim\; \mathrm{N}\left( \frac{\rho \sum_{j=1}^n w_{ij} \phi_j}{\rho \sum_{j=1}^n w_{ij} + 1 - \rho}, \; \frac{\tau^2}{\rho \sum_{j=1}^n w_{ij} + 1 - \rho} \right). \tag{1.4}$$

Here $\rho = 1$ corresponds to the intrinsic CAR model for strong spatial autocorrelation, while $\rho = 0$ corresponds to independent random effects with a constant variance. A uniform prior on the unit interval is typically specified for $\rho$, and the joint distribution for $\boldsymbol{\phi}$ has a GMRF form with a precision matrix given by $\mathbf{Q}(\mathbf{W}) = \rho[\mathrm{diag}(\mathbf{W1}) - \mathbf{W}] + (1 - \rho)\mathbf{I}_n$, where $\mathbf{I}_n$ is the $n \times n$ identity matrix.

## 1.4   Controlling for unmeasured residual spatial confounding

Unmeasured spatial autocorrelation is modelled by the random effects $\boldsymbol{\phi}$, which are forced to be globally spatially smooth by CAR priors such as (1.3) or (1.4). Multivariate Gaussian theory shows that for model (1.4) the partial correlation between $(\phi_i, \phi_j)$ conditional on the remaining random effects $\boldsymbol{\phi}_{-ij}$ is given by

$$\mathrm{Corr}[\phi_i, \phi_j | \boldsymbol{\phi}_{-ij}] \;=\; \frac{\rho w_{ij}}{\sqrt{(\rho \sum_{k=1}^n w_{ik} + 1 - \rho)(\rho \sum_{l=1}^n w_{jl} + 1 - \rho)}}. \tag{1.5}$$

As $w_{ij} = 1$ for all pairs of adjacent areal units, then if $\rho$ is close to one then all pairs of adjacent random effects are spatially autocorrelated, leading to a globally smooth surface. Conversely, if $\rho$ is close to zero then all pairs of adjacent random effects are close to being conditionally independent given all other random effects, leading to no spatial smoothing anywhere in the random effects surface. In either case the random effects exhibit a single global level of spatial smoothness throughout the study region, which is likely to be inappropriate for two reasons. First, the residual spatial structure in the data after removing the covariate effects is unlikely to be globally spatially smooth, and is instead likely to exhibit localised smoothness, with strong spatial autocorrelation between some pairs of adjacent areal units whilst others exhibit abrupt step changes. The first reason for this is that the SMR is not globally smooth, as is evidenced empirically by the left panel of Figure 1.1. In that figure some pairs of adjacent local authorities exhibit similar disease risks, while between others there are large step changes. Second, the air pollution covariate is spatially smooth (see the left panel of Figure 1.2 for an example), so after removing its effect on disease risk

the residual variation is unlikely to be globally smooth.

The second reason for the inappropriateness of (1.4), existing research ([8]) has shown the potential for collinearity between these random effects and any covariate in the model that is also spatially smooth. This potential collinearity can lead to variance inflation and instability in the estimation of the air pollution effect, and the simple solution of omitting the random effects from (1.1) is not appropriate as ignoring residual spatial autocorrelation can lead to similar problems. Therefore two main approaches to solving these problems have been proposed in the literature to-date, creating random effects that are orthogonal to the covariates, such as [36] and [19], or relaxing the global smoothing restrictions of the CAR prior to allow for localised spatial smoothness such as [33] and [29].

### 1.4.1   Orthogonal smoothing

Orthogonal smoothing approaches replace $\boldsymbol{\phi}$ with a set of random effects that are orthogonal to the covariates. The first approach in this vein ([36]) does not enforce this new set of random effects to be spatially autocorrelated where as more recent work does ([19]), so we describe the latter here. Their model is based on the residual projection matrix from a normal linear model, which given the extended covariate matrix $\tilde{\mathbf{U}} = (\mathbf{x}, \mathbf{U})$ containing both pollution and the other confounding factors is given by

$$\mathbf{P} \;=\; \mathbf{I}_n - \tilde{\mathbf{U}}(\tilde{\mathbf{U}}^{\mathrm{T}}\tilde{\mathbf{U}})^{-1}\tilde{\mathbf{U}}^{\mathrm{T}}.$$

The proposed approach is based on the matrix product $\mathbf{PWP}$, where $\mathbf{W}$ is the binary neighbourhood matrix determining the spatial adjacency structure of the areal units. Thus this matrix product combines spatial information via $\mathbf{W}$ with covariate orthogonality via $\mathbf{P}$. It is shown ([19]) that the eigenvectors of $\mathbf{PWP}$ correspond to all possible mutually distinct patterns of spatial clustering orthogonal to the covariates $\tilde{\mathbf{U}}$ accounting for the spatial structure in the data via $\mathbf{W}$. Furthermore, the eigenvectors for all positive eigenvalues correspond to positive spatial correlation, while those eigenvectors relating to negative eigenvalues capture negative spatial dependence. Additionally, the magnitude of the $j$th eigenvalue $\lambda_j$ also determines the relative importance of the spatial pattern in the $j$th eigenvector, so [19] suggest choosing the first $q << n$ eigenvectors corresponding to positive and decreasing eigenvalues. Denote this $n \times q$ matrix of eigenvectors by $\mathbf{M}$, where $\mathbf{m}_i^{\mathrm{T}} = (m_{i1}, \ldots, m_{iq})$ is the $i$th row. Here $q$ is a tuning parameter in the model, with larger values leading to less dimension reduction. The model proposed by [19] is given by

$$
\begin{aligned}
y_i &\sim \text{Poisson}(e_i\theta_i) \quad \text{for} \;\; i = 1, \ldots, n, &\text{(1.6)}\\
\ln(\theta_i) &= \beta_0 + x_i\beta_x + \boldsymbol{u}_i^{\mathrm{T}}\boldsymbol{\beta}_u + \mathbf{m}_i^{\mathrm{T}}\boldsymbol{\delta},\\
\boldsymbol{\beta} = (\beta_0, \beta_x, \boldsymbol{\beta}_u) &\sim \text{N}(\boldsymbol{\mu}_\beta, \Sigma_\beta),\\
\boldsymbol{\delta} &\sim \text{N}(\mathbf{0}, \tau^2 \mathbf{Q}(\mathbf{W})_s^{-1}),
\end{aligned}
$$

where the new random effects $\boldsymbol{\delta}$ have a precision matrix given by $\mathbf{Q}(\mathbf{W})_s = \mathbf{M}^{\mathrm{T}}\mathbf{Q}(\mathbf{W})\mathbf{M}$ and $\mathbf{Q}(\mathbf{W})$ is as defined for the intrinsic CAR prior. As before, an inverse gamma prior can be specified for $\tau^2$. This model can be implemented in the *ngspatial* package for the statistical software $R$, and further details can be found in [19].

