

Quantifying the long-term effects of air pollution on health in England using space time data

Duncan Lee

Geomed 2015

10th September 2015

- This is joint work with Alastair Rushworth from the University of Glasgow, Sujit Sahu and Sabyasachi Mukhopadhyay from the University of Southampton, and Paul Agnew from the UK Met Office.
- The work is funded by the EPSRC grants EP/J017442/1 and EP/J017485/1.

The logo for EPSRC consists of the letters 'EPSRC' in a bold, purple, sans-serif font. The letters are framed by two horizontal teal lines, one above and one below.

Engineering and Physical Sciences
Research Council

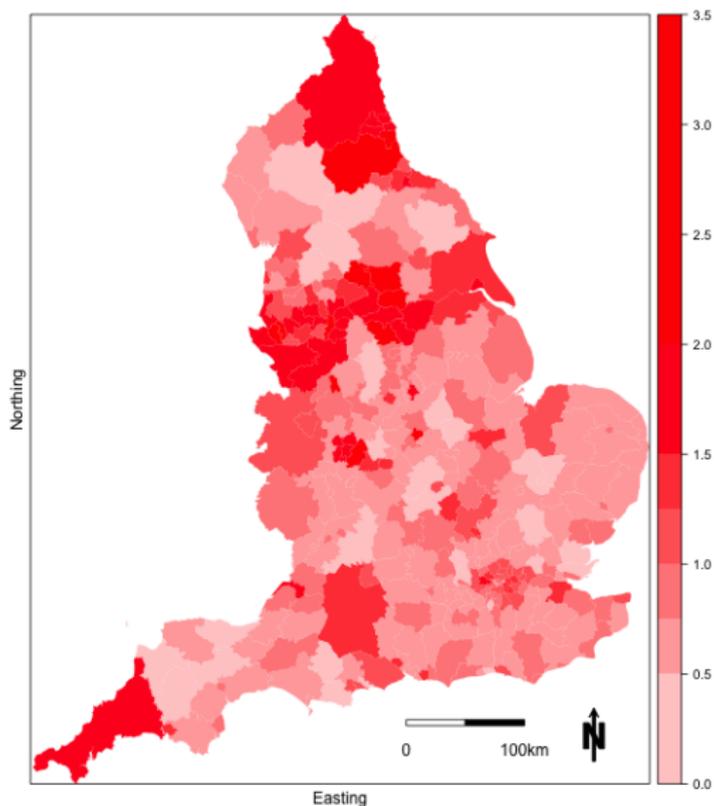
- Air pollution has long been known to adversely affect public health, in both the developed and developing world.
- Recent reports by the UK government and the World Health Organisation estimate that:
 - particulate matter reduces life expectancy by 6 months, with a health cost of £19 billion per year.
 - 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 2030 in the urban areas of Greater London, West Midlands and West Yorkshire.

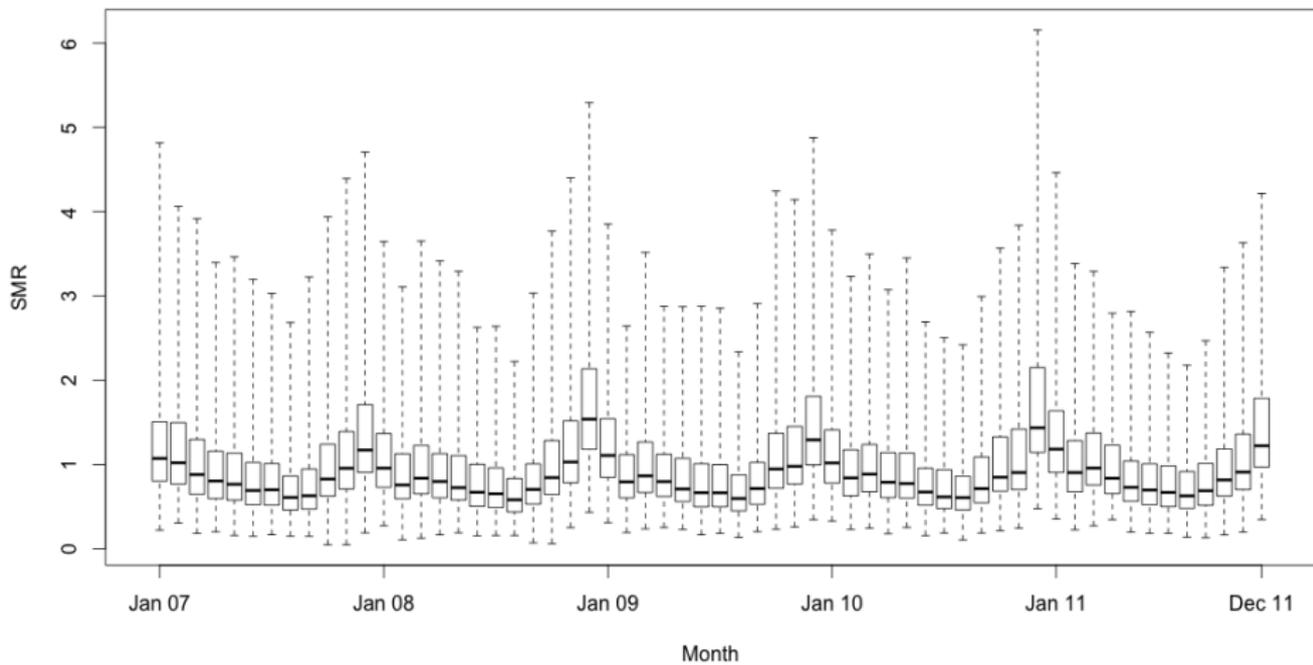
There are two main study designs when investigating the effects of long-term exposure to air pollution.

