

# Using small-area spatial data for statistical and epidemiological research

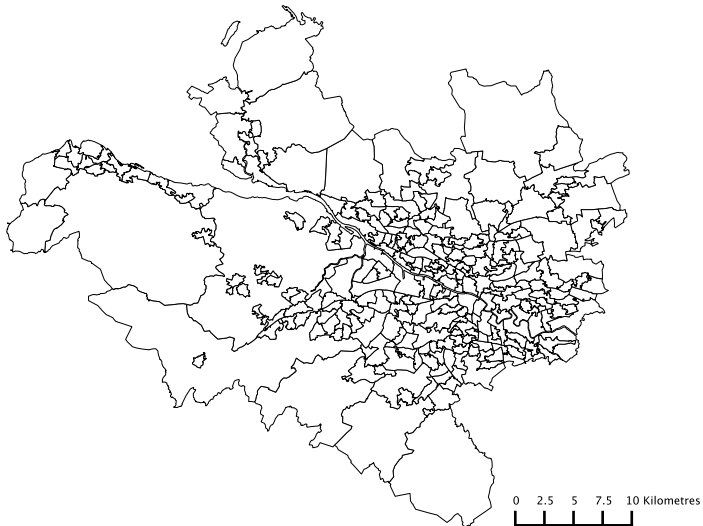
Duncan Lee [Duncan.Lee@glasgow.ac.uk](mailto:Duncan.Lee@glasgow.ac.uk)

Scottish Government 16th December 2015

This is joint work with colleagues from two main projects:

- Andrew Lawson, Gary Napier, Kevin Pollock and Chris Robertson on an MRC funded project focusing on modelling spatio-temporal patterns in disease risk.
- Guowen Huang and Marian Scott on an EPSRC funded project focusing on the effects of air pollution on human health.

- The availability of small-area spatial data has dramatically increased in the last decade or so, including:
  - Scottish Neighbourhood Statistics (SNS).
  - Health and Social Care Information Centre (HSCIC).
- This increase has been accompanied by the widespread development of statistical methodology and software for mapping and modelling these data.
- The latter include Geographical Information Systems software and statistical modelling software, such as ArcGIS / QGIS, and packages in the statistical software R such as CARBayes and RINLA.
- The analysis of small-area data is performed in many fields, such as econometrics, epidemiology, and social science.



The key feature when modelling spatial data is that of **spatial autocorrelation**, which is summarised by Tobler's first law of geography which says

Everything is related to everything else, but near things are more related than distant things.

Which means that standard regression models that assume independence in the residuals are likely to be inappropriate and potentially result in misleading conclusions.

There are a number of commonly used models for capturing spatial autocorrelation in data / residuals from a regression model, including:

- Conditional AutoRegressive (CAR) models.
- Simultaneous AutoRegressive (SAR) models.
- Point level (Geostatistical) models based on each areas centroid.

CAR and SAR models are most commonly used for small-area data.

All models require the spatial closeness between each pair of areal units to be defined, and CAR and SAR models use an  $K \times K$  neighbourhood or adjacency matrix  $\mathbf{W}$ , where  $K$  is the number of areal units in the data set. Then, the  $kj$ th element of  $\mathbf{W}$  is typically defined as:

$$w_{kj} = \begin{cases} 1 & \text{Areas (k,j) share a border} \\ 0 & \text{Otherwise} \end{cases}$$

so that if  $w_{kj} = 1$  data in areal units  $(k, j)$  are modelled as spatially autocorrelated, while if  $w_{kj} = 0$  they are assumed to be conditionally independent.

Most of my research is interested in developing new statistical methodology for epidemiological applications:

- Identifying trends in disease risk over time and how those trends vary in space.
- Identifying the locations of clusters of small-areas that exhibit substantially higher disease risk than their neighbours.
- Estimating the effects of air pollution on human health.

However, I am also involved in a social science led project AQMeN (<http://aqmen.ac.uk/>) looking at changes in urban segregation over time.



