

# Downloadable estimates of air pollution for England and Wales and estimation of their health effects

Sujit Sahu

<http://www.soton.ac.uk/~sks/>

UNIVERSITY OF  
Southampton

- 1 Mukhopadhyay, S. and Sahu, S. K. (2017) *J. of the Royal Statistical Society, Series A*, doi:10.1111/rssa.12299.
- 2 Lee, D., Mukhopadhyay, S., Rushworth, A. and Sahu, S. K. (2016) *Biostatistics*, doi:10.1093/biostatistics/kxw048.

RSS Webinar, February 2018

**EPSRC**

Engineering and Physical Sciences  
Research Council



University  
of Glasgow

# Pollution is still a problem today!

**BBC** Sign In News Sport Weather iPlayer TV Ra

**NEWS UK**

Home World **UK** England N. Ireland Scotland Wales Business Politics Health Education Sci/En

3 April 2014 Last updated at 22:15

Share f t e

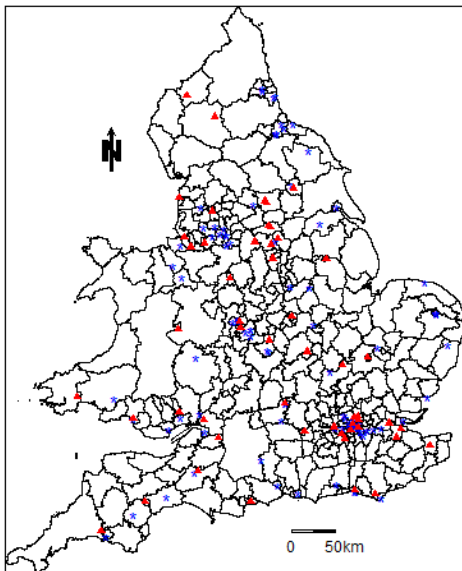
## Air pollution: Forecasters hope for cleaner air on Friday



People with lung and heart problems have been advised to avoid strenuous outdoor activity

- London kids on high air pollution: 'Our eyes start stinging' BBC News, 29 January 2017.
- Traffic pollution kills 5,000 a year in UK, says study. BBC News, 17 April 2012.

# Automatic Urban and Rural Network (AURN)



- Map of 346 local and unitary authorities (LUA) in Eng & Wales.
- 144 AURN monitoring locations are blue \* and red Δ.
- **Statistical modelling challenges:** how do we estimate air pollution at *any* new point or LUA?

- How do we relate health outcome data and pollution?

# UK air pollution data

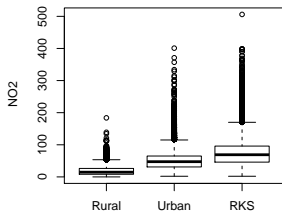
- Monitoring data can be downloaded from the website hosted by DEFRA (Dept for Env Food and Rural Affairs).
- Data are very sparse with a lot of missing data.

Pollutant	2007	2008	2009	2010	2011	Overall
NO <sub>2</sub>	31311	31356	31815	31828	33224	159,534
O <sub>3</sub>	22528	19015	18561	18786	19738	98,628
PM <sub>10</sub>	17783	16939	15240	13968	15297	79,227
PM <sub>2.5</sub>	1754	4121	16725	17667	17910	58,177

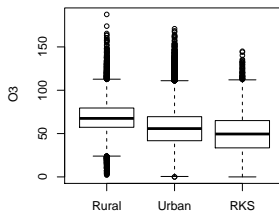
**Table:** Number of available observations out of the total number of observations in a year, which is 52704 ( $366 \times 144$ ) for 2008 and 52560 ( $365 \times 144$ ) for the other years. A 2008 EU directive triggered PM<sub>2.5</sub> monitoring in 2009. What will happen after Britain leaves EU?

- The 144 sites were classified into three types: Rural, Urban and RKS (Road and Kerbside).

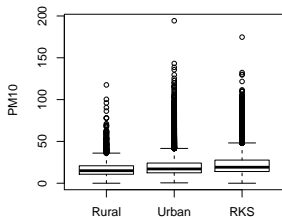
# How do the data look like?