Cohort studies e.g. The Six Cities study by Dockery *et al* (1993) and the Escape study by Beelen *et al* (2014), which relate average air pollution concentrations to the health status of a large pre-defined cohort of people.

Ecological studies e.g. Elliot *et al* (2007) and Lee *et al* (2009), which relate average air pollution concentrations in contiguous small areas (such as electoral wards), against yearly numbers of health events from the population living in that area.

- 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 $K = 323$ local and unitary authorities (LUA) that make up mainland England, and data are collected for $T = 60$ months between 2007 and 2011.
- For LUA k and month 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 mortality ratio is given by $SMR_{kt} = Y_{kt}/E_{kt}$, an exploratory measure of disease risk.





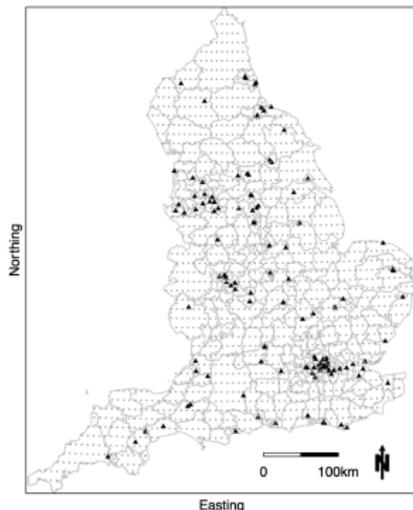
- The SMR exhibit a strong seasonal pattern, and this was accounted for by seasonally adjusting the expected disease counts E_{kt} by a monthly correction factor.
- The main confounding factor is socio-economic deprivation, and this is accounted for using two proxy variables, the percentage of people in each LUA and month who are in receipt of Job Seekers Allowance (JSA), and the average property price in each LUA and month.

Recall from the previous talk that we consider concentrations of the following pollutants:

- Nitrogen dioxide (NO_2).
- Particulate matter measured as $\text{PM}_{2.5}$ and PM_{10} .
- Ozone (O_3).

For simplicity we only consider the first 3 here. Data on these pollutants come from the AURN monitoring network and the AQUM modelled concentrations provided by the Met Office.

Recall from the previous talk that pollution concentrations were predicted on a 12km by 12km regular grid across England using a Bayesian hierarchical model, resulting in LUA k and month t having matrix of $l = 1, \dots, 5000$ pollution predictions $\{z^{(l)}(\mathbf{v}_{kj}, t)\}$ at prediction locations $(\mathbf{v}_{k1}, \dots, \mathbf{v}_{kn_k})$.



For LUA k and month t we have $5000 \times n_k$ predictions from the pollution model, and the simplest approach is to average them, that is

$$\hat{z}_{kt} = \frac{1}{Ln_k} \sum_{j=1}^{n_k} \sum_{l=1}^L z^{(l)}(\mathbf{v}_{kj}, t)$$

This ignores two different sources of uncertainty in the pollution concentrations:

- Spatial variation in pollution within each LUA.
- Posterior uncertainty in the pollution predictions at each prediction location.

A typical Bayesian hierarchical model for the disease data is

$$\begin{aligned}
 Y_{kt} | E_{kt}, R_{kt} &\sim \text{Poisson}(E_{kt} R_{kt}) \\
 \ln(R_{kt}) &= \beta_0 + \hat{z}_{kt} \beta + \mathbf{x}_{kt}^\top \boldsymbol{\beta}_x + \phi_{kt},
 \end{aligned}$$

where

Y_{kt}	health counts	E_{kt}	expected cases
R_{kt}	health risk	ϕ_{kt}	random effect
$\mathbf{x}_{kt}^\top \boldsymbol{\beta}_x$	other covariate effects	β_0	intercept
\hat{z}_{kt}	air pollution	β	pollution effect

1. Unmeasured confounding: Air pollution and the other covariates do not account for all variation. Adding a set of spatio-temporal random effects, ϕ_{kt} can offer a solution.