My choice of research topics is motivated by two main goals:

- 1 Finding an important public health problem for which existing statistical models are inadequate, and hence providing an innovative modelling solution.
- 2 Developing free to use software in R for others to be able to apply both standard and novel statistical models for small-area data to their own problems.

Examples of the latter include the `CARBayes` and `CARBayesST` packages in R for spatial and spatio-temporal modelling of areal unit data.

In this talk I give 2 examples of statistical research you can do focusing on addressing key questions in epidemiology, and using small-area data such as that from the Scottish Neighbourhood Statistics database. The examples are:

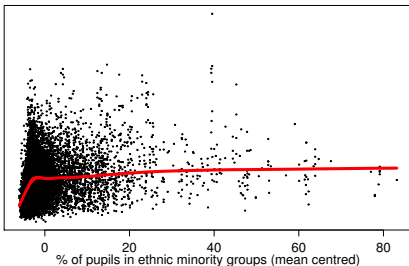
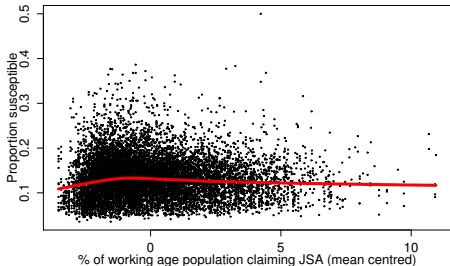
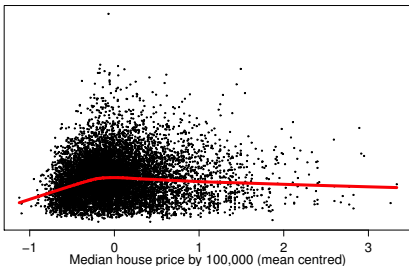
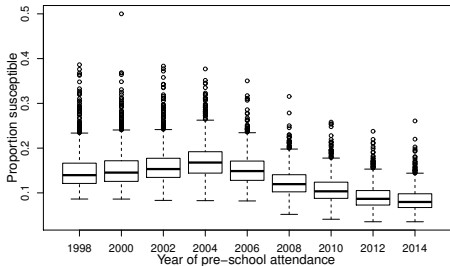
- 1 Identifying temporal trends in measles susceptibility over time.
- 2 Estimating the long-term effects of air pollution on health in Scotland.

- An article published in 1998 by Andrew Wakefield *Wakefield et al. (1998)* linked the measles, mumps and rubella (MMR) vaccine with an increased risk of autism.
- This scare led to a reduction in vaccination rates, which by 2003 reached a low of around 80% in some parts of the UK.
- Repercussions from these decreased vaccination rates were felt in 2013, with large outbreaks of measles occurring in England and Wales (*Pollock et al. 2014*).
- The Wakefield article was partially retracted in 2004 before being fully discredited in 2010 after multiple epidemiological studies failed to find any association.

- Since 1988 individual vaccination records of all children in Scotland are kept in the Scottish Immunisation & Recall System (SIRS) .
- The data analysed here are based on the estimated numbers of children eligible to attend pre-school (aged 2-4) from non-overlapping two-year birth cohorts.
- We have the estimated number susceptible to measles ( $Y_{kt}$ ), and the total number of pre-school children ( $N_{kt}$ ) in the 1235 (indexed by  $k$ ) intermediate geographies in Scotland between 1998 and 2014 (indexed by  $t$ ).
- Thus  $\hat{\theta}_{kt} = Y_{kt}/N_{kt}$  is the estimated proportion of children who are susceptible to measles, where measles susceptibility is based upon the receipt of one or two vaccinations that each have a failure rate of 10%.

We use three covariate factors here to capture the potential impacts of socio-economic deprivation and ethnicity:

- Median House Price (MHP) in each IG and year.
- Percentage of working age people in receipt of Job Seekers Allowance (JSA) in each IG and year.
- Percentage of school children from ethnic minorities (EM).



