# Global modelling of air pollution using multiple data sources

Matthew Thomas – M.L.Thomas@bath.ac.uk

Supervised by Dr. Gavin Shaddick
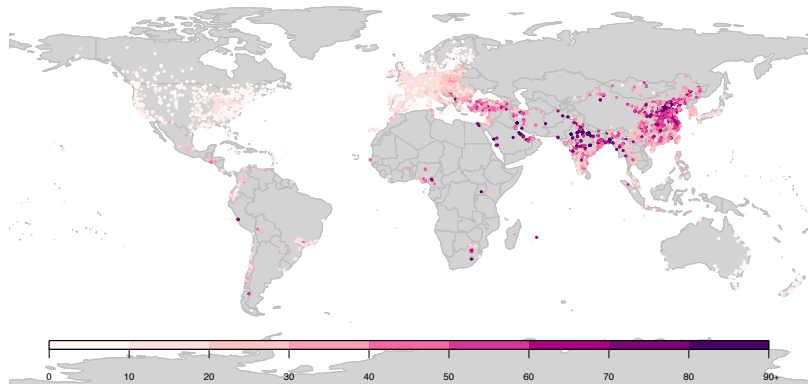In collaboration with IHME and WHO

July 17, 2018

# MOTIVATION

- ▶ Air pollution is an important determinant of health and poses a significant threat globally.
- ▶ It is known to trigger cardiovascular and respiratory diseases in addition to some cancers.
- ▶ Particulate Matter ($PM_{2.5}$) is estimated to be
  - ▶ 4th highest health risk factor in the world
  - ▶ attributable to 5.5 million premature deaths
- ▶ There is convincing evidence for the need to model air pollution effectively.

## REQUIREMENTS

- ▶ WHO and other partners plan to strengthen air pollution monitoring globally.
- ▶ This will produce accurate and convincing evidence of risks posed.

# GROUND MONITORING



Figure: World map with ground monitor locations, coloured by the estimated level of PM$_{2.5}$ in $\mu g m^{-3}$.

# REQUIREMENTS

- ▶ WHO and other partners plan to strengthen air pollution monitoring globally.
- ▶ This will produce accurate and convincing evidence of risks posed.
- ▶ Allow data integration from different sources.
- ▶ This will allow borrowing from each methods respective strengths.
- ▶ Currently, three methods are considered:
  - ▶ Ground Monitoring,
  - ▶ Satellite Remote Sensing and
  - ▶ Atmospheric Modelling
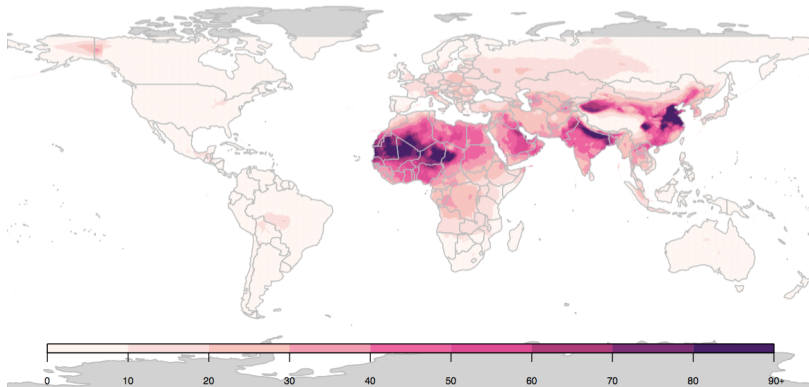
# SATELLITE REMOTE SENSING



