



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

February 22, 2016

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 East Asia
 - ▶ 6th in South Asia and
 - ▶ 7th in Africa and the Middle East
- ▶ There is convincing evidence for the need to model air pollution effectively.

MOTIVATION

- ▶ WHO and other partners plan to strengthen air pollution monitoring globally.
- ▶ Aim is to 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

GROUND MONITORING

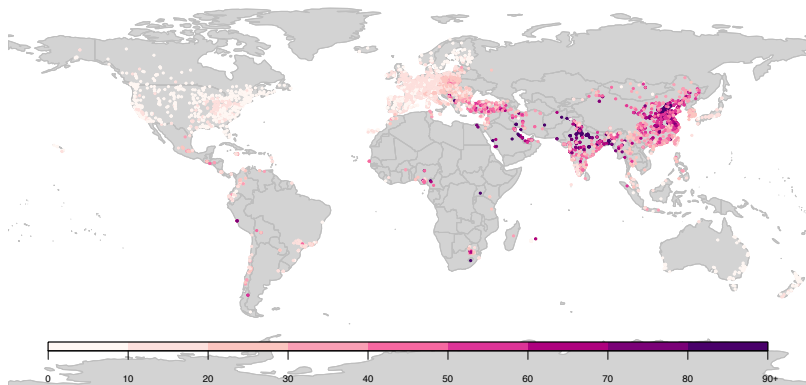


Figure: World map with ground monitor locations, coloured by the estimated level of PM_{2.5} in μgm^{-3} .

SATELLITE REMOTE SENSING

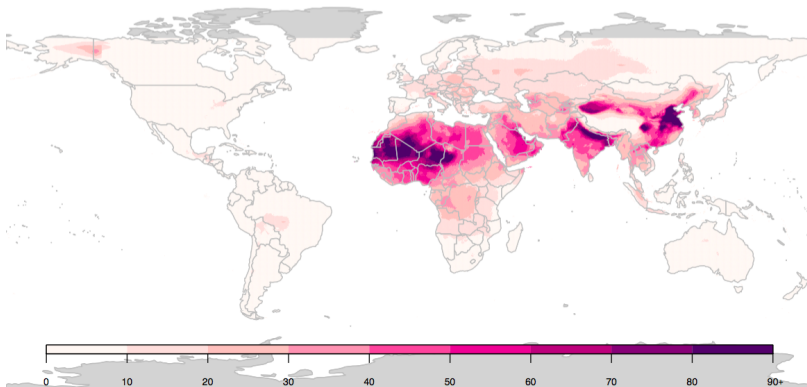


Figure: Global satellite remote sensing estimates of PM_{2.5} in $\mu\text{g m}^{-3}$ for 2014 used in GBD2015

ATMOSPHERIC MODELLING

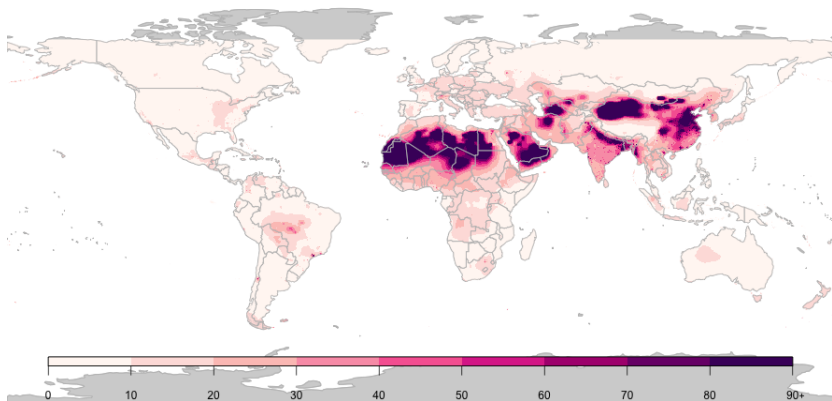


Figure: Global chemical transport model estimates of PM_{2.5} in $\mu\text{g m}^{-3}$ for 2014 used in GBD2015