### 1.4.2   Localised smoothing

Localised smoothing approaches retain the class of CAR priors but allow for localised smoothing by modelling $\mathcal{W} = \{w_{jk} | j \sim k\}$, the elements of the neighbourhood matrix $\mathbf{W}$ corresponding to geographically adjacent areal units, as binary random quantities rather than keeping those fixed at the value 1. From (1.5) it is clear that if $w_{ij} \in \mathcal{W}$ equals one then $(\phi_i, \phi_j)$ are spatially autocorrelated and are smoothed over in the modelling process, while if $w_{ij} \in \mathcal{W}$ is estimated as zero then they are conditionally independent and no such spatial smoothing is enforced. Thus estimating the elements in $\mathcal{W}$ allows the random effects surface to exhibit localised smoothness between some pairs of adjacent random effects but not between others. A number of approaches have been proposed for estimating $\mathcal{W}$ (See [4], [33], [34], [32], [28] and [29]), and here we discuss the approach proposed by [29] because [33], [34] and [32] are set in a disease mapping rather than regression context, while [4] and [28] are in the spatio-temporal rather than spatial domain.

The approach taken by [29] is to specify a joint prior distribution for $(\tilde{\phi}, \mathcal{W})$, an extended vector of random effects and the neighbourhood adjacency elements $\mathcal{W}$. They decompose the joint prior distribution as $f(\tilde{\phi}, \mathcal{W}) = f(\tilde{\phi} | \mathcal{W}) f(\mathcal{W})$, and term their approach a Localised Conditional Autoregressive (LCAR) model. The first of these distributions is $f(\tilde{\phi} | \mathcal{W})$, a random effects model given a fixed neighbourhood structure $\mathcal{W}$. The intrinsic model (1.3) is not appropriate in this context of allowing $\mathcal{W}$ to be estimated, because $\sum_{j=1}^{n} w_{ij}$ could be estimated as zero for some areal unit $i$ leading to an infinite mean and variance. Therefore an augmented random effects vector $\tilde{\phi} = (\phi, \phi_*)$ is specified, where $\phi_*$ is a global random effect that is potentially common to all areal units and prevents the infinite mean and variance problem described above. An extended $(n+1) \times (n+1)$ neighbourhood matrix $\tilde{\mathbf{W}}$ is created for this augmented vector $\tilde{\phi}$, where there is a one-to-one relationship between a particular set of values in $\mathcal{W}$ and its matrix representation $\tilde{\mathbf{W}}$. The adjacency relation between $(\phi_i, \phi_*)$ is denoted by $w_{i*}$, and is equal to zero if all of the adjacency elements in $\mathcal{W}$ relating to $\mathcal{A}_i$ equal one. Otherwise $w_{i*} = 1$.

Based on $\tilde{\mathbf{W}}$ an intrinsic CAR prior is specified for $\tilde{\phi}$, whose precision matrix is given by $\mathbf{Q}(\tilde{\mathbf{W}}, \epsilon) = \text{diag}(\tilde{\mathbf{W}}\mathbf{1}) - \tilde{\mathbf{W}} + \epsilon\mathbf{I}$. Here $\epsilon\mathbf{I}$, with $\epsilon = 0.001$, is added to make the precision matrix diagonally dominant and hence invertible. The full conditional distribution $f(\phi_i | \tilde{\phi}_{-i})$ corresponding to this joint distribution is given by

$$\phi_i | \tilde{\phi}_{-i}, \tau^2, \tilde{\mathbf{W}} \;\sim\; \mathrm{N} \left( \frac{\sum_{j=1}^{n} w_{ij}\phi_j + w_{i*}\phi_*}{\sum_{j=1}^{n} w_{ij} + w_{i*} + \epsilon}, \; \frac{\tau^2}{\sum_{j=1}^{n} w_{ij} + w_{i*} + \epsilon} \right). \qquad (1.7)$$

This shows that if $\sum_{j=1}^{n} w_{ij} = 0$ then $w_{i*} = 1$ and the prior mean and variance simplify to $\phi_*/(1+\epsilon)$ and $\tau^2/(1+\epsilon)$, which corresponds to $\phi_i$ being independent of its neighbouring random effects. The next part of the model is the prior distribution $f(\mathcal{W})$, which is specified via a prior on its neighbourhood representation $f(\tilde{\mathbf{W}})$. A discrete uniform prior of the form

$$\tilde{\mathbf{W}} \sim \text{Discrete Uniform}(\tilde{\mathbf{W}}^{(0)}, \tilde{\mathbf{W}}^{(1)}, \dots, \tilde{\mathbf{W}}^{(N_{\mathcal{W}})}), \tag{1.8}$$

is specified. The set of $\tilde{\mathbf{W}}$ contains $N_{\mathcal{W}} = |\mathcal{W}| = \mathbf{1}^T \mathbf{W} \mathbf{1}/2$ binary elements, so its sample space has dimensionality $N_{\mathcal{W}}$ and size $2^{N_{\mathcal{W}}}$. Therefore the discrete uniform prior (1.8) is specified to vastly simplify the size and structure of the sample space for $\mathcal{W}$, which has dimension 1 and size $N_{\mathcal{W}} + 1$. This dimension reduction of the $\tilde{\mathbf{W}}$ space is undertaken because the number of elements to estimate $N_{\mathcal{W}}$ is much larger than $n$, and existing research ([32]) has shown that the elements are only weakly identifiable from the data. Element $\tilde{\mathbf{W}}^{(j)}$ has $j$ elements in $\mathcal{W}$ equal to one and $N_{\mathcal{W}} - j$ elements equal to zero, so that $\tilde{\mathbf{W}}^{(N_{\mathcal{W}})}$ is the intrinsic CAR model for strong spatial smoothing while $\tilde{\mathbf{W}}^{(0)}$ has all elements of $\mathcal{W}$ equal to zero and corresponds to independence. The set of candidate values $(\tilde{\mathbf{W}}^{(1)}, \dots, \tilde{\mathbf{W}}^{(N_{\mathcal{W}}-1)})$ are elicited from disease data prior to the study period using a Gaussian approximation, and further details are given by [29]. The full model proposed is given by

$$
\begin{aligned}
y_i | e_i, \theta_i &\sim \text{Poisson}(e_i, \theta_i) \quad \text{for } i = 1, \dots, n, \\
\log(\theta_i) &= \beta_0 + x_i \beta_x + \boldsymbol{u}_i^{\mathrm{T}} \boldsymbol{\beta}_u + \phi_i, \\
\tilde{\boldsymbol{\phi}} &\sim \text{N}(\mathbf{0}, \tau^2 \mathbf{Q}(\tilde{\mathbf{W}}, \epsilon = 0.001)^{-1}), \\
\tilde{\mathbf{W}} &\sim \text{Discrete Uniform}(\tilde{\mathbf{W}}^{(0)}, \tilde{\mathbf{W}}^{(1)}, \dots, \tilde{\mathbf{W}}^{(N_{\mathcal{W}})}), \\
\boldsymbol{\beta} = (\beta_0, \beta_x, \boldsymbol{\beta}_u) &\sim \text{N}(\boldsymbol{\mu}_\beta, \Sigma_\beta),
\end{aligned}
\tag{1.9}
$$

with an inverse gamma prior distribution specified for $\tau^2$. A software package to implement this model is provided in the supplementary material accompanying [29].