Figure: Global satellite remote sensing estimates of $PM_{2.5}$ in $\mu g m^{-3}$ for 2014 used in GBD2015
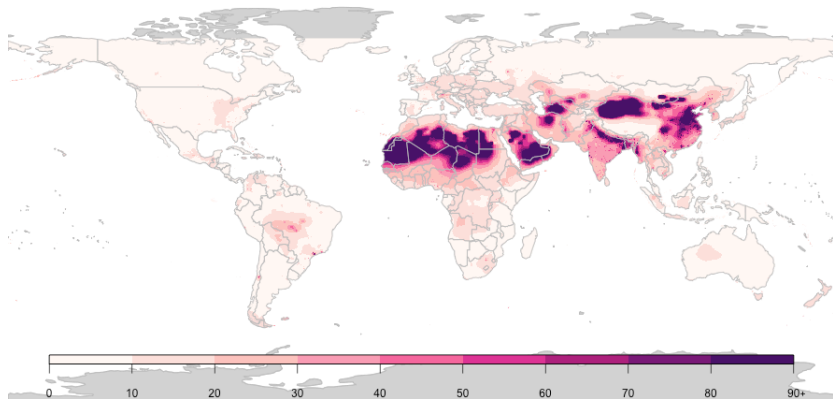
# ATMOSPHERIC MODELLING



Figure: Global chemical transport model estimates of $PM_{2.5}$ in $\mu gm^{-3}$ for 2014 used in GBD2015
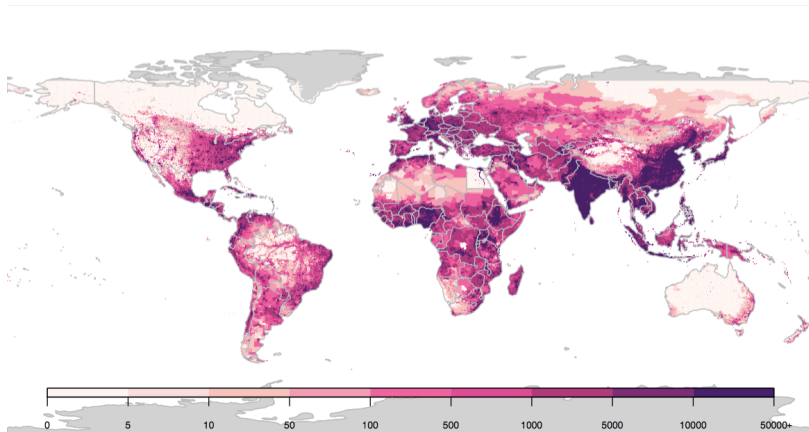
# POPULATION ESTIMATES



Figure: Estimate of population density per $0.1^o \times 0.1^o$ grid location for 2014 used in GBD2015

# PREVIOUS APPROACH

▶ The current GBD approach to modelling combines estimates from atmospheric models and satellites into a 'fused' estimate.

▶ Let $x_i^{am}$ and $x_i^{sat}$ be atmospheric model and satellite estimates for grid cell $i$, then the fused estimate is defined as:

$$x_i^{fus} = \frac{x_i^{sat} + x_i^{am}}{2}.$$

▶ The ground monitors and grid data are calibrated, logged and fused data is used as an explanatory variable in a linear model to determine ground level $PM_{2.5}$:

$$\log\left(y_i^{gm}\right) = \beta_0 + \beta_1 \log\left(x_i^{fus}\right) + \epsilon_i \quad i = 1, \ldots, n.$$

▶ Ground level $PM_{2.5}$ is then estimated using tradition linear modelling techniques.

# A MULTILEVEL RANDOM EFFECT MODEL

▶ Suppose that a ground monitor at location $s$ is situated in grid cell $B_j$.

▶ To avoid non-negativity and skew we consider the estimates of $PM_{2.5}$ on the log-scale

▶ We then assume that the log estimates of $PM_{2.5}$ from ground monitors, $y_s$ are normally distributed

$$y_s = z_{B_i} + \sum_{j=1}^{n} \gamma_j x_{s,j} + \epsilon_s$$

where

   ▶ $x_{s,j}$ are covariate information for ground monitor at location $s$,
   ▶ $z_{B_i}$ is a mean trend for grid cell $B_i$
   ▶ $\epsilon_s \sim N(0, \sigma_\epsilon^2)$ is measurement error.