POPULATION ESTIMATES

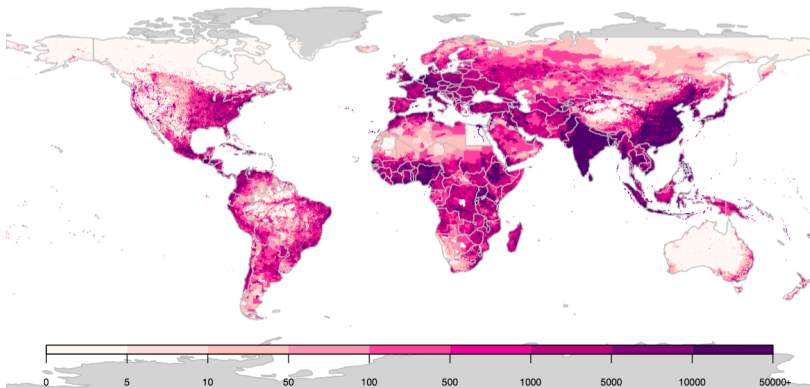


Figure: Estimate of population density per $0.1^\circ \times 0.1^\circ$ grid location for 2014 used in GBD2015

LINEAR MODELLING

- ▶ 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 $PM_{2.5}$:

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

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

PREDICTIONS

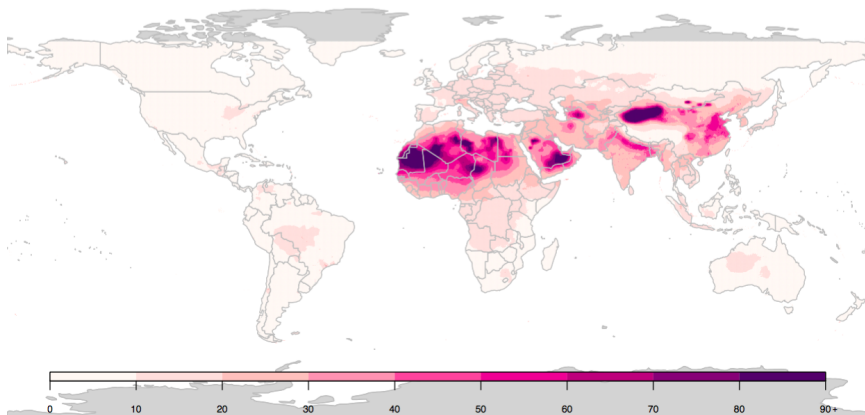


Figure: Predictions of PM_{2.5} in $\mu\text{g m}^{-3}$ for 2014, from existing WHO/GBD model.

PREDICTIONS, BY REGION

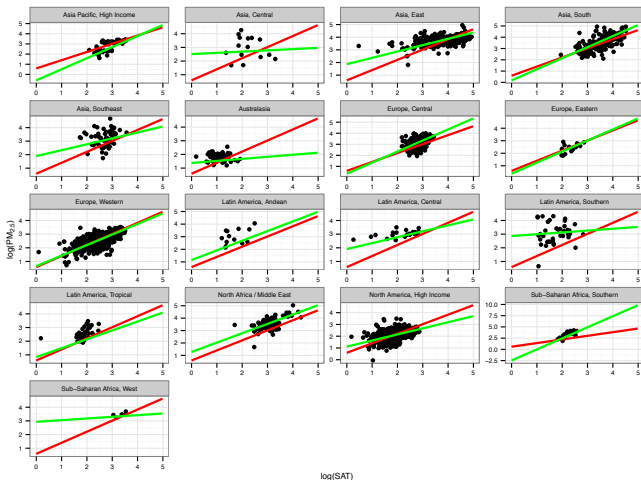


Figure: $\text{PM}_{2.5}$ measurements against satellite estimates on the log-scale, for 2014, split by region. The red and green lines denote the single 'global' and a region specific model respectively, estimated using all of the data.

HIERARCHICAL MODELLING

- ▶ **Observation Level:** We assume the ground monitor data y_{ijkl}^{gm} comes from a measurement error model, on the log-scale:

$$\log\left(y_{ijkl}^{(gm)}\right) = z_{ijkl}^{(gm)} + \epsilon_{ijkl} \quad \epsilon_{ijkl} \sim N(0, \sigma_{\epsilon}^2)$$

- ▶ **Process Level:** Let x_{ijkl}^{sat} , x_{ijkl}^{am} and x_{ijkl}^{pop} denote the satellite, atmospheric model and population estimates respectively. The underlying process is modelled as follows:

$$z_{ijkl} = \tilde{\beta}_{0jkl} + \tilde{\beta}_{1jkl} \log(x_{ijkl}^{sat}) + \tilde{\beta}_{3jkl} \log(x_{ijkl}^{pop}) + \beta_4 x_i^{elev} + \beta_5 x_i^{dust} + \beta_6 x_i^{sanoc}$$

$$\tilde{\beta}_{mjkl} = \beta_m^G + \beta_{mj}^{SR} + \beta_{mjk}^R + \beta_{mjkl}^C + \beta_m^P P_i + \beta_m^A A_i + \beta_m^U U_i, \quad (m = 0, 1, 2, 3)$$

- ▶ **Prior Level:** Vague priors were used, default in R-INLA to exploit conjugacy and therefore allow efficient computation.