## 1.5 Estimating representative pollution concentrations

The disease and pollution data are spatially misaligned, because the geographical scales at which the data are measured are different. The disease data are available as a summary measure for each areal unit $\mathcal{A}_i$, which are typically defined by administrative boundaries and are of irregular shapes and sizes. In contrast, monitored and modelled pollution data are available at point and grid locations within the study region $\mathcal{A}$, and are typically irregularly spaced (monitored) and on a regular grid (modelled) respectively. This spatial misalignment has been termed the *change of support problem* by [14], who argue that the desired pollution concentration for areal unit $\mathcal{A}_i$ is

$$x_i = \frac{1}{|\mathcal{A}_i|} \int_{\mathbf{s} \in \mathcal{A}_i} x(\mathbf{s}) \mathrm{d}\mathbf{s}, \tag{1.10}$$

where $x_i(\mathbf{s})$ is the true unobserved concentration at location $\mathbf{s}$. Equation (1.10) is the average concentration across areal unit $\mathcal{A}_i$, and an alternative is the population weighted average pollution concentration given by

$$x_i = \int_{\mathbf{s} \in \mathcal{A}_i} p(\mathbf{s}) x(\mathbf{s}) \mathrm{d}\mathbf{s}, \tag{1.11}$$

where the population density at point $\mathbf{s}$, $p(\mathbf{s})$, is scaled so that $\int_{\mathcal{A}_i} p(\mathbf{s}) \mathrm{d}\mathbf{s} = 1$. This latter measure attempts to adjust for varying population density within an areal unit, so that the pollution measure is representative of the average concentration to which the population might be exposed. However, both (1.10) and (1.11) are unknown quantities, and their estimation can be based on either the monitored or modelled data or both. The monitored data are likely to be measured with little error, but they are spatially sparse and may not be available in many of the $n$ areal units at which the disease data are recorded. In contrast, the modelled concentrations have been predicted on a regular grid such as 1 kilometre squares, and thus have complete spatial coverage of the study region. However, they are known to contain biases and calibration problems, and are less accurate than the monitoring data.

The simplest approach to estimating (1.10) or (1.11) is to average the modelled concentrations within each areal unit, an approach adopted by [27] and [18], the latter using a population weighted average. An alternative approach has been used by [53], [13], [7] and [25], who estimate (1.10) using Monte Carlo integration. Following [25], they set up a regular grid of prediction points $\mathbf{s}_{i1}^*, \ldots, \mathbf{s}_{iN_i}^*$ within $\mathcal{A}_i$, and estimate (1.10) by

$$x_i = \frac{1}{N_i} \sum_{j=1}^{N_i} x(\mathbf{s}_{ij}^*), \tag{1.12}$$

where $x(\mathbf{s}_{ij}^*)$ is a prediction of the pollution concentration at location $\mathbf{s}_{ij}^*$ from a statistical model. They proposed a spatio-temporal model since their disease data were spatio-temporal as well. Here we adopt their approach in our spatial only case as follows. Let $z(\boldsymbol{s}_j)$ denote the observed air pollution concentration at point location $\boldsymbol{s}_j$ for $j = 1, \ldots, J$. Then we assume the following hierarchical model.

$$\begin{aligned} z(\mathbf{s}_j) &= x(\mathbf{s}_j) + \epsilon(\mathbf{s}_j), \quad \epsilon(\mathbf{s}_j) \sim \mathrm{N}(0, \sigma_\epsilon^2), \\ x(\mathbf{s}_j) &= \boldsymbol{v}(\mathbf{s}_j)^\mathrm{T} \boldsymbol{\alpha} + \eta(\mathbf{s}_j), \end{aligned} \tag{1.13}$$

for $j = 1, \ldots, J$. The true concentration $x(\mathbf{s}_j)$ is represented by a spatial process $\eta(\mathbf{s}_j)$ and covariates, $\boldsymbol{v}(\mathbf{s}_j)$, the latter of which include measures of meteorology, spatial trend terms and any other relevant information including the modelled concentrations in the grid square containing location $\mathbf{s}_j$. The vector representing the spatial process $\eta(\mathbf{s}_j)$ at all the $J$ spatial locations is denoted by $\boldsymbol{\eta}$, and is modelled by

$$\boldsymbol{\eta} \sim \mathrm{N}(\mathbf{0}, \sigma_\eta^2 S_\eta(\rho)). \tag{1.14}$$

Here $\mathbf{0}$ is a vector of zeros, and $S_\eta(\rho)$ is a spatial correlation matrix which has elements $S_\eta(\rho)_{jk} = \exp(-\rho||\mathbf{s}_j - \mathbf{s}_k||)$, where $||.||$ denotes Euclidean distance. From this model

Bayesian spatial prediction is used to estimate (1.12) for all $n$ areal units, where samples are drawn from the posterior predictive distribution $f(x(\mathbf{s}_{ij}^*)|\mathbf{z})$ for each prediction location, where $\boldsymbol{z}$ is the vector of all data points $z(\mathbf{s}_j)$. MCMC simulation is used to sample from the posterior predictive distribution, and samples can be combined over the $N_i$ prediction locations to estimate (1.12). This process can be repeated for a number of MCMC samples, giving a posterior predictive distribution for (1.12) from which a single point estimate could be computed or the entire posterior distribution could be used. The latter approach would allow for the inherent variation in the vector of areal unit specific pollution concentrations $\mathbf{x}$, and approaches to propagate this variation into the health model are discussed in the next section.

In this chapter we propose to use a recently developed downscaler model ([41], and then generalised by [2], [1], and [54]). The downscaled method is implemented by assuming the single covariate $v(\mathbf{s}_j)$ to be the modelled concentration for the grid cell $\mathcal{B}$ that contains the location $\boldsymbol{s}_j$. The modelled concentrations, we use are those from the AQUM, developed by the UK Met Office ([44]) on a 12 Kilometre square grid cell. The grid cells are denoted by the red dots in the right panel of Figure 1.1. The monitoring data we use are obtained from the publicly available data from the AURN network, for details see *http://uk-air.defra.gov.uk/*.

