

Statistica Sinica Preprint No: SS-2016-0162R2

Title	Statistical-Physical Estimation of Pollution Emission
Manuscript ID	SS-2016-0162.R2
URL	http://www.stat.sinica.edu.tw/statistica/
DOI	10.5705/ss.202016.0162
Complete List of Authors	Emre Barut Youngdeok Hwang and Kyongmin Yeo,
Corresponding Author	Emre Barut
E-mail	barut@gwu.edu
Notice: Accepted version subject to English editing.	

STATISTICAL-PHYSICAL ESTIMATION OF POLLUTION EMISSION

Youngdeok Hwang, Emre Barut and Kyongmin Yeo

IBM Thomas J. Watson Research Center and George Washington University

Abstract: Air pollution is driven by non-local dynamics, in which air quality at a site is determined by transport of pollutants from distant pollution emission sources to the site by atmospheric processes. In order to understand the underlying nature of pollution generation, it is crucial to employ a physical knowledge to account for the pollution transport by wind. However, in most cases, it is not possible to utilize the physics models to obtain useful information, as it requires massive calibration and computation. In this paper, we propose a method to estimate the pollution emission from the domain of interest, by using both the physical knowledge and observed data. The proposed method uses an efficient optimization algorithm to estimate the emission from each of the spatial locations, while incorporating the physics knowledge. We demonstrate the effectiveness of the new method through a simulation study that mimics the real application.

Key words and phrases: alternating direction method of multipliers, dispersion, inverse model, penalized regression.

1. Introduction

Air pollution is produced by both natural and anthropogenic emissions and transported via physical processes driven by wind. The pollutant can be from a variety of sources, including traffic, fossil fuel uses, or burning natural biomass (World Health Organization, 2005). There have been a substantial progress in environmental science and engineering in developing computational models to forecast the evolution of physical processes. The notable examples include Weather Research and Forecasting coupled with Chemistry (WRF-Chem, Fast et al., 2006) and Community Multi-scale Air Quality Model (CMAQ, Byun and Schere, 2006) among many others. These models use knowledge from the physical and chemical processes to construct a system of partial differential equations which compute

the transport of the pollutant particles from where they are emitted to its nearby areas under the given weather conditions.

The quality of the physical model prediction heavily depends on the accuracy of the information as to how much the pollution is emitted from each geographic locations, among many model input parameters. However, such pollution emission information is often inaccurate, incomplete, or even unavailable. For example, Fu et al. (2012) compared the smoke emission rates from wild fire estimated by two most widely used models, to find that the difference can be from a factor of five to eight. To estimate the emission information more accurately, it is critical to reduce such a large uncertainty in the emission information. Conceptually, this is the ‘inverse’ problem, as opposed to the ‘forward’ problem which generates the data using the given model parameters; the interest is estimating the physical model parameters using the model outputs and real observation data. In applied mathematics and computational physics, it is often referred as *Inverse problems* (e.g., Biegler et al., 2011).

From a statistical point of view, this is closely related to *calibration* of a physics model. In this direction, Kennedy and O’Hagan (2001) solved the calibration problem using Bayesian framework with Gaussian process. Higdon et al. (2008) proposed an extension of Kennedy and O’Hagan (2001) to incorporate the high dimensionality of the computer model outputs. Liu and West (2009) proposed combining Bayesian multivariate dynamic linear models with Gaussian processes for a model with time-dependent functional outputs.

Another effort in different direction to solve the air pollution problem is to build statistical air quality models using the measurement data. Carroll et al. (1997) built a spatial-temporal model for hourly ozone level using Gaussian random field. Paciorek et al. (2009) developed a practical modeling approach to solve the epidemiological problem. Ghosh et al. (2010) studied formation and deformation of atmospheric concentrations of total nitrate using the empirical chemical relationship and dynamic statistical models. Williams et al. (2011) estimated the pollution source direction using the dispersion model, treating the computer model outputs as given. Although proven to be useful, these efforts do not incorporate enough physical knowledge as to the pollution generation process.

In utilizing a computational model to solve the air quality problem, an important interest is to make use of observed data obtained from a monitoring network located over the spatial domain to estimate the emission source information. It is crucial to include the knowledge from physics into the modeling framework to incorporate the continuously changing dynamic nature of weather conditions and pollution emission. Physics model can introduce such components into the modeling framework, but it is only possible when it is paired with the appropriate statistical modeling.

We propose a statistical framework to exploit physical assumptions of the model linking the computational physics model to the prediction, while obtaining the critical information about the nature of pollution generation. Although it is developed to link the physical dispersion model and the empirical observations, this interdisciplinary framework can also be used for solving the challenging problem of using both physical and statistical knowledge. Practicality is the major concern of our approach, otherwise the solution adds substantial amount of computational and operational complexity, which eventually overkills the problem in practice. At a conceptual level, our work is aligned with Malmberg et al. (2008) in that we combine the statistical model and physical knowledge.