# A MULTILEVEL RANDOM EFFECT MODEL

▶ The mean trend $z_{B_i}$ for grid cell $B_i$ is modelled using the following,

$$z_{B_i} = \tilde{\beta}_0 + \sum_{j=1}^{k} \tilde{\beta}_j u_{B_i,j} + \sum_{j=k+1}^{m} \beta_j u_{B_i,j} + e_{B_i}$$

where

  ▶ $u_{B_i,j}$ are covariate information for grid cell $B_i$,
  ▶ $e_s \sim N(0, \sigma_e^2)$ is the within cell variability.

# A MULTILEVEL RANDOM EFFECT MODEL

▶ To allow for more local variation we allow a series random effects

$$\tilde{\beta}_j = \beta_{j0} + \sum_{k=1}^{K} \beta_{jk}$$

▶ We propose these random effects to have a nested hierarchy to allow borrowing between levels.

▶ Here we aggregate countries into regions and regions into superregions

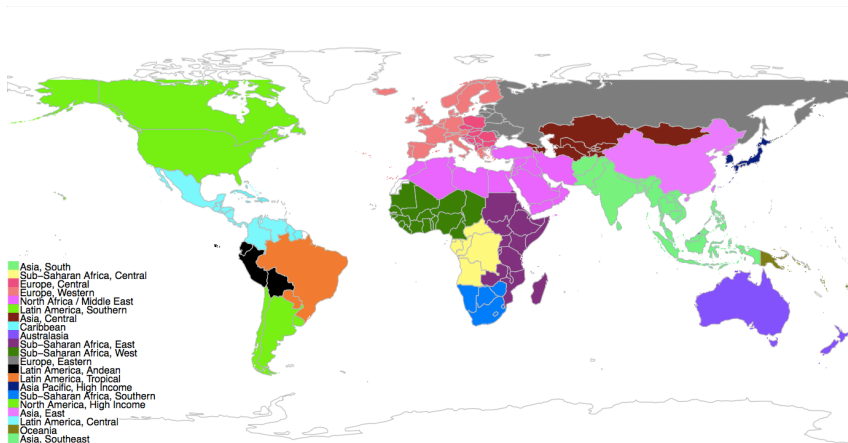  ▶ Using country level mortality levels and causes of death

# DEFINED REGIONS



Asia, South
Sub-Saharan Africa, Central
Europe, Central
Europe, Western
North Africa / Middle East
Latin America, Southern
Asia, Central
Caribbean
Australasia
Sub-Saharan Africa, East
Sub-Saharan Africa, West
Europe, Eastern
Latin America, Andean
Latin America, Tropical
Asia Pacific, High Income
Sub-Saharan Africa, Southern
North America, High Income
Asia, East
Latin America, Central
Oceania
Asia, Southeast

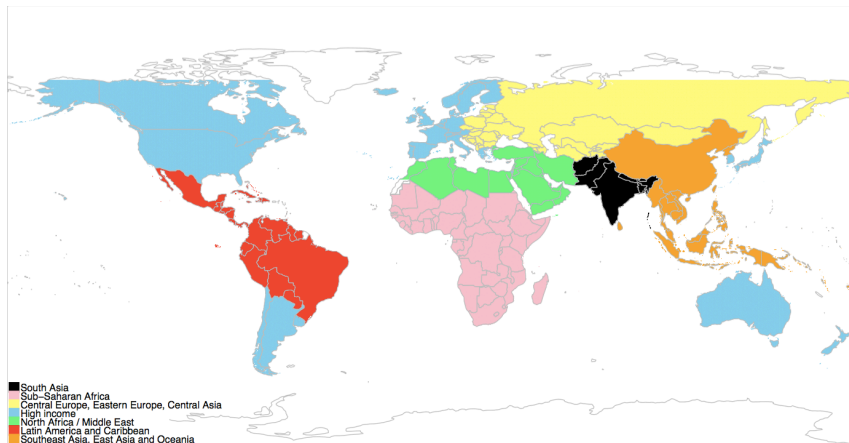Figure: World map coloured by GBD defined Regions

# DEFINED SUPER REGIONS



Figure: World map coloured by GBD defined Super Regions

## COMPUTATION

▶ Bayesian models of this complexity do not have analytical solutions.