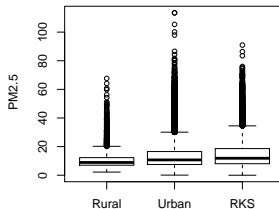
(a): NO<sub>2</sub>



(b): O<sub>3</sub>



(c): PM<sub>10</sub>



(d): PM<sub>2.5</sub>

# Aims and objectives of our work

- ① To model daily levels of four major pollutants namely,  $\text{NO}_2$ ,  $\text{O}_3$ ,  $\text{PM}_{10}$  and  $\text{PM}_{2.5}$ , for the period 2007–2011.
- ② To build up a process based suitable spatio-temporal model that
  - ① can handle highly variable and sparse air pollution data.
  - ② is more accurate than recently developed methods.
  - ③ is based on a spatio-temporal process which allows us to interpolate at any unobserved point location, which in turn lets us to aggregate pollution levels in both space and time.
- In our model, we use as a covariate the output from a computer simulation model AQUM (Air Quality Unified Model) interpolated on a 1-kilometer grid.

# Spatio-temporal auto-regressive models

- General form of spatio-temporal model (books by Cressie and Wikle, 2011 and Banerjee, Carlin and Gelfand, 2015):

$$\begin{aligned}Z(\mathbf{s}, t) &= \mu(\mathbf{s}, t) + \epsilon(\mathbf{s}, t), \\ \mu(\mathbf{s}, t) &= \mathbf{x}(\mathbf{s}, t)' \boldsymbol{\beta} + \eta(\mathbf{s}, t), \\ \eta(\mathbf{s}, t) &= \rho \eta(\mathbf{s}, t-1) + \omega(\mathbf{s}, t),\end{aligned}$$

- $Z(\mathbf{s}, t)$  is the **square-root** of observed data at site  $\mathbf{s}$  and time  $t$ .
- $\boldsymbol{\beta}$  is the regression parameter, and  $\mathbf{x}(\mathbf{s}, t)$  is the covariate vector.
- $\epsilon(\mathbf{s}, t)$  is the white noise  $N(0, \sigma_{\epsilon}^2)$ , (nugget), accounting for measurement error.
- $\eta(\mathbf{s}, t)$  is the space-time interaction term, modelled by an auto-regressive Gaussian Process model.

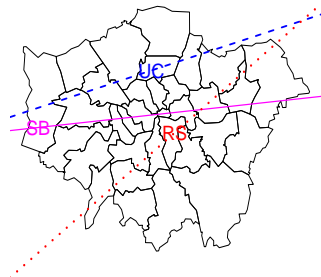
# Modelling the regression part, $\mathbf{x}(\mathbf{s}, t)' \beta$

- We allow site-wise regression lines for **Rural**, **Urban** and **RKS**. With  $x(\mathbf{s}_i, t)$  as the AQUM value, we assume:

$$\mathbf{x}(\mathbf{s}_i, t)' \beta = \sum_{k=0}^2 \delta_k(\mathbf{s}_i) (\beta_{0k} + \beta_{1k} X(\mathbf{s}_i, t)),$$

where

- $\delta_0(\mathbf{s}_i) = 1$  for all  $\mathbf{s}_i$ , ( $\beta_{00}$  and  $\beta_{10}$  overall intercept and slope.)
- and for  $k = 1, 2$ ,  $\delta_k(\mathbf{s}_i) = 1$ , if  $\mathbf{s}_i$  is of  $k$ -th type of site,  $\delta_k(\mathbf{s}_i) = 0$ , otherwise.
- Three different regression lines are obtained: one each for **Rural**, **Urban**, **RKS**.



- The model is fitted separately for each of the four pollutants,  $\text{NO}_2$ ,  $\text{O}_3$ ,  $\text{PM}_{10}$  and  $\text{PM}_{2.5}$



# Modelling the space-time interaction term, $\omega(\mathbf{s}, t)$

- We use an extended space-time model based on Gaussian Predictive Processes (GPP).
- We have added further flexibility into the model by improving the knot-selection process in the GPP method.
- The extension allowed us to have more knots in the densely populated areas leading to better estimation in those neighbourhoods.
- Details are omitted but all models are implemented by **extending** the R package `spTimer` publicly available from CRAN.

# Results for NO<sub>2</sub> and O<sub>3</sub> model validation

<b>NO<sub>2</sub>: Fitting N = 92,440, validation N=67,094, SD=37.19</b>					
<b>Model</b>	<b>RMSPE</b>	<b>MAPE</b>	<b>Bias</b>	<b>Coverage (%)</b>	<b>R<sup>2</sup></b>
Simple Kriging	32.87	22.88	2.56	69.59	0.53
Linear model	30.46	19.63	-5.09	94.43	0.60
Best model	17.65	12.99	0.41	97.42	0.89

<b>O<sub>3</sub>: Fitting N = 58,900, validation N=39,728, SD=22.23</b>					
Simple Kriging	13.30	9.86	-2.95	78.25	0.80
Linear model	16.0	12.42	8.47	93.86	0.69
Best model	10.17	7.59	0.07	91.72	0.89

**Table:** Assessment of predictive performance for a range of models for NO<sub>2</sub> and O<sub>3</sub>.  $R^2$  denotes the sample correlation coefficient between the predictions and actual observations. Our modelling roughly halves the error variability in out of sample validation!

