

A Data Fusion Approach for Space-Time Analysis of Speciated PM_{2.5}

Colin Rundel

Duke University

August 3, 2014

Fine particulate matter ($\text{PM}_{2.5}$) is an EPA regulated air pollutant linked to a variety of adverse health effects

- Classified based on particle size ($< 2.5 \mu\text{m}$ diameter)
- Major species: Sulfate, Nitrate, Ammonium, Soil, Carbon.
- Minor species: trace elements (K, Mg, Ca), heavy metals (Cu, Fe), etc.
- Complex spatio-temporal dependence between species

Data (2007)

Speciated PM_{2.5} Sources

- Chemical Speciation Network (CSN) - 221 stations
- Interagency Monitoring of Protected Visual Environments (IMPROVE) - 172 stations

Total PM_{2.5} Sources

- Federal Reference Method (FRM) - 949 stations

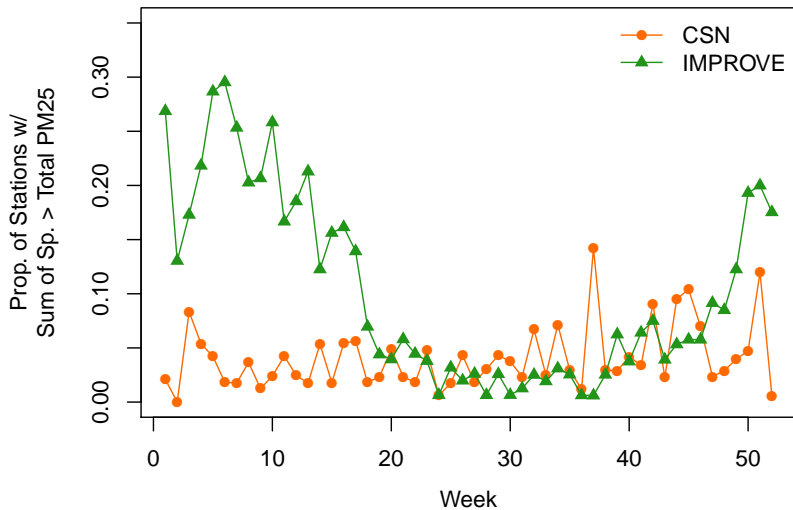
Model Output

- Community Multi-scale Air Quality (CMAQ) - 12 km grid

Data Issues

- Monitoring frequency
- Total vs Sum of Species

Total PM_{2.5} vs Sum of Species



Species Model Details

For the 5 major species (Sulfate, Nitrate, Ammonium, Soil, Carbon) and the two networks (CSN, IMPROVE):

$$C_t^i(\mathbf{s}) = Z_t^i(\mathbf{s}) + \epsilon_{C,t}^i(\mathbf{s})$$

$$I_t^i(\mathbf{s}) = Z_t^i(\mathbf{s}) + \epsilon_{I,t}^i(\mathbf{s})$$

where $Z_t^i(\mathbf{s})$ are the latent “true” concentrations of species i at time t and locations \mathbf{s} ,

$$Z_t^i(\mathbf{s}) = \max \left(0, \tilde{Z}_t^i(\mathbf{s}) \right)$$

$$\tilde{Z}_t^i(\mathbf{s}) = \beta_{0,t}^i + \beta_{0,t}^i(\mathbf{s}) + \beta_{1,t}^i Q_t^i(B_{\mathbf{s}})$$

Total PM_{2.5} Model Details

For total PM_{2.5} from the three networks (CSN, IMPROVE, FRM):

$$C_t^{tot}(\mathbf{s}) = Z_t^{tot}(\mathbf{s}) + \epsilon_{C,t}^{tot}(\mathbf{s})$$

$$I_t^{tot}(\mathbf{s}) = Z_t^{tot}(\mathbf{s}) + \epsilon_{I,t}^{tot}(\mathbf{s})$$

$$F_t^{tot}(\mathbf{s}) = Z_t^{tot}(\mathbf{s}) + \epsilon_{F,t}^{tot}(\mathbf{s})$$

where $Z_t^{tot}(\mathbf{s})$ are the latent “true” concentration of total PM_{2.5} at time t and locations \mathbf{s} , which is given by the sum of the major species and the “other” species concentrations.