▶ 'Big' data means traditional MCMC techniques are impractical.

▶ Recent advances in approximate Bayesian inference provide fast and efficient methods for modelling, such as Integrated Nested Laplace Approximations (INLA).

▶ INLA performs numerical calculations of posterior densities using Laplace Approximations hierarchical latent Gaussian models:

$$p(\theta_k|\boldsymbol{y}) = \int p(\boldsymbol{\theta}|\boldsymbol{y})d\boldsymbol{\theta}_{-k} \quad p(z_j|\boldsymbol{y}) = \int p(z_j|\boldsymbol{\theta},\boldsymbol{y})p(\boldsymbol{\theta}|\boldsymbol{y})d\boldsymbol{\theta}$$

▶ A latent Gaussian process allows for sparse matrices, and therefore efficient computation.

# COMPUTATION

- ▶ Already suite of programs to implement these (R-INLA).
- ▶ However, while INLA is computationally more attractive, R-INLA still requires huge computation and memory usage.
- ▶ Unable to run this model on standard computers (4-8GB RAM).
- ▶ Required the use of a High-Performance Computing (HPC) service.
    - ▶ Balena cluster at University of Bath.
    - ▶ $2 \times 512$GB RAM nodes ($32 \times 32$GB RAM cores).
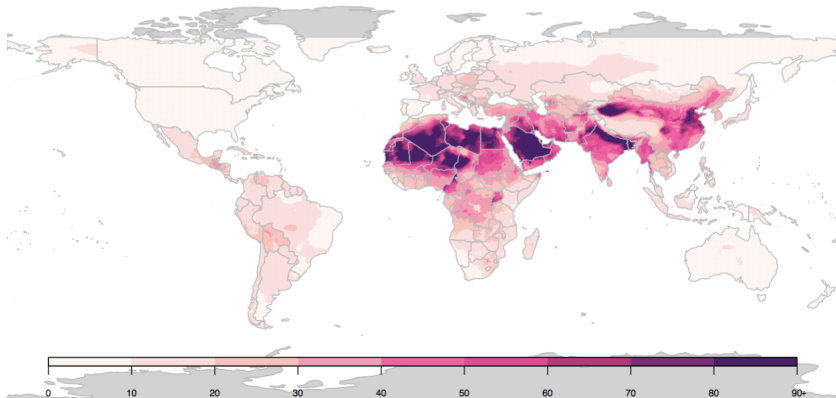- ▶ Took an iterative approach to prediction.

# PREDICTIONS



Figure: Predictions of $PM_{2.5}$ in $\mu g m^{-3}$, from hierarchical model for 2014.
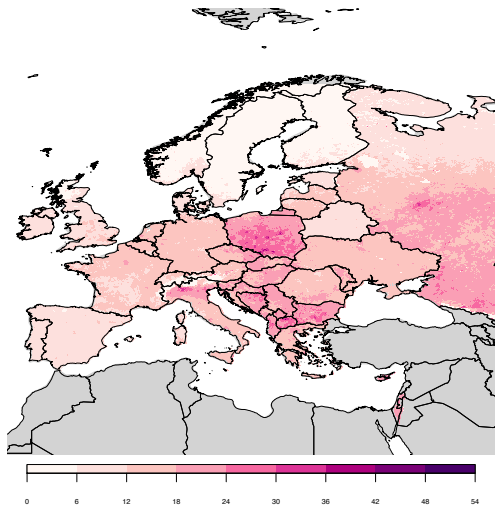
# PREDICTIONS: REGIONAL