## 1.6 Propagating pollution uncertainty into the health model

The Poisson log-linear model (1.1) treats the vector of estimated pollution concentrations $\mathbf{x} = (x_1, \ldots, x_n)$ as known constants, so that $x_i$ is the known and constant pollution concentration for areal unit $\mathcal{A}_i$. This assumption ignores two different sources of uncertainty or variation in $\mathbf{x}$ when estimating its health effects. The first source of uncertainty is due to *measurement error*, which occurs because the true constant pollution exposure in each areal unit is unknown and its estimated value is subject to error and uncertainty. The second source of variation is that the pollution concentration is not constant across each areal unit, meaning that there is within-area variability in exposure. This within-area variability in exposure means that the population level risk model (1.1) has a different algebraic form compared to what one would obtain by aggregating an individual level risk model to the population scale. The bias in the health effect estimate resulting from this is known as *ecological bias*, and models to overcome this problem and that of measurement error are discussed below.

### 1.6.1 Measurement error models

Measurement error models are yet to be applied in air pollution and health studies with an areal unit design, but have been extensively investigated in the context of both time series (see [52] and [6]) and cohort studies (see [17] and [46]). Both *Classical* and *Berkson* measurement error models have been applied in the literature, and in some cases in combination in the same model ([52] and [46]). Measurement error models have been widely applied to account for many different types of errors, including spatial misalignment as described in the previous section ([17]), and the difference between outdoor concentrations and personal exposures ([11]). The simplest measurement error set up is that the true constant pollution

concentration in areal unit $\mathcal{A}_i$ is unknown and denoted by $x_i$, and needs to be estimated by error-prone measurements $x_{i1}, \ldots, x_{ig_i}$. Then a simple *Classical* measurement error model fitted in a Bayesian setting is given by

$$
\begin{aligned}
x_{i1}, \ldots, x_{ig_i} &\sim \mathrm{N}(x_i, \sigma^2) \quad \text{for } i = 1, \ldots g_i, \\
x_i &\sim \mathrm{N}(\mu_x, \sigma_x^2), \\
\sigma^2 &\sim \text{Inverse-Gamma}(a, b),
\end{aligned}
\tag{1.15}
$$

where weakly informative priors are typically specified for $(x_i, \sigma^2)$. This model could be added as an extra level in a hierarchical health model such as (1.1), (1.6) or (1.9). Thus $\mathbf{x} = (x_1, \ldots, x_n)$ would be treated as a set of parameters to be estimated in the model, and the uncertainty in its value would be propagated through the health model.

### 1.6.2  Ecological bias models

A causal relationship between air pollution and health can only be estimated from individual level disease data such as that modelled in a cohort study, and not from the population level disease data used in areal unit studies. Naively interpreting the population level association found in these studies in terms of individual-level cause and effect is incorrect, and is known as the *ecological fallacy*. The difference between the estimated individual and population level relationships is known as *ecological bias*, and has been the subject of extensive study ([48]). This bias occurs when there is within area variability in the pollution concentrations, because a non-linear risk model changes its form under aggregation from the individual to the population level. Consider the ideal situation of having individual level data on disease presence and pollution exposure for all individuals in the study region. Let $y_{ik}$, for $k = 1, \ldots, n_i$, be the binary observation denoting whether individual $k$ in areal unit $\mathcal{A}_i$ has the disease under study or not, and let $x_{ik}$ denote that individuals pollution exposure. Then a simple individual-level risk model is given by

$$
\begin{aligned}
y_{ik} &\sim \text{Bernoulli}(p_{ik}) \quad \text{for } k = 1, \ldots n_i, \ i = 1, \ldots, n, \\
\ln(p_{ik}) &= \beta_0 + x_{ik}\beta_I,
\end{aligned}
\tag{1.16}
$$

where a log rather than logit link is used because the likelihood of disease presence in any single individual is small. However, for areal unit studies the individual $y_{ik}$'s are unknown, and only the total number of disease cases from the population living in each areal unit $y_i = \sum_{k=1}^{n_i} y_{ik}$ are known. If there is no within area variability in exposure, that is if $x_{ik} = x_k$ for $k = 1, \ldots, n_i$, then $p_{ik} = p_k$ and there is no ecological bias as (1.16) aggregates to a binomial model, or a Poisson approximation to it as in (1.1). However, when there is within area variability in exposure then computing the expectation of the aggregated $y_i$ from the individual level model gives:

$$
\mathbb{E}[y_i] = \mathbb{E}\left[\sum_{k=1}^{n_i} y_{ik}\right] = \exp(\beta_0)\sum_{k=1}^{n_i} \exp(\beta_I x_{ik}) = \exp(\beta_0^*)\mathbb{E}[\exp(x_{ik}\beta_I)].
\tag{1.17}
$$

Thus, comparing (1.1) and (1.17) shows that ecological bias occurs because $\mathbb{E}[\exp(x_{ik}\beta_I)] \neq \exp(\mathbb{E}[x_{ik}\beta_x])$, so that in general $\beta_I \neq \beta_x$. Two main approaches have been proposed to solve this problem, the first of which uses a sample of exposures $x_{i1}, \ldots, x_{ig_i}$ for areal unit $\mathcal{A}_i$. These exposures could be $g_i$ data points in the same areal unit, or $g_i$ samples from the posterior predictive distribution of (1.10) obtained from a first stage pollution model such as (1.13). Based on such a sample of exposures [49] propose a Poisson log-linear convolution model for $y_i$, where the risk model for $\theta_i$ corresponding to (1.1) would be changed to

$$\theta_i = \exp(\beta_0^* + \boldsymbol{u}_i^{\mathrm{T}}\boldsymbol{\beta}_u + \phi_i)\sum_{k=1}^{g_i}\exp(x_{ik}\beta_I), \tag{1.18}$$

which essentially approximates $\mathbb{E}[\exp(x_{ik}\beta_I)]$ by its sample average. The second approach was first proposed by [38], who suggested representing the within area distribution of exposures by a parametric distribution. Let $X_i$ denote the random variable characterising the within-area exposure distribution for areal unit $\mathcal{A}_i$. Then the desired quantity from (1.17) is $\mathbb{E}[\exp(X_i\beta_I)]$, the moment generating function of $X_i$. Assuming that the within-area exposure distribution is $X_i \sim \mathrm{N}(\mu_i, \sigma_i^2)$ leads to the risk model

$$\theta_i = \exp(\beta_0 + \boldsymbol{u}_i^{\mathrm{T}}\boldsymbol{\beta}_u + \phi_i + \mu_i\beta_I + \sigma_i^2\beta_I^2/2), \tag{1.19}$$