How should ϕ_{kt} be structured in space and time?

2. Spatial Misalignment: The air pollution model estimates the true exposure surface $Z(v_{kj}, t)$ at grid locations, $\{v_{kj}\}$ not as LUA regional averages.

How can we reconcile these quantities?

3. Uncertainty: The posterior distribution of $Z(v_{kj}, t)$ is available via MCMC samples, and therefore uncertainty in air pollution is quantified.

How should this source of uncertainty be incorporated into the health model?

Rushworth et al. (2014) propose the ‘global’ autoregressive conditional autoregressive (CAR) model for the random effects $\phi_t = (\phi_{1t}, \dots, \phi_{Kt})$ at time t :

$$\begin{aligned}\phi_1 &\sim N(\mathbf{0}, \sigma^2 \mathbf{Q}(\mathbf{W}, \rho)^{-1}) \\ \phi_t | \phi_{t-1} &\sim N(\alpha \phi_{t-1}, \sigma^2 \mathbf{Q}(\mathbf{W}, \rho)^{-1}) \quad \text{for } t \geq 2\end{aligned}$$

where

$$\begin{aligned}\mathbf{Q}(\mathbf{W}, \rho) &= \rho [\text{diag}(\mathbf{W}\mathbf{1}) - \mathbf{W}] + (1 - \rho)\mathbf{I} \\ \mathbf{W} &= \text{spatial (binary) neighbours matrix.}\end{aligned}$$

\mathbf{W} is a $K \times K$ matrix that encodes neighbourhood relationships in the study region such that

$$\begin{aligned}
 w_{kj} = 1 & \iff \text{units } k \text{ and } j \text{ share a common border} \\
 w_{kj} = 0 & \text{ otherwise, or if } k = j
 \end{aligned}$$

If $T = 1$ then the conditional distribution for ϕ_{k1} is

$$\phi_{k1} | \boldsymbol{\phi}_{-k1} \sim \mathbf{N} \left(\frac{\rho \sum_{j=1}^n w_{kj} \phi_{j1}}{1 - \rho + \rho \sum_{j=1}^n w_{kj}}, \frac{\tau^2}{1 - \rho + \rho \sum_{j=1}^n w_{kj}} \right)$$

$\mathbf{Q}(\mathbf{W}, \rho)$ restricts the range of surfaces that can be fitted, as ρ controls the level of spatial autocorrelation globally across the entire region.

Therefore we treat the non-zero elements of \mathbf{W} as random variables $w_{kj}^+ \in [0, 1]$. We control model complexity using a normal shrinkage prior on the logit transformed w_{kj}^+ :

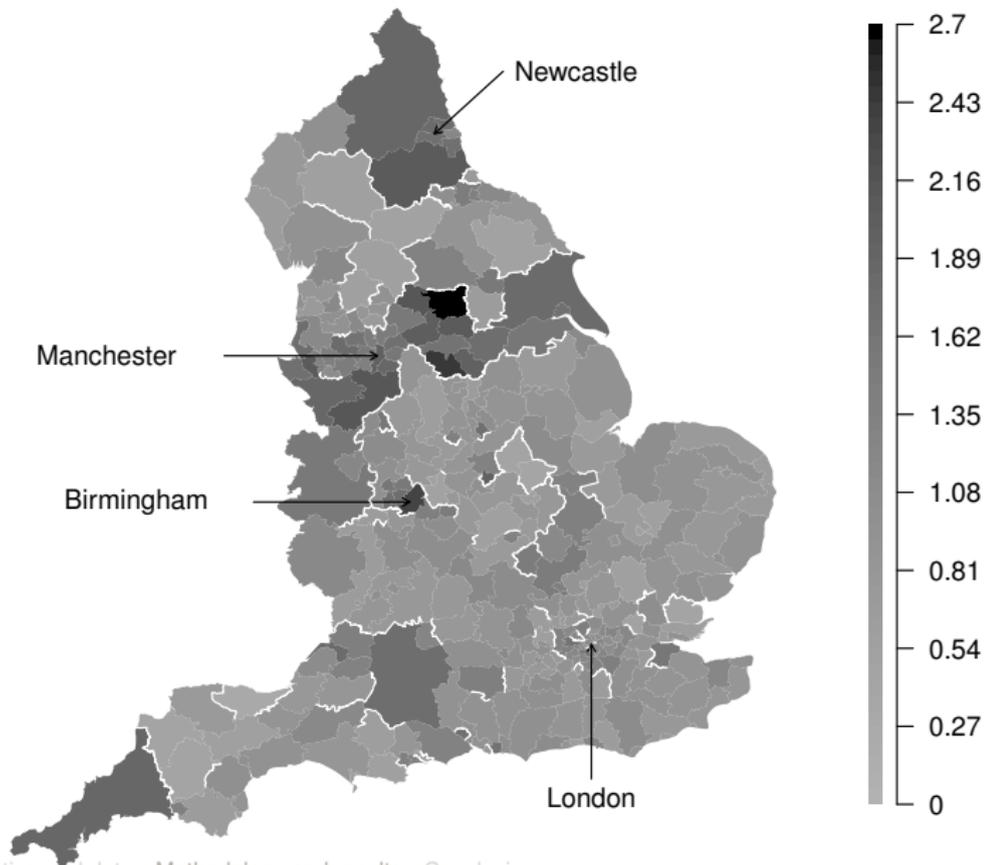
$$\ln \left(\frac{w_{kj}^+}{1 - w_{kj}^+} \right) \sim \mathbf{N}(\mu, \tau^2)$$