BAYESIAN HIERARCHICAL MODELLING

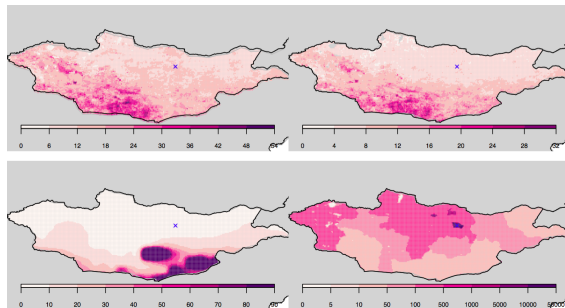
- ▶ Many spatial or spatio-temporal models that involve data inherently have a hierarchical structure.
- ▶ Hierarchical models extremely useful and flexible framework in which to model complex relationships and dependencies in data.
- ▶ Bayesian hierarchical models are commonly written:
 1. The observation level $\mathbf{y}|z, \theta$ - Data \mathbf{y} , are assumed to arise from an underlying latent process z , which is unobservable but measurements with error can be taken.
 2. The underlying process level $z|\theta$ - The latent process z assumed to drive the observable data and is the true underlying quantity of interest.
 3. The prior level θ - This level describes known prior information about the model parameters θ
- ▶ Bayesian techniques to statistical modelling allow us to interpret levels in the model that weren't measured such as the underlying latent process.

APPROACH TO DATA INTEGRATION

- ▶ Data integration in the current framework uses a fused estimate.
- ▶ Atmospheric model estimates are numerically simulated data from a specified PDE.
- ▶ Satellite estimates are modelled from images.
- ▶ Both estimation methods are very different; as they should provide different perspectives on the modelled system and have very different error structures.
- ▶ So, initially terms were fitted separately within the model.
- ▶ However, both data sources were highly collinear and numerical estimates were removed from the model.

ADDITION OF EXTRA COVARIATES

- ▶ In many areas of the world air pollution estimates weren't very accurate.
- ▶ Example: Ulan Bator, Mongolia



- ▶ Other pollutant levels were not available
- ▶ Population was added into the model (on the log-scale)
- ▶ Other covariates were added into the model

RANDOM EFFECTS

- ▶ Linear models used by WHO, assume a single global relationship.
- ▶ This is a massive assumption, that is unlikely to hold.
- ▶ Each country is assigned to a 'Region' and 'Super Region' (Nested Hierarchy).