where $(\mu_i, \sigma_i^2)$ can be estimated by their sample equivalents from the $g_i$ samples $x_{i1}, \ldots, x_{ig_i}$. However, if the within-area exposure distribution is skewed then a normal approximation may be inappropriate, and a log-normal distribution could be used instead. The moment generating function of a log-normal distribution does not exist, so [42] propose approximating it with a three term Taylor expansion leading to the risk model
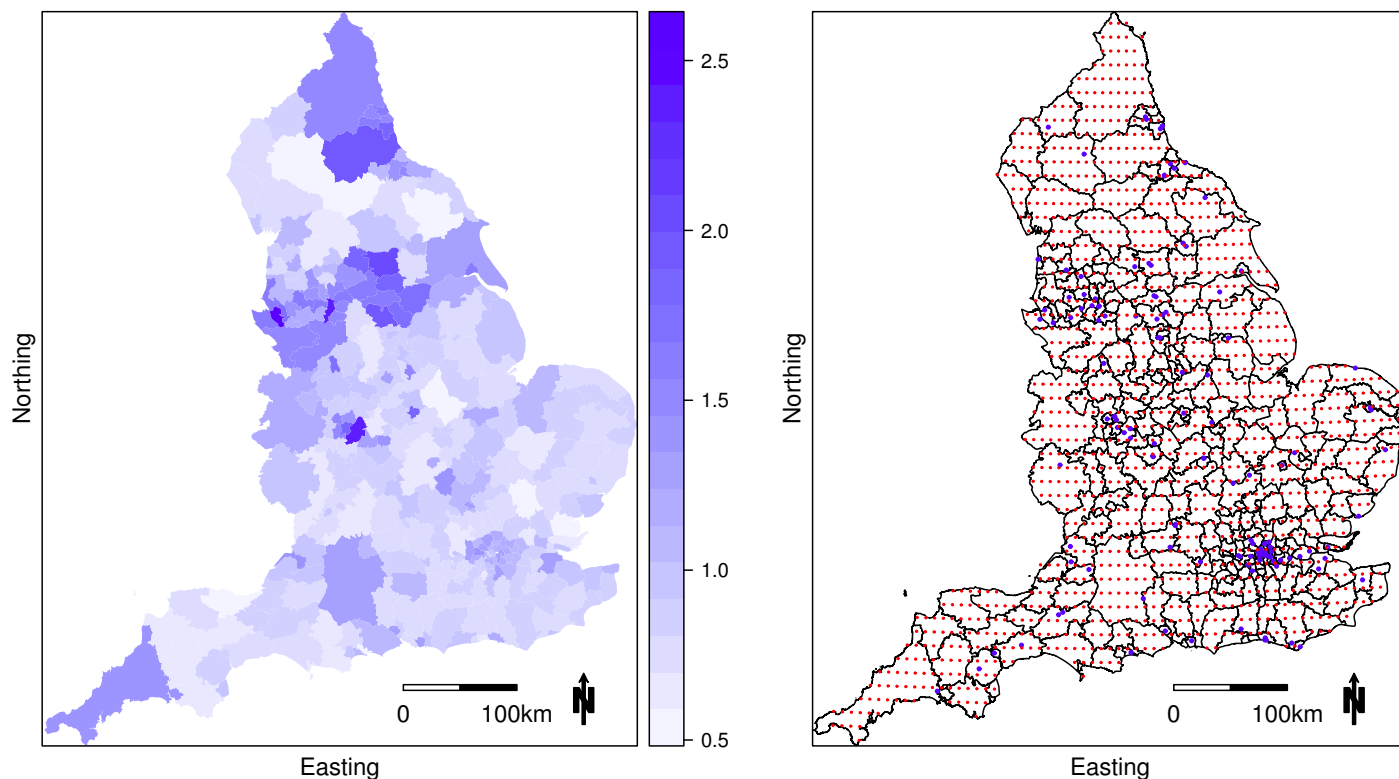
$$\theta_i = \exp(\beta_0 + \boldsymbol{u}_i^{\mathrm{T}}\boldsymbol{\beta}_u + \phi_i + \mu_i\beta_I + \sigma_i^2\beta_I^2/2 + \eta_i^3\beta_I^3/6), \tag{1.20}$$

where $\eta_i^3$ is the third central moment and is given by $\eta_i^3 = \sigma_i^2/(\mu_i[(\sigma_i^2/\mu_i^2) + 3])$.

## 1.7 Illustrative example

We illustrate the methods discussed in this chapter by presenting a new study examining the impact of long-term exposure to $PM_{2.5}$ on respiratory hospitalisation risk in England in 2010. For this study England is partitioned into $n = 323$ local and unitary authorities, which are typically either individual cities or larger rural areas. The disease data are counts of the numbers of hospital admissions due to respiratory disease in 2010, and the spatial pattern in the standardised morbidity ratio ($SMR_i = y_i/e_i$) is displayed in the left panel of Figure 1.1. The figure shows that the highest risk areas are cities in the north and central parts of England, such as Liverpool, Birmingham and Manchester. In contrast, the lowest risk areas are typically rural, and include Rutland, West Somerset and Richmondshire. The SMR map shows evidence of localised spatial smoothness, with some pairs of neighbouring areal units exhibiting similar risks while other pairs have vastly different values.

Figure 1.1: The left map displays the standardised morbidity ratio (SMR) for hospital admissions due to respiratory disease in 2010, while the right map shows the locations of the pollution monitors (blue dots) and the corners of the 12 Kilometre square grid cells (red dots).



The pollution metric considered in this study is the annual average PM$_{2.5}$ concentration levels for 2009, where the monitoring data come from the AURN network of sites while the modelled concentrations come from the AQUM model. Their locations are displayed in the right panel of Figure 1.1, where the blue dots represent the monitor locations while the red dots are the corners of the 12 kilometre square grid cells for which the modelled data are available. The figure shows that the monitors are clustered mainly in the cities, while the rural areas such as the south west of England have very poor spatial coverage. In contrast, the modelled concentrations are calculated on a regular grid of size 12 kilometres, and thus provide complete spatial coverage of the study region. Finally, the confounding effects of socio-economic deprivation on disease risk were accounted for by including the English index of multiple deprivation into the model.