$$Z_t^{tot}(\mathbf{s}) = \sum_{i=1}^5 Z_t^i(\mathbf{s}) + Z_t^o(\mathbf{s})$$

$$Z_t^o(\mathbf{s}) = \max \left(0, \tilde{Z}_t^o(\mathbf{s}) \right) \quad \tilde{Z}_t^o(\mathbf{s}) = \beta_{0,t}^o + \beta_{0,t}^o(\mathbf{s}) + \beta_{1,t}^o Q_t^o(B_s)$$

Spatial Dependence

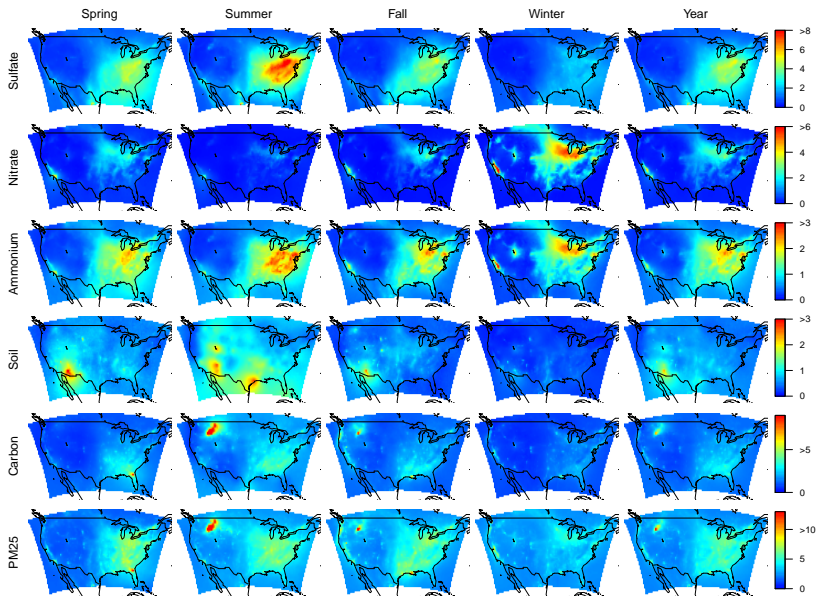
Spatial dependence enters the model through the $\beta_{0,t}^i(\mathbf{s})$ parameters for $i \in \{0, 1, 2, 3, 4, 5\}$.

$$\beta_{0,t}^i(\mathbf{s}) = \sigma_t^i w_t^i(\mathbf{s})$$

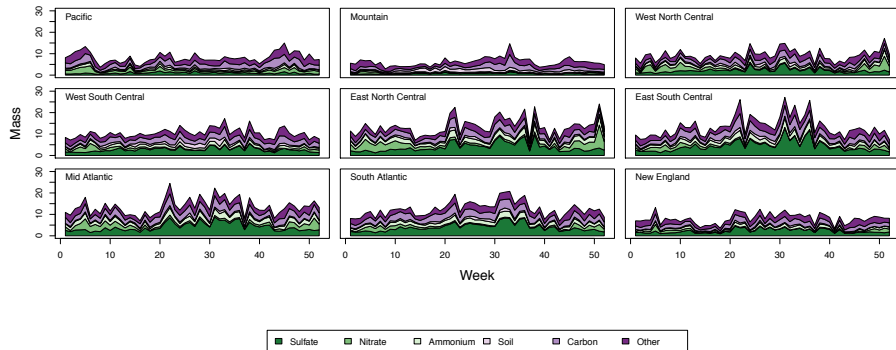
where $w_t^i(\mathbf{s})$ are zero mean, variance 1, Gaussian processes with exponential correlation given by

$$\text{corr}(w_t^i(\mathbf{s}), w_t^i(\mathbf{s}')) = \exp(-\phi_t^i |\mathbf{s} - \mathbf{s}'|)$$