DEFINED REGIONS

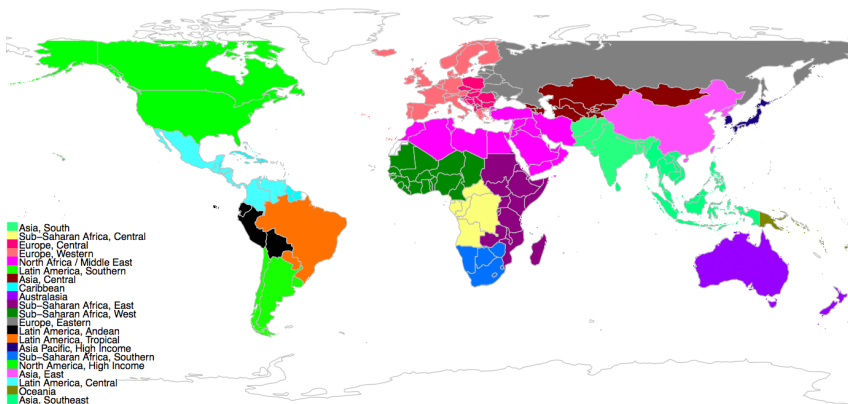


Figure: World map coloured by GBD defined Regions

DEFINED SUPER REGIONS

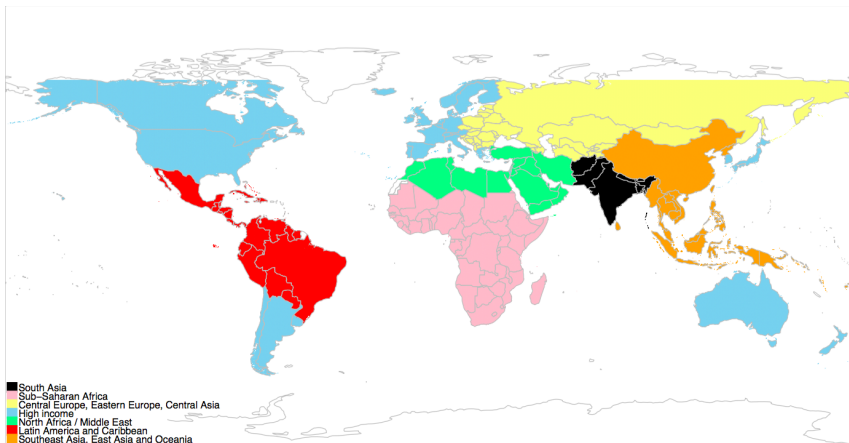


Figure: World map coloured by GBD defined Super Regions

RANDOM EFFECTS

- ▶ Linear models used by WHO, assume a single global relationship.
- ▶ This is a massive assumption, that is unlikely to hold.
- ▶ Each country is assigned to a 'Region' and 'Super Region' (Nested Hierarchy).
- ▶ Could like earlier fit models by Super Region, Region or Country to look at more local relationships. However this comes with issues.
- ▶ Instead we added IID random effects for Super Region, Region and Country on Satellite and Population
- ▶ This allows borrowing from hierarchy when there is limited data.
- ▶ It proved useful for further borrowing at lower levels in the hierarchy by using a Conditional Autoregressive model (CAR) on the population coefficient.

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:

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

- ▶ 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.
 - ▶ $2 \times 512\text{GB}$ RAM nodes ($32 \times 32\text{GB}$ RAM cores).
- ▶ Unable to use INLA as parallelised code.
- ▶ Restricted to $1 \times 32\text{GB}$ RAM node.
- ▶ Took an iterative approach to prediction.

PREDICTIONS

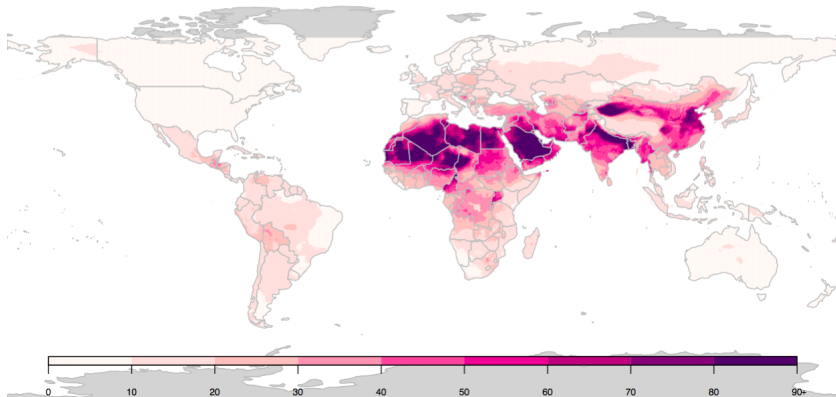


Figure: Predictions of PM_{2.5} in $\mu\text{g m}^{-3}$, from hierarchical model for 2014.

PREDICTIONS: REGIONAL

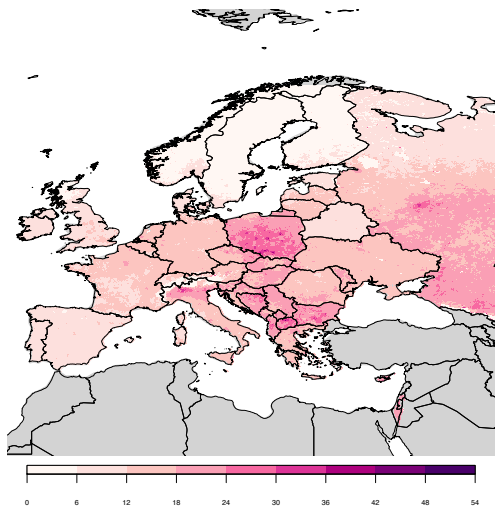


Figure: Predictions of PM_{2.5} in $\mu\text{g m}^{-3}$, from hierarchical model for 2014 in Europe