The modelling of these data was undertaken in two stages. In the first stage, using MCMC we fitted (1.13) to annual PM$_{2.5}$ data from $J = 166$ monitoring sites for the year 2009. At the $t$th iteration ($t = 1, \ldots, 5000$) we then obtained $x_i^{(t)}$ using (1.12), where $x(\boldsymbol{s}_{ij}^{*(t)})$ was a draw from the predictive distribution as mentioned above. Then in stage two the health impact of PM$_{2.5}$ was estimated, using the posterior distribution of $x_i^{(t)}$. Three different health models were fitted to the data. Model A is given by (1.1) in conjunction with the CAR prior (1.4), and is routinely used in studies of this type. In this model the pollution concentrations

are assumed to be fixed, and have been estimated as the mean of the posterior predictive distribution of (1.12) based on 5000 posterior samples. Model B also assumes the pollution concentrations are fixed, but extends model A by using a more flexible spatial correlation model for the random effects, namely the LCAR model given by (1.9). Finally, Model C also uses the LCAR model (1.9), but extends model B by allowing for within-area variation in the pollution concentrations via the log-normal model given by (1.20). The first three moments of this log-normal approximation are computed based on their sample equivalents from the posterior predictive distribution of (1.12). A log-normal model was used rather than a normal one as the posterior predictive distributions for (1.12) exhibited small amounts of right skewness.
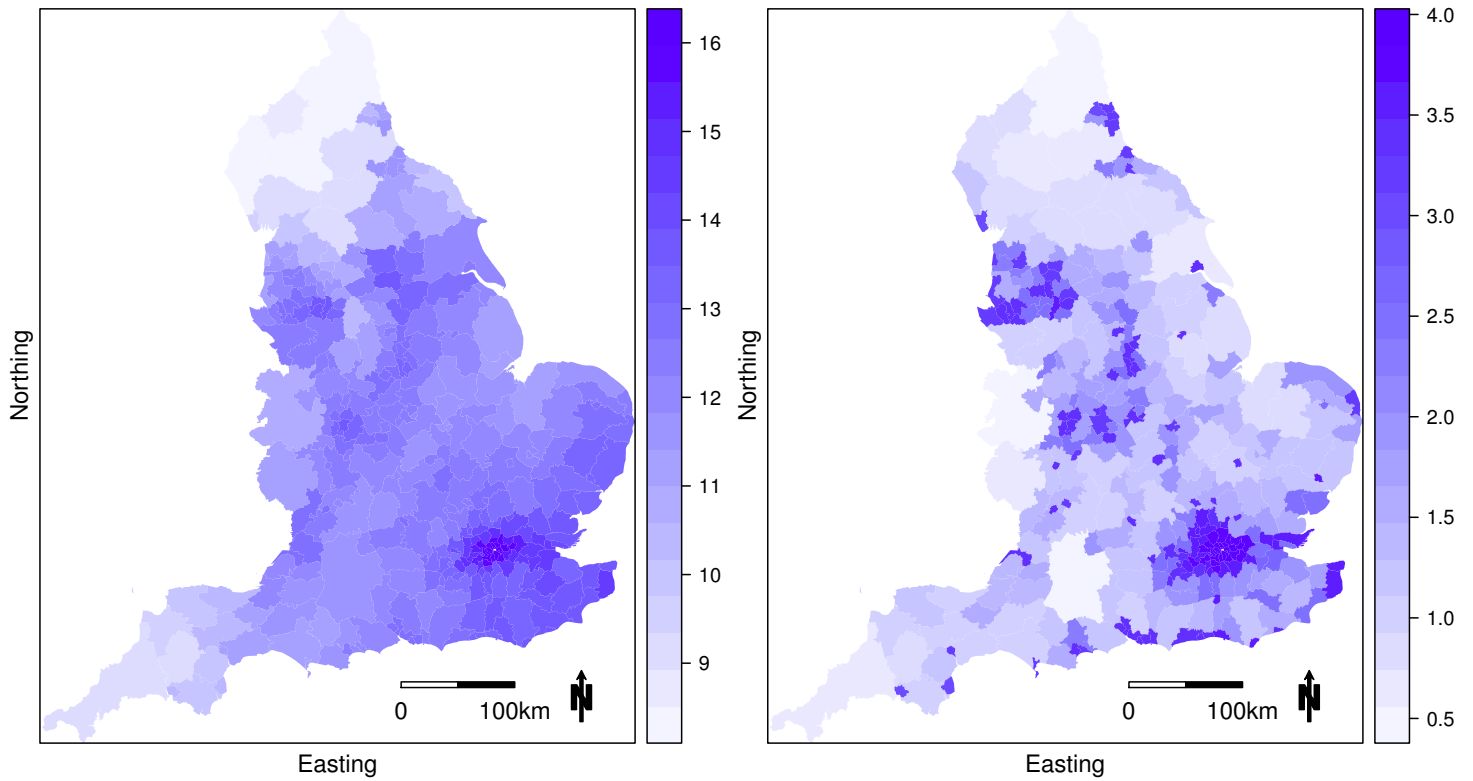
The left panel of Figure 1.2 displays the posterior predictive mean average $PM_{2.5}$ concentration for each local or unitary authority obtained from the pollution model in stage one, and shows that the highest concentrations are observed in the city of London in the south east of England. The other highly polluted areas are the large cities of Birmingham and Manchester, while the rural areas have the lowest pollution concentrations, particularly in the far south west and north of England where population density is relatively low. The estimated $PM_{2.5}$ concentrations are spatially smooth, which illustrates the potential for collinearity with spatially smooth random effects highlighted by [8]. The posterior predictive standard deviations in $PM_{2.5}$ are displayed in the right panel of Figure 1.2, and show substantial uncertainty in exposure and a clear mean-variance relationship, as the largest values coincide with higher mean concentrations.

The estimated relative risks and 95% credible intervals corresponding to a 1 $\mu gm^{-3}$ increase in $PM_{2.5}$ concentrations are: Model A - *1.032 (1.001, 1.079)*; Model B - *1.040 (1.011, 1.067)*; and Model C - *1.034 (1.011, 1.057)*. All three models suggest that areal units with higher concentrations of $PM_{2.5}$ exhibit higher risks for respiratory disease, with increases ranging between 3.2% and 4%. The three models exhibit differences in their estimated risks, as replacing a globally smooth set of random effects (Model A) with a locally smooth set (Model B) has inflated the risk. In addition, the 95% credible intervals from Model A are wider than those from Model B, and both these effects may be due to the collinearity between the globally smooth random effects in Model A and $PM_{2.5}$. Moving from Model B to Model C allows for the inherent within-area variation in the average $PM_{2.5}$ concentrations in each areal unit, which has led to an attenuation of the risk by 0.6%. Thus the estimated risks for Models A and C are similar, which may be due to the two biases described above working in opposite directions. Thus, overall we believe the most reliable estimate comes from Model C, as it can capture more flexible spatial autocorrelation structures than Model A, while correctly allowing for the variation in the pollution concentrations when estimating its health effects unlike Model B.