Figure: Predictions of PM$_{2.5}$ in $\mu g m^{-3}$, from hierarchical model for 2014 in Europe
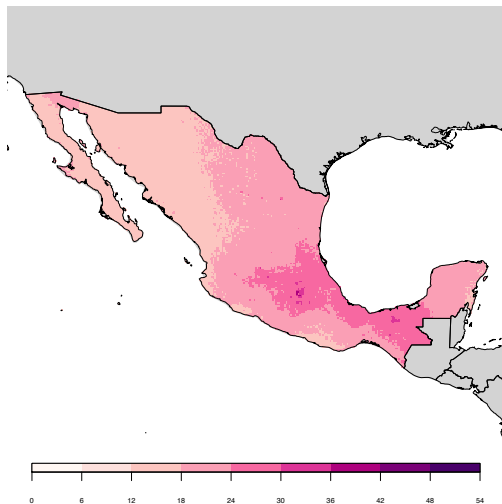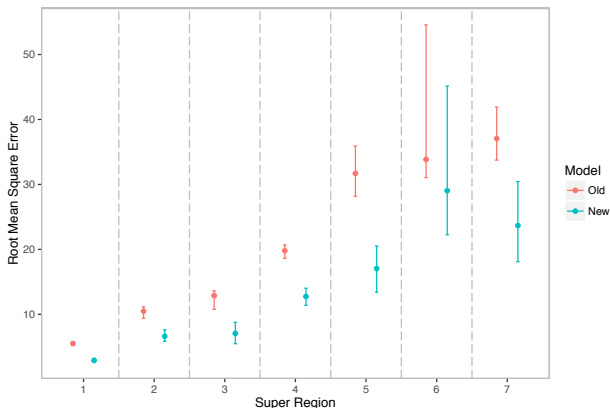
# PREDICTIONS: LOCAL



Figure: Predictions of PM$_{2.5}$ in $\mu gm^{-3}$, from hierarchical model for 2014 in Mexico

# EVALUATION: CROSSVALIDATION



Figure: Comparison of RMSE between approaches. Dots denote the median of the distribution from 25 training/evaluation sets and the vertical lines the range of values. Super-regions are 1: High income, 2: Central Europe, Eastern Europe and Central Asia, 3: Latin America and Caribbean, 4: Southeast Asia, East Asia and Oceania, 5: North Africa / Middle East, 6: Sub-Saharan Africa and 7: South Asia.

# BAYESIAN MELDING

- Bayesian melding assumes there is one latent process $z_s$ that drives all sources of data.
- **Data Level:** Ground monitor data is assumed to be a measurement error model i.e.

$$y_s^{gm} = z_s + \epsilon_s \quad \epsilon_s \sim N(0, \sigma_\epsilon^2)$$

- The grid data is then modelled at point locations as a function of the true underlying process

$$y_s^{grid} = f(z_s) + \delta_s \quad \delta_s \sim N(0, \sigma_\delta^2).$$

- As we cannot model grid data with a point process, we integrate and get the following integral:

$$y_{B_j}^{grid} = \int_{B_j} f(z_s) + \delta_s d\mathbf{s}, j = 1, 2, \ldots, m$$

# BAYESIAN MELDING

▶ **Latent Process Level:** In the second stage of the model, the true underlying process $z_s$ is assumed to follow the model

$$z_s = \mu_s + m_s$$

where $\mu_s$ is a spatial trend and the $m_s$ is a spatial process for location $s$.

▶ **Inference:** It will be quantify the true levels of $PM_{2.5}$

$$p(z_s | \mathbf{y}^{gm}, \mathbf{y}^{grid}) = \int p(z_s | \mathbf{y}^{gm}, \mathbf{y}^{grid}, \boldsymbol{\theta}) p(\boldsymbol{\theta} | z_s) d\boldsymbol{\theta}$$

# BAYESIAN MELDING

- ▶ Makes use of a flexible and coherent framework
- ▶ Allows user to assume one underlying process driving the
- ▶ Treats estimation methods as different quantities but are intrinsically linked
- ▶ To implement this framework on large-scale problems!
- ▶ Look at approximate Bayesian inference (INLA) for more efficient computation
- ▶ Allow for time effects.

# Thank you for listening!

# ANY QUESTIONS?