# Results for PM<sub>10</sub> and PM<sub>2.5</sub> model validation

PM <sub>10</sub> : Fitting N = 46,894, validation N=32,333, SD=11.98					
Model	RMSPE	MAPE	Bias	Coverage (%)	$R^2$
Simple Kriging	7.34	4.75	-0.75	64.96	0.77
Linear model	9.98	6.74	-1.74	93.70	0.61
Best model	5.48	3.56	-0.65	90.03	0.81
PM <sub>2.5</sub> : Fitting SS = 35,791, validation SS=22,386, SD=9.52					
Model	RMSPE	MAPE	Bias	Coverage (%)	$R^2$
Simple Kriging	4.63	2.96	-0.72	67.84	0.81
Linear model	8.03	5.30	-1.87	92.73	0.60
Best model	4.30	2.66	-0.97	82.38	0.85

**Table:** Assessment of predictive performance for a range of models for PM<sub>10</sub> and PM<sub>2.5</sub>.  $R^2$  denotes the sample correlation coefficient between the predictions and actual observations. Our modelling roughly halves the error variability in out of sample validation!

# Summary of RMSEs for daily data for London only

- Difficult to compare our error rates with those from other articles as there is none modelling UK data for a 5-year period!
- Closest is Pirani, Gulliver, Fuller and Blangiardo (2014) who modelled daily  $PM_{10}$  data for London for 728 days in 2002-2003.

Model	RMSPE	MAPE	Bias	$R^2$	Cover (%)
PM <sub>10</sub> : Fitting N = 11,828, validation N=1,393					
Our best model	3.81	2.85	0.87	0.85	89.37
Pirani et al 2014	4.75	—	—	0.63	—

**Table:** Our model validation measures for  $PM_{10}$  using 24 fitting and 5 validation sites within London for years in 2007-2011. The second row quotes the results from Pirani et al for their 2002-2003 data.

# Aggregating to local authority area, $\mathcal{A}_k$

- We define average pollution:

$$v(\mathcal{A}_k, t) = \frac{1}{|\mathcal{A}_k|} \int_{\mathbf{s} \in \mathcal{A}_k} \mu^2(\mathbf{s}, t) d\mathbf{s}, \quad (1)$$

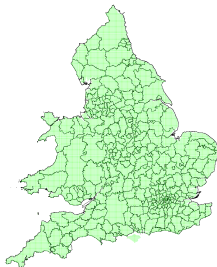
where  $\mu^2(\mathbf{s}, t)$  is the true unobserved concentration at location  $\mathbf{s}$  and at time  $t$ .

- We estimate it by the block average of  $N_k$  grid corners within the LA  $\mathcal{A}_k$ :

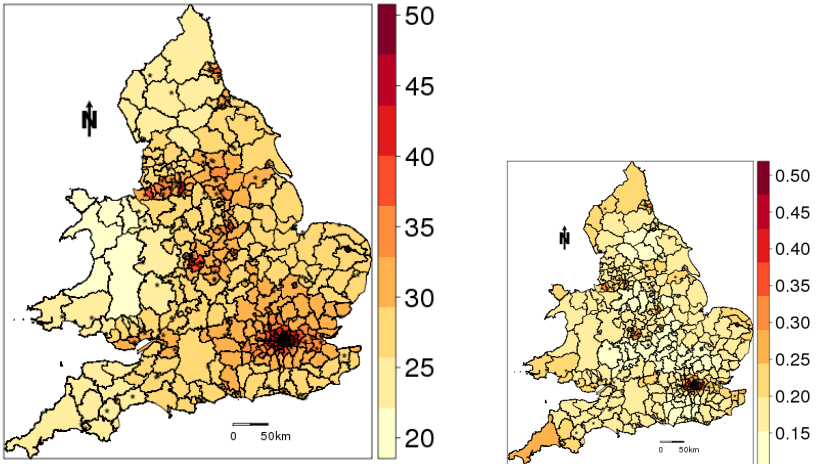
$$\hat{v}(\mathcal{A}_k, t) = \frac{1}{N_k} \sum_{j=1}^{N_k} \mu^2(\mathbf{s}_{kj}^*, t), \quad (2)$$

where  $\mathbf{s}_{kj}^*, j = 1, \dots, N_k$  is a grid of sites covering the areal unit  $\mathcal{A}_k$ .