Initially a simple binomial logistic regression model was fitted to these data:

$$Y_{kt} \sim \text{Binomial}(N_{kt}, \theta_{kt}),$$
$$\ln \left( \frac{\theta_{kt}}{1 - \theta_{kt}} \right) = \beta_0 + S_1(\text{MHP}_{kt}) + S_2(\text{JSA}_{kt}) + S_3(\text{EM}_{kt}),$$

where  $S_i(\cdot) = \sum_{j=1}^3 B_i(\cdot)\beta_{ij}$ ,  $i = 1, 2, 3$  is a natural cubic spline of each covariate with 3 degrees of freedom to allow for the non-linear relationships observed above.

Spatial autocorrelation in the residuals was assessed by computing Moran's I statistics and performing permutation tests on the residuals from the above model separately for each year.

The values of the Moran's I statistics obtained ranged between 0.12 and 0.38, with the statistics generated from 10 000 random permutations of the data yielding  $p$ -values of less than 0.0001 in all cases. Thus spatial autocorrelation has to be modelled.



The goals of the analysis are as follows:

- 1 Estimate the magnitude of the increased measles susceptibility associated with the MMR vaccination scare linked to the Wakefield article and assess whether susceptibility has decreased in recent years.
- 2 Assess whether the magnitude of the inequality, as measured by the spatial variability, in measles susceptibility in Scotland increased at the same time and whether spatial variation has now decreased.
- 3 Determine whether any area based covariates, such as deprivation, have any impact on measles susceptibility in Scotland.

We propose a Bayesian hierarchical model to answer these questions, which is given by:

$$Y_{kt} | N_{kt}, \theta_{kt} \sim \text{Binomial}(N_{kt}, \theta_{kt}),$$
$$\ln \left( \frac{\theta_{kt}}{1 - \theta_{kt}} \right) = \beta_0 + S_1(\text{MHP}_{kt}) + S_2(\text{JSA}_{kt}) + S_3(\text{EM}_{kt}) + \phi_{kt} + \delta_t.$$

where  $\delta_t$  is the Scotland-wide temporal trend at time  $t$  while  $\phi_{kt}$  is a spatio-temporal effect for area  $k$  and time  $t$ .

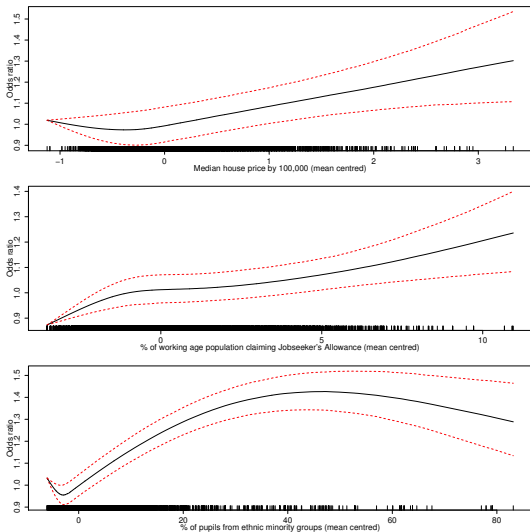
- The overall temporal trend is modelled by a first order random walk process:

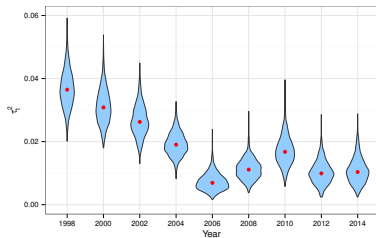
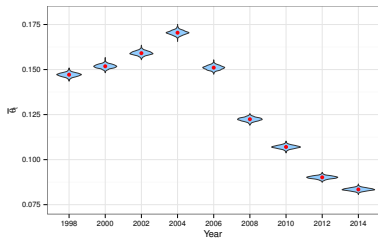
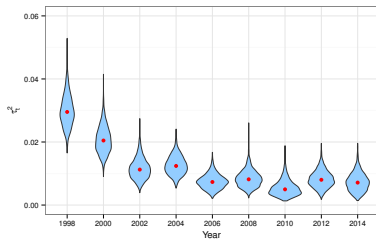
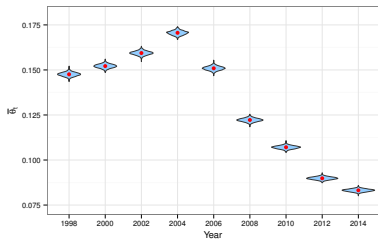
$$\delta_t \sim \text{N}(\delta_{t-1}, \sigma^2).$$