## 1.8 Discussion

This chapter has critiqued the statistical challenges involved in estimating the long-term health impact of air pollution using an ecological areal unit study design, and has provided a worked example to show the potential impact of an inappropriate model specification. The modelling approach taken in the latter was implemented in two stages, a first stage pollution model, whose results were then used in a second stage health model. Two-stage

Figure 1.2: The maps display the posterior mean (left panel) and standard deviation (right panel) of the annual average $PM_{2.5}$ concentrations in 2009 for each local and unitary authority in England

approaches such as this are becoming common in the general air pollution and health literature (see [6] and [25]), but as highlighted by [46] induce their own biases into the health effects. Therefore future research in this area will need to consider a single integrated model, that simultaneously estimates the spatio-temporal pattern in air pollution concentrations as well as its resulting effects on disease. An important issue in an integrated Bayesian model is that of feedback, namely should information in the disease counts be allowed to effect the estimated pollution concentrations, when it is the relationship in the opposite direction that is of primary interest. In time series studies this feedback has been prevented (see [30]), but an interesting area of work would be to examine the impact of this in an areal unit study.

The other future research direction is to extend the methodology discussed here into the spatio-temporal domain. The development of locally smooth/orthogonal random effects models and allowing for within area variation in exposure have mainly been considered in the purely spatial domain, and a number of studies are now utilising spatio-temporal data ([20], [16] and [25]). In addition to the challenges outlined here for spatial studies, spatio-temporal studies are likely to throw up a number of additional modelling challenges for which methodological development is required. The most obvious of these is developing a model for spatio-temporal autocorrelation, which has to be flexible enough to allow for non-stationarity, non-separability and allow for varying levels of smoothness in both space and time. The use of a spatio-temporal study also naturally leads to questions about lag times of air pollution effects, a subject that has only been investigated in a time series context to date.

# References

[1] V J Berrocal, A E Gelfand, and D M Holland. A bivariate space-time downcsaler under space and time misalignment. *Annals of Applied Statistics*, 4:1942–1975, 2010.

[2] V J Berrocal, A E Gelfand, and D M Holland. A spatio-temporal downcsaler for outputs from numerical models. *Journal of Agricultural, Biological and Environmental Statistics*, 15:176–197, 2010.

[3] J Besag, J York, and A Mollie. Bayesian image restoration with two applications in spatial statistics. *Ann I Stat Math*, 43:1–59, 1991.

[4] A Brezger, L Fahrmeir, and A Hennerfeind. Adapative Gaussain Markov random fields with applications in human brain mapping. *Journal of the Royal Statistical Society Series C*, 56:327–345, 2007.

[5] G Cesaroni, F Forastiere, M Stafoggia, Z Andersen, C Badaloni, R Beelen, B Caracciolo, U de Faire, R Erbel, K Eriksen, L Fratiglioni, C Galassi, R Hampel, M Heier, F Hennig, A Hilding, B Hoffmann, D Houthuijs, K Jckel, M Korek, T Lanki, K Leander, P Magnusson, E Migliore, C Ostenson, K Overvad, N Pedersen, J Pekkanen, J Penell, G Pershagen, A Pyko, O Raaschou-Nielsen, A Ranzi, F Ricceri, C Sacerdote, V Salomaa, V Swart, A Turunen, P Vineis, G Weinmayr, K Wolf, K andde Hoogh, G Hoek, B Brunekreef, and A Peters. Long term exposure to ambient air pollution and incidence of acute coronary events: prospective cohort study and meta-analysis in 11 European cohorts from the ESCAPE Project. *British Medical Journal*, 348, 2014.

[6] H Chang, R Peng, and F Dominici. Estimating the acute health effects of coarse particulate matter accounting for exposure measurement error. *Biostatistics*, 12:637–652, 2011.

[7] J Choi, M Fuentes, and B Reich. Spatial-temporal association between fine particulate

matter and daily mortality. *Computational Statistics and Data Analysis*, 53:2989–3000, 2009.

[8]  D Clayton, L Bernardinelli, and C Montomoli. Spatial Correlation in Ecological Analysis. *International Journal of Epidemiology*, 22:1193–1202, 1993.

[9]  D Dockery, C Pope, X Xu, J Spengler, J Ware, M Fay, B Ferris, and F Speizer. An Association Between Air Pollution And Mortality In Six U.S. Cities. *The New England Journal of Medicine*, 329:1753–1759, 1993.

[10]  F Dominici, M Daniels, S Zeger, and J Samet. Air Pollution and Mortality: Estimating Regional and National Dose-Response Relationships. *Journal of the American Statistical Association*, 97:100–111, 2002.

[11]  F Dominici and S Zeger. A measurement error model for time series studies of air pollution and mortality. *Biostatistics*, 1:157–175, 2000.

[12]  P Elliott, G Shaddick, J Wakefield, C Hoogh, and D Briggs. Long-term associations of outdoor air pollution with mortality in Great Britain. *Thorax*, 62:1088–1094, 2007.

[13]  M Fuentes, H Song, S Ghosh, D Holland, and J Davis. Spatial Association between Speciated Fine Particles and Mortality. *Biometrics*, 62:855–863, 2006.

[14]  A Gelfand, L Zhu, and B Carlin. On the change of support problem for spatio-temporal data. *Biostatistics*, 2:31–45, 2001.

[15]  A Gelman. Prior distributions for variance parameters in hierarchical models. *Bayesian Analysis*, 1:515–534, 2006.

[16]  S Greven, F Dominici, and S Zeger. An Approach to the Estimation of Chronic Air Pollution Effects Using Spatio-Temporal Information. *Journal of the American Statistical Association*, 106:396–406, 2011.

[17]  A Gryparis, C Paciorek, A Zeka, J Schwartz, and B Coull. Measurement error caused by spatial misalignment in environmental epidemiology. *Biostatistics*, 10:258–274, 2009.

[18]  R Haining, G Li, R Maheswaran, M Blangiardo, J Law, N Best, and S Richardson. 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, 2010.

[19]  J Hughes and M Haran. Dimension reduction and alleviation of confounding for spatial generalized linear mixed models. *Journal of the Royal Statistical Society Series B*, 75:139–160, 2013.

[20]  H Janes, F Dominici, and S Zeger. Trends in Air Pollution and Mortality: An Approach to the Assessment of Unmeasured Confounding. *Epidemiology*, 18:416–423, 2007.

[21]  M Jerrett, M Buzzelli, R Burnett, and P DeLuca. Particulate air pollution, social confounders, and mortality in small areas of an industrial city. *Social Science and Medicine*, 60:2845–2863, 2005.

[22]  J Kaufman, S Adar, R Allen, G Barr, M Budoff, G Burke, A Casillas, M Cohen, C Curl, M Daviglus, A Roux, D Jacobs, R Kronmal, T Larsen, S Liu, T Lumley, A Navas-Acien, D O'Leary, J Rotter, P Sampson, L Sheppard, D Siscovick, J Stein, A Szpiro, and R Tracy. Prospective Study of Particulate Air Pollution Exposures, Subclinical Atherosclerosis, and Clinical Cardiovascular Disease. *American Journal of Epidemiology*, 176:825–837, 2012.