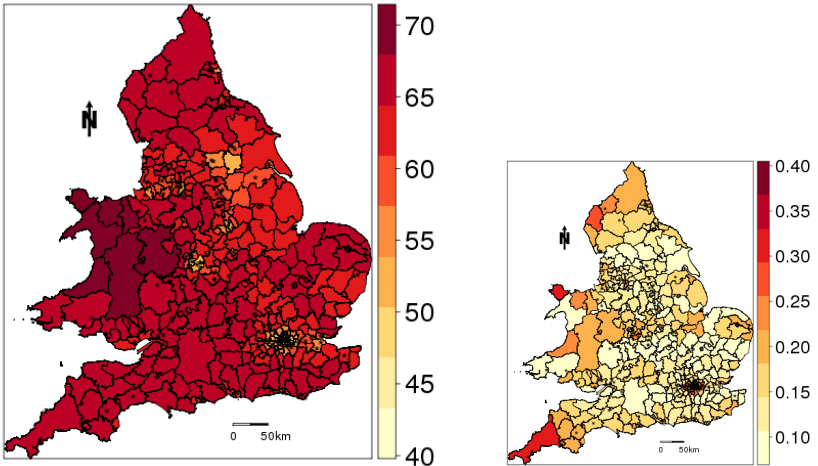
# Aggregating to LUAs...



- Map of 346 LUAs in England and Wales.
- A 1-kilometer square grid (151,248 **green dots**) is superimposed.
- Average air pollution in an LUA is the block average of the pollution levels in the **green dots** falling within that LUA.
- Our best Bayesian model is used to interpolate (model based Kriging) the air pollution at the **green dots**.
- Thus we produce air pollution estimate at any LUA at any time point (daily, monthly, annual)!

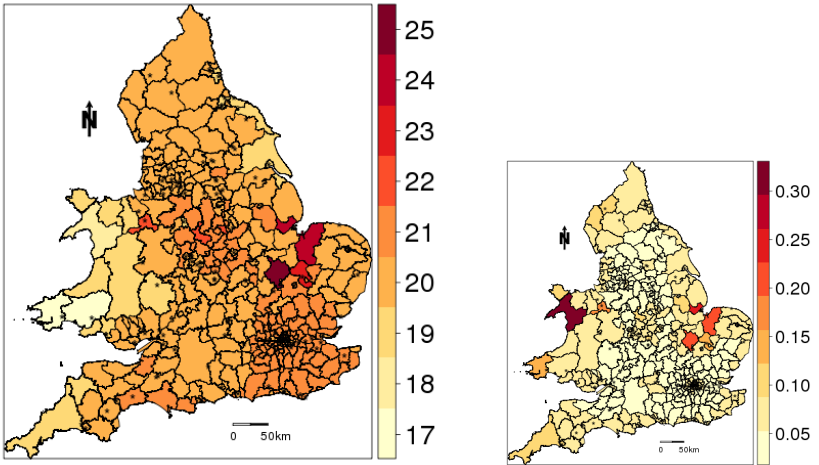


**Figure:** Local authority-wise annual prediction plot for NO<sub>2</sub> and their standard deviations (right panel) for 2011. Annual limit value of 40 is exceeded in most cities.

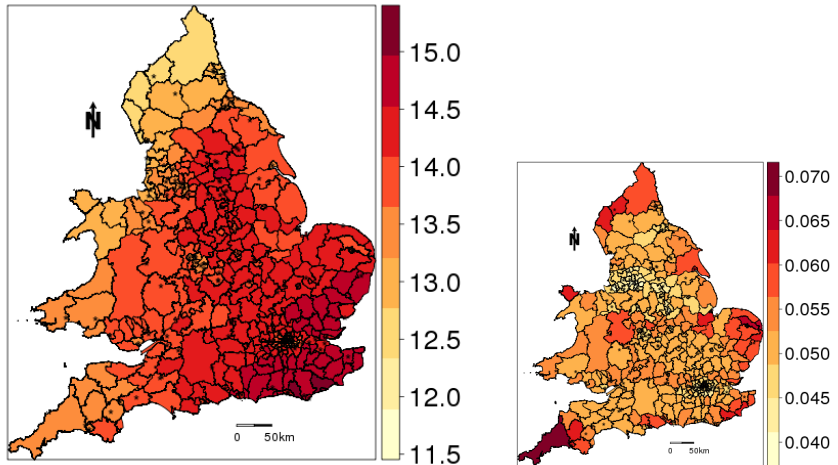


**Figure:** Local authority-wise annual prediction plot for O<sub>3</sub> and their standard deviations (right panel) for 2011. Rural areas have higher O<sub>3</sub> levels than urban areas.