- The spatial trend in year  $t$ ,  $\phi_t = (\phi_{1t}, \dots, \phi_{Kt})$  is modelled by a conditional autoregressive (CAR) model:

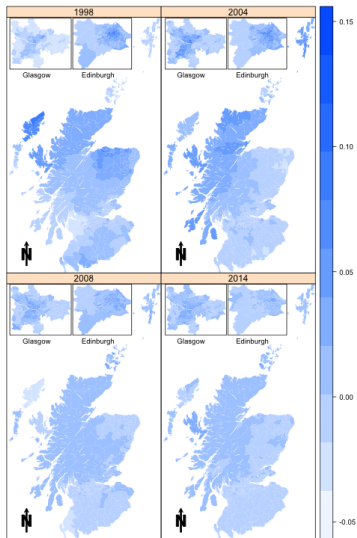
$$\phi_{kt} | \phi_{-kt}, \mathbf{W}, \rho, \tau_t^2 \sim \text{N} \left( \frac{\rho \sum_{j=1}^K w_{kj} \phi_{jt}}{\rho \sum_{j=1}^K w_{kj} + 1 - \rho}, \frac{\tau_t^2}{\rho \sum_{j=1}^K w_{kj} + 1 - \rho} \right),$$

where each year has its own variance  $\tau_t^2$  to allow for temporally varying spatial variability.





The top row is with covariates and the bottom panel is without.

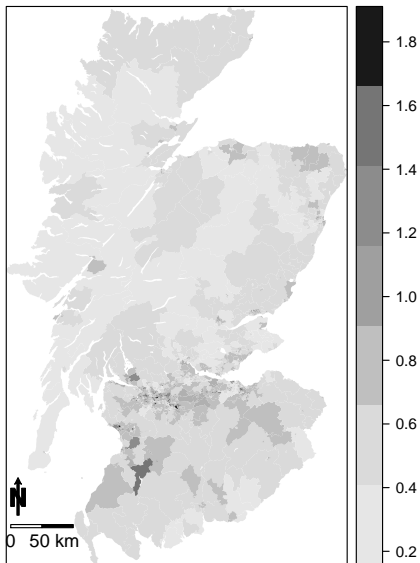


- Measles susceptibility increased to a peak in 2004 coinciding with the media coverage surrounding the Wakefield article, before dropping dramatically until the present day.
- Spatial variation in measles susceptibility decreased until 2006, whereafter it has stayed relatively constant over time.
- Socio-economic deprivation appears to have a *U-shaped* relationship with measles susceptibility, with increasing susceptibility for very poor and very affluent communities.
- Spatially, the rural northwest part of Scotland appears to have the highest rates of susceptibility, as it has stayed consistently higher than other parts of the country for all time periods.

- Air pollution has long been known to adversely affect public health, in both the developed and developing world.
- Recent reports by the UK government and the World Health Organisation estimate that:
  - particulate matter reduces life expectancy by 6 months, with a health cost of £19 billion per year.
  - there were estimated to be over 23,000 premature deaths from air pollution in 2010.
- Air pollution will remain a key health problem for some time, as nitrogen dioxide emissions are predicted to exceed European Union limits until after 2020 in key parts of the UK, including in Glasgow.



- In ecological studies the data relate to populations living in a set of  $k = 1, \dots, K$  non-overlapping areal units for  $t = 1, \dots, T$  time periods, rather than to individuals.
- In this study we have  $K = 1207$  Intermediate Geographies that make up mainland Scotland, and data are collected for  $T = 5$  years between 2007 and 2011.
- For IG  $k$  and year  $t$  the observed number of hospital admissions due to respiratory disease is denoted by  $Y_{kt}$ , while the expected number of admissions based on population demographics is denoted by  $E_{kt}$ .
- The standardised morbidity ratio is given by  $SMR_{kt} = Y_{kt}/E_{kt}$ , an exploratory measure of disease risk.