PREDICTIONS: LOCAL

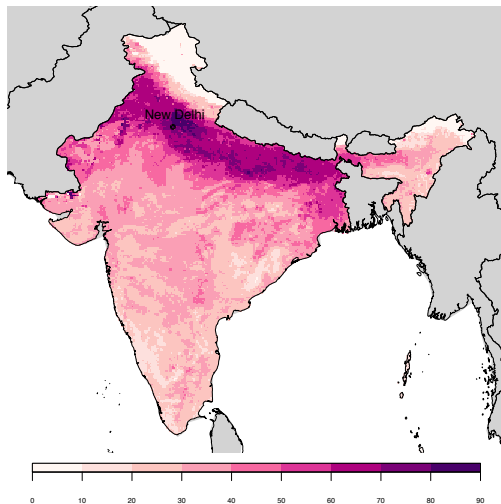


Figure: Predictions of PM_{2.5} in $\mu\text{g m}^{-3}$, from hierarchical model for 2014 in India

PREDICTIONS: LOCAL

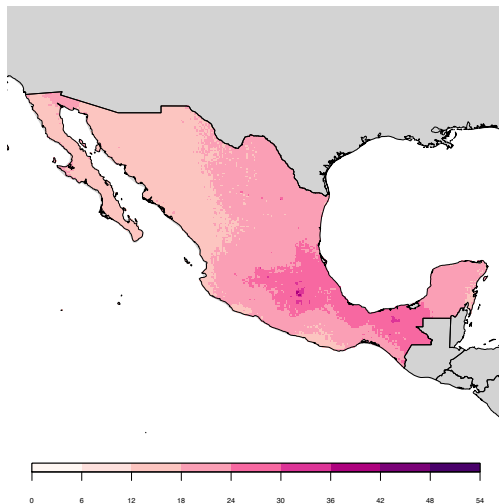


Figure: Predictions of PM_{2.5} in $\mu\text{g m}^{-3}$, from hierarchical model for 2014 in Mexico

UNCERTAINTY

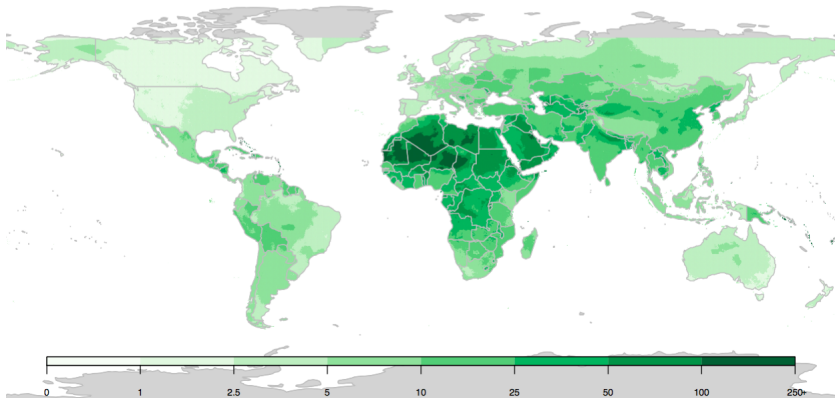


Figure: Uncertainty of PM2.5 predictions for 2010, for hierarchical model; half length of estimated 95% credible intervals

EXCEEDANCE PROBABILITIES

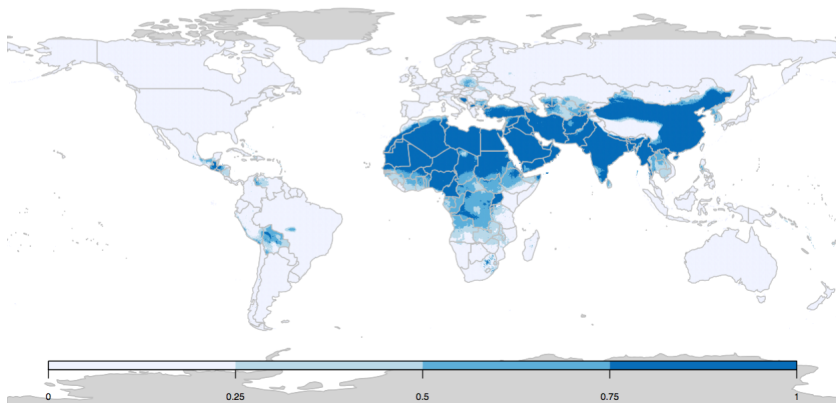


Figure: Probability that level of PM_{2.5} in each cell exceeds 25 μg m⁻³ in 2010, for hierarchical model