**Figure:** Local authority-wise annual prediction plot for  $PM_{10}$  and their standard deviations for 2011. Cities and suburbs have higher levels.



**Figure:** Local authority-wise annual prediction plot for PM<sub>2.5</sub> and their standard deviations for 2011. Cities and suburbs, especially in the South-East, have higher levels.

# Estimating health effects (Lee et al 2016, Biostatistics)

- Let  $Y_{kt}$  denote the number of hospitalisation in the  $k$ th local authority  $\mathcal{A}_k$  in the  $t$ th month.
- $k = 1, \dots, 323$  local authorities in England
- $t = 1, \dots, 60$  months in five years, 2007-2011.

$$Y_{kt} \sim \text{Poisson}(E_{kt}R_{kt})$$

$$\log(R_{kt}) = \alpha + \beta_1 \hat{v}_{kt} + \beta_2 \text{j} \text{sa}_{kt} + \beta_3 \text{house}_{kt} + \psi_{kt}$$

- $E_{kt}$  is directly standardised hospitalisation (age and sex) counts nationally.
- $R_{kt}$ : Relative risk,
- $\hat{v}_{kt}$ : pollution estimate from our model.
- $\text{j} \text{sa}_{kt}$ : Average job seekers allowance: Confounder
- $\text{house}_{kt}$ : Average house price: Confounder
- $\psi_{kt}$ : space-time random effect for which we had elaborate models.

# Results from the health outcome model

	RR	Lower 2.5%	Upper 97.5%	Pollutant SD
NO <sub>2</sub>	1.028	1.021	1.033	16.07
PM <sub>10</sub>	1.026	1.011	1.039	4.90
PM <sub>2.5</sub>	1.006	0.993	1.020	4.11
O <sub>3</sub>	0.997	0.994	0.999	7.30

**Table:** Estimated health effects from each of the four pollutants. All results are presented as relative risks for a one standard deviation increase in pollution.

- An estimated 2.8% increased risk of hospitalisation due to one sd increase in exposure to NO<sub>2</sub>.
- Implies 17,000 extra hospital admissions per year, as there are around 613,000 admissions per year in England.
- This implies a potential annual spending of **£4.76 million** assuming a week's hospital stay on average for each patient.

# Conclusions

- 1 We have developed pollutant specific models which worked well for **all four** important pollutants,  $PM_{10}$ ,  $PM_{2.5}$ ,  $O_3$ ,  $NO_2$ .
- 2 Our models fill up the sparsity of the observed air quality data by integrating output from the **AQUM** which are available over a fine grid.
- 3 We are able to estimate pollution levels, along with their uncertainties, at any desired level of administrative geography.
- 4 We are able to **measure long term exposure** since we have modelled daily data for a 5 year period, for all four pollutants.
- 5 We are not aware of any similar study offering high quality air pollution estimates along with their **individual error bars**.

- ➊ Exposure estimates, and their uncertainties, from our best model:
  - ➊ for all four pollutants
  - ➋ at both daily and annual time scales
  - ➌ for the five years 2007-2011
  - ➍ at the 151,248 1-kilometer grid points
  - ➎ and also for all the local authorities in England and Wales
- ➋ are available online. Total size is about 64GB.
- ➌ From my website <http://www.soton.ac.uk/~sks/>.
- ➍ Thus we provide the most accurate empirically verified air pollution estimates at 1-kilometer grid in England and Wales for the five years 2007-2011.

## Possibilities are endless!

- Government and regulatory bodies can use the data to evaluate post-hoc compliance to air pollution standards in even un-monitored areas all over England and Wales.
- Compliance can be evaluated at any socio-economic-politico geographic scale: i.e. post-code, local authority area, LSOA, electoral wards etc.
- Researchers from both academic and government agencies such as the [Public Health England](#) can link air pollution to a range of health out-come data.
- For example, colleagues in UCL are associating air pollution levels with the millenium cohort data on children's mental health.

## For example

- Improve the models by further methodological development, e.g. multivariate models for the four pollutants.
  - Obtain similar exposure estimates for 2012-2017 possibly using new and improved models.
  - Develop on-line tools (apps) to deliver data sets on the fly!
    - for user defined geographies and coarser time domains, e.g. monthly, quarterly etc.
  - Evaluate health impact using rigorous epidemiological studies.
- 
- Please email me ([S.K.Sahu@soton.ac.uk](mailto:S.K.Sahu@soton.ac.uk)) if you have queries about the data sets.