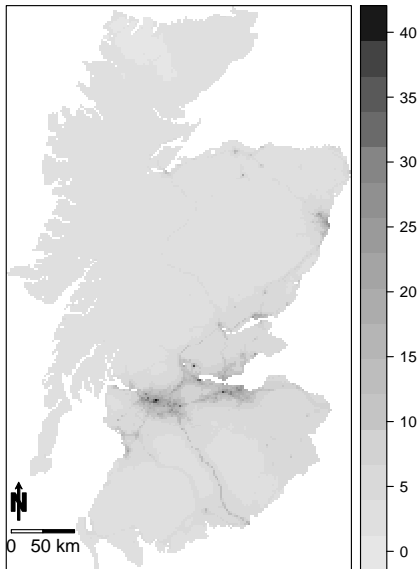


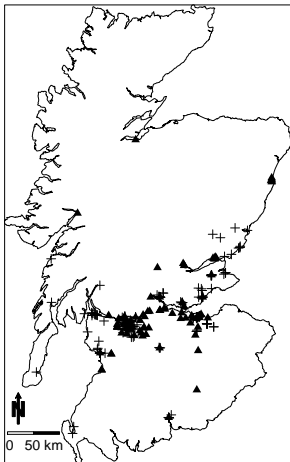
Air pollution concentration data are available from two sources:

- Measured data from a small number of monitoring sites such as the Automatic Urban and Rural Network (AURN).
- Modelled concentrations from a dispersion model on a 1kilometre regular grid, such as the data provided by DEFRA.

Neither data source is ideal, as the measured data are spatially sparse and the modelled data are only estimates and not real measurements.

# Modelled nitrogen dioxide (NO<sub>2</sub>) data





Triangles are monitors and + are diffusion tubes.

The first goal in this study is to estimate representative concentrations of nitrogen dioxide for each IG from these two sets of data. However, there are two natural questions one has to overcome.

- How should the different spatial scales of the three data sets be resolved (point, grid square, IG)?
- How can we use both the modelled and monitored data to get the *best* estimate of NO<sub>2</sub> concentrations.

- The accepted approach in the literature is a statistical fusion model, which essentially regresses the measured pollution data against the modelled pollution data.
- Let  $\mathbf{X}_t = (X_t(\mathbf{s}_1), \dots, X_t(\mathbf{s}_{n_t}))$  denote the vector of  $n_t$  measured  $\text{NO}_2$  concentrations (on the natural log scale) at sites  $(\mathbf{s}_1, \dots, \mathbf{s}_{n_t})$  in year  $t$
- These measured concentrations are related to an  $n_t \times p$  design matrix of covariates  $\mathbf{Z}_t$ , such as the modelled concentrations, site type, etc.

Then the regression component of the model is given by:

$$X_t \sim N(\mathbf{Z}_t \boldsymbol{\beta}_t, \sigma_t^2 \mathbf{I}_t) \quad t = 1, \dots, T,$$

The regression parameters  $\boldsymbol{\beta}_t$  and the variance parameter  $\sigma_t^2$  vary over time via autoregressive processes:

$$\begin{aligned} \boldsymbol{\beta}_t &\sim N(\boldsymbol{\beta} + \kappa(\boldsymbol{\beta}_{t-1} - \boldsymbol{\beta}), \tau^2 \mathbf{V}) \quad t = 2, \dots, T, \\ \ln(\sigma_t^2) &\sim N(\ln(\sigma_{t-1}^2), \sigma^2) \quad t = 2, \dots, T, \end{aligned}$$

Note, no spatial autocorrelation is allowed for here because the modelled concentrations are spatially autocorrelated and account for it.



- The model is fitted in a Bayesian setting, using Markov Chain Monte Carlo (MCMC) simulation methods.
- It predicts  $\text{NO}_2$  concentrations at each 1km grid square in mainland Scotland.
- We aggregate to the IG level by computing the mean and maximum of the set of 1km grid squares in each IG.
- The mean is commonly used in these studies, while the maximum may better represented densely populated parts of each IG.