[23]  I Kleinschmidt, M Hills, and P Elliott. Smoking behaviour can be predicted by neighbourhood deprivation measures. *J Epidemiol Community Health*, 49:S72–S77, DOI:10.1136/jech.49.Suppl2.S72, 1995.

[24]  J Krall, G Anderson, F Dominici, M Bell, and R Peng. Short-term exposure to particulate matter constituents and mortality in a national study of U.S. urban communities. *Environmental Health Perspectives*, 121:1148–1153, 2013.

[25]  A Lawson, J Choi, B Cai, M Hossain, R Kirby, and J Liu. Bayesian 2-Stage Space-Time Mixture Modeling With Spatial Misalignment of the Exposure in Small Area Health Data. *Journal of Agricultural, Biological and Environmental Statistics*, 17:417–

441, 2012.

[26] D Lee. Using spline models to estimate the varying health risks from air pollution across Scotland. *Statistics in Medicine*, 31:3366–3378, 2012.

[27] D Lee, C Ferguson, and R Mitchell. Air pollution and health in Scotland: a multicity study. *Biostatistics*, 10:409–423, 2009.

[28] D Lee and R Mitchell. Controlling for localised spatio-temporal autocorrelation in long-term air pollution and health studies. *Statistical Methods in Medical Research*, to appear:DOI: 10.1177/0962280214527384, 2014.

[29] D Lee, A Rushworth, and S Sahu. A Bayesian Localized Conditional Autoregressive Model for Estimating the Health Effects of Air Pollution. *Biometrics*, 70:419–429, 2014.

[30] D Lee and G Shaddick. Spatial modeling of air pollution in studies of its short term health effects. *Biometrics*, 66:1238 – 1246, 2010.

[31] B Leroux, X Lei, and N Breslow. *Estimation of disease rates in small areas: A new mixed model for spatial dependence*, chapter Statistical Models in Epidemiology, the Environment and Clinical Trials, Halloran, M and Berry, D (eds), pages 135–178. Springer-Verlag, New York, 1999.

[32] P Li, S banerjee, and A McBean. Mining boundary effects in areally referenced spatial data using the Bayesian information criterion. *Geoinformatica*, 15:435–454, 2011.

[33] H Lu, C Reilly, S Banerjee, and B Carlin. Bayesian areal wombling via adjacency modelling. *Environ Ecol Stat*, 14:433–452, 2007.

[34] H Ma, B Carlin, and S Banerjee. Hierarchical and Joint Site-Edge Methods for Medicare Hospice Service Region Boundary Analysis. *Biometrics*, 66:355–364, 2010.

[35] R Maheswaran, R Haining, P Brindley, J Law, T Pearson, P Fryers, S Wise, and M Campbell. Outdoor air pollution and stroke in Sheffield, United Kingdom. *Stroke*, 36:239–243, 2005.

[36] B Reich, J Hodges, and V Zadnik. Effects of Residual Smoothing on the Posterior of the Fixed Effects in Disease-Mapping Models. *Biometrics*, 62:1197–1206, 2006.

[37] E Richardson and R Mitchell. Gender differences in relationships between urban green space and health in the united kingdom. *Social Science and Medicine*, 71:568–575, 2010.

[38] S Richardson, I Stucker, and D Hemon. Comparison of Relative Risks Obtained in Ecological and Individual Studies: Some Methodological Considerations. *International Journal of Epidemiology*, 16:111–120, 1987.

[39] C Robert and G Casella. *Introducing Monte Carlo Methods with R*. Springer, New York, 1st edition, 2010.

[40] H Rue, S Martino, and N Chopin. Approximate bayesian inference for latent gaussian models using integrated nested laplace approximations (with discussion). *J Roy Stat Soc B*, 71:319–392, 2009.

[41] S K Sahu, S Yip, and D M Holland. Improved space-time forecasting of next day ozone concentrations in the eastern u.s. *Atmospheric Environment*, 43:494–501, 2009.

[42] R Salway and J Wakefield. A hybrid model for reducing ecological bias. *Biostatistics*, 9:1–17, 2008.

[43] E Samoli, G Touloumi, J Schwartz, R Anderson, C Schindler, B Forsberg, M Vigotti, J Vonk, M Kosnik, J Skorkovsky, and K Katsouyanni. Short-Term Effects of Carbon Monoxide on Mortality: An Analysis within the APHEA Project. *Environmental Health Perspectives*, 115:1578–1583, 2007.

[44] N H Savage, P Agnew, L S Davis, C Ordonez, R Thorpe, C E Johnson, F M O'Connor, and Dalvi M. Air quality modelling using the met office unified model (aqum os24-26): model description and initial evaluation. *Geoscientific Model Development*, 6:353–372, 2013.

[45] J Schwartz and A Marcus. Mortality and Air Pollution in London: A Time Series Analysis. *American Journal of Epidemiology*, 131:185–194, 1990.

[46] A Szpiro and C Paciorek. Measurement error in two-stage analyses, with application to

air pollution epidemiology. *Environmetrics*, 24:501–517, 2013.

[47] J Wakefield. Disease mapping and spatial regression with count data. *Biostatistics*, 8:158–183, 2007.

[48] J Wakefield and R Salway. A statistical framework for ecological and aggregate studies. *Journal of the Royal Statistical Society Series A*, 164:119–137, 2001.

[49] J Wakefield and G Shaddick. Health-exposure modeling and the ecological fallacy. *Biostatistics*, 7:438–455, 2006.

[50] L Wymer and A Dufour. A model for estimating the incidence of swimming-related gastrointestinal illness as a function of water quality indicators. *Environmetrics*, 13:669–678, 2002.

[51] L Young, C Gotway, J Yang, G Kearney, and C DuClos. Linking health and environmental data in geographical analysis: It's so much more than centroids. *Spatial and Spatio-temporal Epidemiology*, 1:73–84, 2009.

[52] S Zeger, D Thomas, F Dominici, J Samet, J Schwartz, D Dockery, and A Cohen. Exposure Measurement Error in Time-Series Studies of Air Pollution: Concepts and Consequences. *Environmental Health Perspectives*, 108:419–426, 2000.

[53] L Zhu, B Carlin, and A Gelfand. Hierarchical regression with misaligned spatial data: relating ambient ozone and pediatric asthma ER visits in Atlanta. *Environmetrics*, 14:537–557, 2003.

[54] J V Zidek, N D Le, and Z Liu. Combining data and simulated data for space-time fields: application to ozone. *Environmental and Ecological Statistics*, 19:37–56, 2012.

# Acknowledgements