If μ is positive then w_{kj}^+ are *a-priori* close to one reflecting prior preference for spatial smoothness.

Table: Posterior median relative risks and 95% credible intervals for a 1-standard deviation increase in each pollutant, which is 16.07, 4.90, and $4.11 \mu\text{gm}^{-3}$ respectively.

Pollutant	No random effects	Non-adaptive ϕ_{kt}	Adaptive ϕ_{kt}
NO ₂	1.151 (1.144, 1.158)	1.057 (1.045, 1.069)	1.048 (1.036, 1.060)
PM ₁₀	1.013 (1.007, 1.020)	1.007 (0.998, 1.015)	1.006 (0.995, 1.015)
PM _{2.5}	1.013 (1.007, 1.019)	1.006 (0.997, 1.014)	1.006 (0.997, 1.016)

Simpler models have a tendency to have larger estimated air-pollution effects.



- The spatial misalignment between the pollution and disease data means there is within LUA variation in pollution which can result in ecological bias.
- Wakefield and Shaddick (2006) among others derive an appropriate aggregate model that overcomes this ecological bias.
- Under simplifying assumptions the bias term can be derived to be of the order of β^2 , thus as β is small here the effect of this bias is likely to be negligible.
- This negligible effect was confirmed empirically by Lee and Sarran (2015) in this air pollution and health context, and the naive ecological model and aggregate model similar to Wakefield and Shaddick (2006) give identical results.

The pollution model yields predictive distributions, based on 5,000 MCMC samples, for average air pollution in each LUA and month.

This uncertainty should be fed into the health model so that the resulting health estimates account for this variation.

3 possible strategies:

- (1)** Treat posterior mean pollution concentrations as true values (no uncertainty).
- (2)** Directly feed samples from the posterior air pollution density through the health model.
- (3)** Treat the posterior pollution densities as prior distributions in the health model (e.g. using a Gaussian approximation).

Table: Posterior median relative risks and 95% credible intervals for a 1-standard deviation increase in each pollutant, which is 16.07, 4.90, and $4.11 \mu\text{gm}^{-3}$ respectively.

Pollutant	No uncertainty	Posterior	Prior
NO ₂	1.048 (1.036, 1.060)	1.001 (0.999, 1.003)	1.035 (1.030, 1.041)
PM ₁₀	1.006 (0.995, 1.015)	1.000 (0.998, 1.003)	1.025 (0.999, 1.043)
PM _{2.5}	1.006 (0.997, 1.016)	1.001 (0.997, 1.004)	1.008 (0.995, 1.062)

- Choices for handling spatio-temporal autocorrelation have important consequences for the estimated effects of air pollution.
- It is important to treat air pollution exposure as uncertain, as it is rarely realistic to assume exposure is observed (or predicted) without error.
- From our study it appears that NO_2 poses the greatest ongoing health risks, and a 3.5% increased risk corresponds to around 21,500 more admissions per year.
- Given that large parts of England are expected to exceed EU emissions targets for the next 15 years, NO_2 continues to be a major health problem.

Air pollution effect estimates are small and can be volatile, so we will conduct a large sensitivity analysis of our study results looking at how sensitive the estimated pollution effects are to:

- the pollution model described in the previous talk.
- the use of monthly mean concentrations rather than monthly extremes or exceedences of a threshold level.
- the choice of variables used to control for the confounding effects of socio-economic deprivation.
- The choice of random effects model to account for spatio-temporal autocorrelation.