BAYESIAN MELDING

- ▶ Bayesian melding makes use of a Bayesian hierarchical model.
- ▶ Assumes a latent process $z(\mathbf{s})$ that represents the true level PM_{2.5}.
- ▶ **Data Level:** Ground monitor data is assumed to be a measurement error model i.e.

$$y^{gm}(\mathbf{s}) = z(\mathbf{s}) + \epsilon(\mathbf{s}) \quad \epsilon(\mathbf{s}) \sim N(0, \sigma_\epsilon^2)$$

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

$$y^{grid}(\mathbf{s}) = f(z(\mathbf{s})) + \delta(\mathbf{s}) \quad \delta(\mathbf{s}) \sim N(0, \sigma_\delta^2).$$

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

$$y^{grid}(B_j) = \int_{B_j} f(z(\mathbf{s})) + \delta(\mathbf{s}) d\mathbf{s}, j = 1, 2, \dots, m$$

BAYESIAN MELDING

- ▶ **Latent Process Level:** In the second stage of the model, the true underlying process $z(\mathbf{s})$ is assumed to follow the model

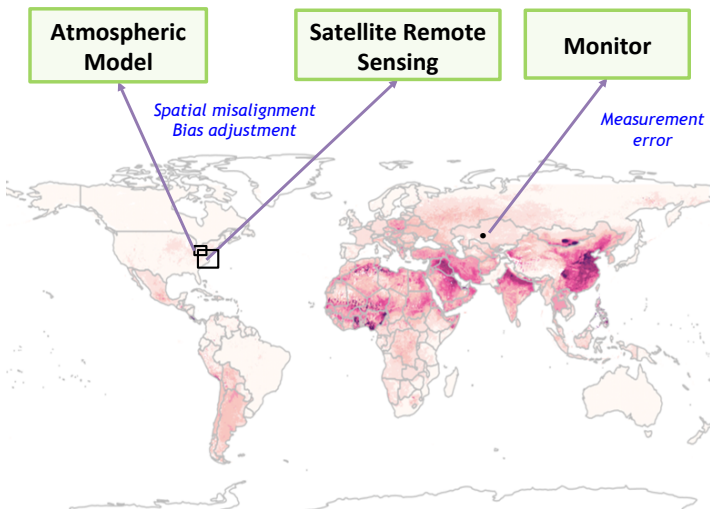
$$z(\mathbf{s}) = \mu(\mathbf{s}) + m(\mathbf{s})$$

where $\mu(\mathbf{s})$ is a spatial trend and the $m(\mathbf{s})$ is a spatial process for location \mathbf{s} .

- ▶ **Prior Level:** It will also be necessary to specify relevant priors for model parameters.
- ▶ **Inference:** It will be quantify the true levels of $\text{PM}_{2.5}$

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

BAYESIAN MELDING



BAYESIAN MELDING

Advantages:

- ▶ 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

Disadvantages:

- ▶ Very computationally demanding (particularly with MCMC)
- ▶ Only implemented in small-scale problems (~20 Monitors)

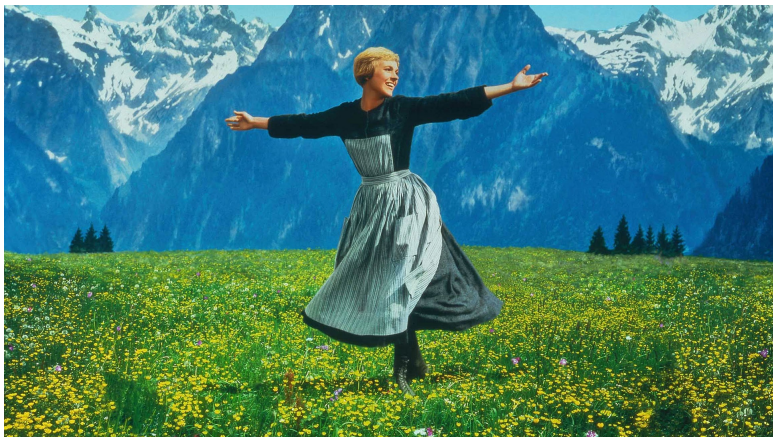
Aims:

- ▶ To implement this framework on large-scale problems!
- ▶ Look at approximate Bayesian inference (INLA) for more efficient computation
- ▶ Allow for time effects.

LIVE A LITTLE LESS LIKE THIS....



... AND MORE LIKE THIS...



ANY QUESTIONS?



Thank you for listening!