Model results



Model results



Model Validation

		Sulfate	Nitrate	Ammonium	Soil	Carbon	PM25
RMSE	Tobit w/o CMAQ	1.347	2.257	0.858	1.363	3.073	6.298
	Tobit	1.151	1.641	0.724	1.307	2.851	5.393
CRPS	Tobit w/o CMAQ	0.639	0.758	0.374	0.468	1.064	3.023
	Tobit	0.554	0.558	0.329	0.438	0.885	2.452
EmpCov	Tobit w/o CMAQ	0.935	0.931	0.933	0.907	0.923	0.924
	Tobit	0.920	0.930	0.924	0.915	0.921	0.906

Validation based on randomly selecting 10% of stations as hold outs.

Run times

Total run time for model fitting (50,000 iterations):

- CPU - 7.7 hours $\times 52$ weeks
- CPU+GPU - 4.8 hours

Total run time for model prediction at 5950 locations (1,000 iterations):

- CPU - 7.2 hours $\times 52$ weeks
- CPU+GPU - 4.3 hours

One run takes about 775 hours total on CPU alone, 473 on CPU and GPU.

Model fitting performance

Parameter	CPU (secs)	CPU+GPU (secs)	Rel. Performance
β_0, β_1	0.00029	0.00030	0.97
$\beta_0(s)$	0.09205	0.09132	1.00
σ	0.00383	0.00385	0.99
ϕ	0.46084	0.25174	1.83
ϵ	0.00003	0.00003	1.00
Total	0.55708	0.34729	1.60

Acknowledgments

- Alan Gelfand - Duke
- Dave Holland - EPA
- Erin Schliep - Duke
- Wyatt Apel - NERL

Disclaimer - The U.S. Environmental Protection Agency through its Office of Research and Development partially collaborated in this research. Although it has been reviewed by the Agency and approved for publication, it does not necessarily reflect the Agency's policies or views.

Email : rundel@gmail.com

Presentation : <http://github.com/rundel/Presentations/>

Paper : Rundel C., Schliep E., Holland D., Gelfand A. (2014) A data fusion approach for space-time analysis of speciated PM_{2.5}.
The Annals of Applied Statistics. In submission

Email : rundel@gmail.com

Presentation : <http://github.com/rundel/Presentations/>

Paper : Rundel C., Schliep E., Holland D., Gelfand A. (2014) A data fusion approach for space-time analysis of speciated PM_{2.5}.
The Annals of Applied Statistics. In submission

Questions?

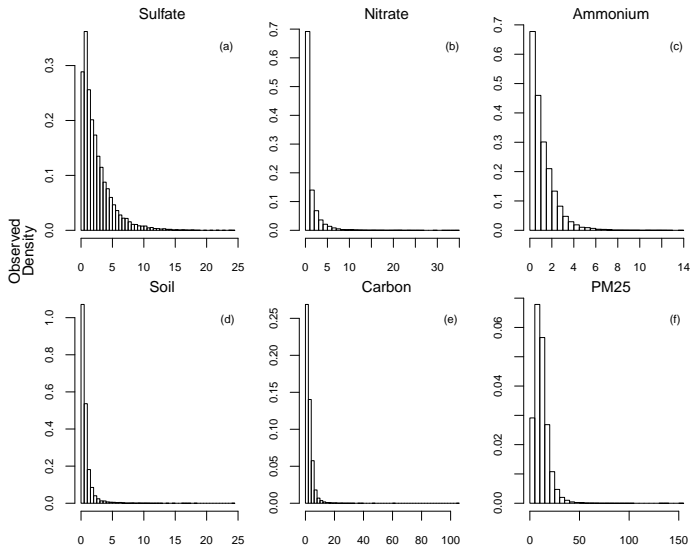
System Specs:

- 4 core Intel i5-2500K @ 3.30 GHz
- 16 GB DDR3 @ 1333 MHz
- GeForce GTX 460

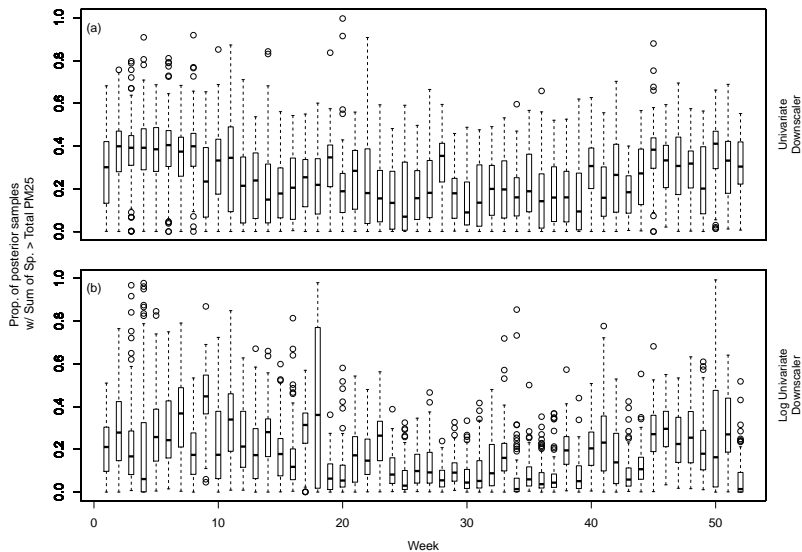
Software Specs:

- Ubuntu 13.10
- OpenBlas 0.2.8
- CUDA 6.0RC
- Magma 1.4.1

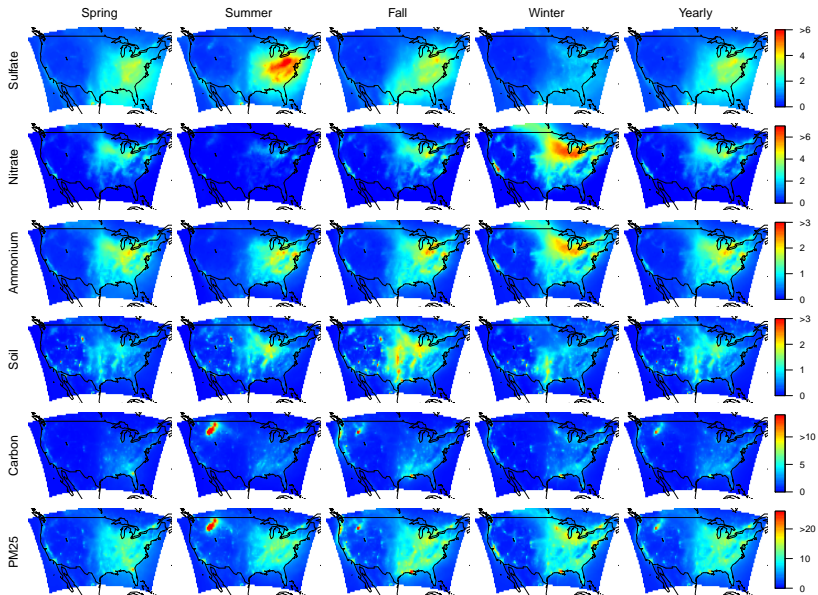
Species Distributions

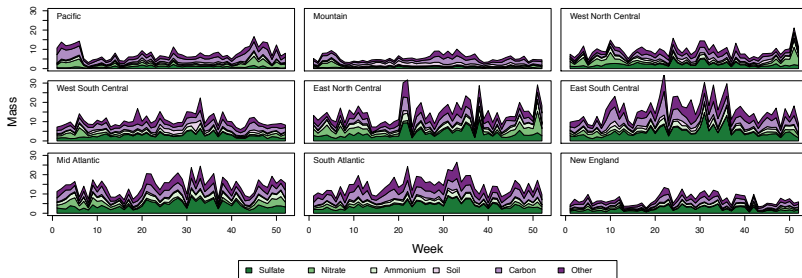


Univariate Model Exceedance



CMAQ Maps





Regions