Recall that  $(Y_{kt}, E_{kt})$  denote the observed and expected numbers of hospital admissions in the  $k$ th IG in year  $t$ . Then a Poisson generalised linear mixed model is typically used:

$$\begin{aligned}
 Y_{kt} \mid E_{kt}, R_{kt} &\sim \text{Poisson}(E_{kt}R_{kt}), \\
 \ln(R_{kt}) &= \mathbf{b}_{kt}^T \boldsymbol{\alpha} + \tilde{X}_{kt} \lambda + \phi_{kt},
 \end{aligned}$$

where  $\tilde{X}_{kt}$  is NO<sub>2</sub> (mean or max) and  $\mathbf{b}_{kt}^T$  are other covariates such as socio-economic deprivation. Finally,  $\phi_{kt}$  is a random effect that accounts for any unmeasured spatio-temporal autocorrelation.

Our basic model for spatial autocorrelation is a combination of an autoregressive time series process of order 1 and a conditional autoregressive spatial process, and is given by

$$\phi_t \mid \phi_{t-1} \sim \mathbf{N}(\gamma\phi_{t-1}, \nu^2\mathbf{Q}(\rho, \mathbf{W})^{-1}),$$

where  $\mathbf{W}$  is a spatial neighbourhood matrix and  $(\gamma, \rho)$  are temporal and spatial autocorrelation parameters respectively. Again, the model is fitted in a Bayesian framework using MCMC simulation, and the R package `CARBayesST` can be used to fit this model.

The results are presented as relative risks for a standard deviation increase in each covariates value, which is  $\text{NO}_2$  6.84  $\mu\text{gm}^{-3}$ , Logprice 0.38, JSA 2.35.

<b>Parameter</b>	<b>Spatial mean <math>\text{NO}_2</math></b>	<b>Spatial max <math>\text{NO}_2</math></b>
$\text{NO}_2$	0.993 (0.980,1.008)	1.021 (1.004,1.037)
Logprice	0.920 (0.909,0.929)	0.921 (0.911,0.930)
JSA	1.200 (1.185,1.215)	1.196 (1.180,1.214)
$\rho$	0.926 (0.891,0.956)	0.911 (0.866,0.946)
$\gamma$	0.832 (0.797,0.867)	0.830 (0.792,0.865)

- In this study the choice of spatially representative measure of  $\text{NO}_2$  had a large impact on the results, with no effect being seen for the spatial mean but a substantial effect being observed for the spatial max.
- Estimating the effect of air pollution on health is a difficult task because of factors such as:
  - Spatial misalignment between the pollution and disease data.
  - The difficult task of controlling for residual spatio-temporal autocorrelation.
  - The lack of real *exposure* data.
- Future work in this field should look at personal exposures and not outdoor measured data at fixed locations.

- In summary, databases of small-area statistics such as SNS provide a valuable resource for addressing key questions of public and policy interest.
- Although this talk has focused on epidemiology, the SNS has data on topics as diverse as property prices, educational attainment, benefit claimants, access to services and demography.
- A key statistical issue is that of spatial autocorrelation, in that simple models ignoring this are likely to produce biased results, particularly in terms of confidence intervals.
- A second issue is that of the modifiable areal unit problem (MAUP), which essentially means do your inferences remain valid if you change spatial scale, e.g. from IG to DZ?

- *An integrated Bayesian model for estimating the long-term health effects of air pollution by fusing modelled and measured pollution data: A case study of nitrogen dioxide concentrations in Scotland*, Guowen Huang, et al., *Spatial and Spatio-temporal Epidemiology*, 2015, 14-15, 63-74.
- *A model to estimate the impact of changes in MMR vaccine uptake on inequalities in measles susceptibility in Scotland*, Gary, Napier, et al, under revision.
- CARBayes and CARBayesST R packages are available from <https://cran.r-project.org>.
- Further details about my research can be found at <http://www.gla.ac.uk/schools/mathematicsstatistics/staff/duncanlee/>.