Main contribution of our paper is two-fold. First, we introduce a reduced order physics model designed such that it can be well paired with a statistical method to estimate the emission intensity surface of the area of interest. We describe how our physics model can be reformulated into a regression problem. Second, we develop an efficient algorithm to estimate the detailed emission information at each location changing over time. Our model also imposes sparsity on the estimated coefficients, so that the dimensionality of the inverse problem as well as prior knowledge on the emission can be well incorporated. Our method provides an insight regarding how much pollution emission is produced so that it can be used for policy or administrative decision purposes.

The remainder of the paper is organized as follows. Section 2 describes the fundamental physical process regarding the transport of a pollution on a spatial domain. Section 3 gives the statistical model that we propose. Section 4 describes optimization methodology to estimate the parameters. Section 5 presents the application of the proposed method to synthetic air-quality monitoring data,

where the results and interpretations are presented together. We conclude with some remarks and discussion in Section 6.

2. Physical Process

In this section, we provide a brief introduction of our physical model, and describe how it is connected to the statistical model in Section 3.

The building block of our model is the *dispersion* process. The dispersion process consists of two major contributions, *advection* and *diffusion*. Advection is transfer of pollutants from one location to another, following the streamline of wind. Diffusion characterizes the movement of pollution from a region of high concentration to low concentration due to mixing by atmospheric turbulence.

Let $\phi(\mathbf{s}, t)$ denote the pollution concentration at location over the computational domain \mathbf{s} at time t . Mathematically, the dispersion process can be expressed as

$$\frac{\partial \phi(\mathbf{s}, t)}{\partial t} = -\nabla \cdot (\mathbf{u}(\mathbf{s}, t)\phi(\mathbf{s}, t)) + \nabla \cdot \{\mathbf{K}(\mathbf{s}, t; \mathbf{u}) \cdot \nabla \phi(\mathbf{s}, t)\} + Q(\mathbf{s}, t), \quad (2.1)$$

which implies that the temporal change of pollution concentration at a location is determined by the advection (the first term on the right hand side) and the diffusion (the second term). Here, $\mathbf{u}(\mathbf{s}, t)$ is the velocity of wind, which can be obtained from a numerical weather prediction model, and \mathbf{K} is a diffusion coefficient matrix, defining the rate of mixing of the pollutant. The last term $Q(\mathbf{s}, t)$ represents the rate of pollutant emission, i.e. the newly added pollution, at location \mathbf{s} and time t .

We illustrate the model in (2.1) using a simple example. Consider a one-dimensional computational domain of three grids, $\mathbf{s} = (s_1, s_2, s_3)$. We assume that these three grids are the mid-subset of a much larger domain, so that boundary conditions do not affect the calculation. The model output over these grids at each time point t can be expressed as a vector of length three. Now, suppose that the initial concentration at $t = 0$ is zero everywhere (including the outside of three grid points), the wind $\mathbf{u} = (1, 1, 1)$, uniform diffusion $K = 1/4$, the spatial grid $\delta_x = 1$ and time step $\delta_t = 1$. The emission at these three grids are the same for all t and are given by $\boldsymbol{\beta}$, that is, $Q(\mathbf{s}, t) = Q(\mathbf{s}) = \boldsymbol{\beta} = (\beta_1, \beta_2, \beta_3)$.

Using the first order Taylor expansion with finite-difference discretization (Moin, 2010), (2.1) at j -th grid s_j at time $t + 1$ can be approximated by

$$\phi(s_j, t+1) \cong \phi(s_j, t) + \phi(s_{j-1}, t) - \phi(s_j, t) + \frac{1}{4} \{ \phi(s_{j-1}, t) - 2\phi(s_j, t) + \phi(s_{j+1}, t) \} + \beta_j,$$

for $j = 1, 2, 3$. Then for $t = 1, 2$,

$$\phi(\mathbf{s}, 1) = \begin{bmatrix} 1/2 & 1/4 & 0 \\ 5/4 & 1/2 & 1/4 \\ 0 & 5/4 & 1/2 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix}; \phi(\mathbf{s}, 2) = \begin{bmatrix} 22/16 & 0 & 1/16 \\ 0 & 22/16 & 0 \\ 25/16 & 0 & 22/16 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix}, \quad (2.2)$$

and further calculations proceed similarly.

Now we link the calculated dispersion to monitoring observations using this formulation. We define $\tilde{X}_{t,ij}$ the dispersion at i th monitoring location, \mathbf{s}_i , at time t , originated from the j th grid computed by (2.1). In (2.2), for example, $(\tilde{X}_{1,11}, \tilde{X}_{1,12}, \tilde{X}_{1,13}) = (1/2, 1/4, 0)$; $(\tilde{X}_{1,21}, \tilde{X}_{1,22}, \tilde{X}_{1,23}) = (5/4, 1/2, 1/4)$. The pollution concentration observation at \mathbf{s}_i at time t and the associated noise are denoted by $y_{t,i}$ and $\epsilon_{t,i}$, respectively, for $t = 1, \dots, T$ and $i = 1, \dots, n$. Denoting the emission from grid j by β_j and assuming there are \tilde{p} grids, we have the following data generating process:

$$y_{t,i} = \sum_{j=1}^{\tilde{p}} \tilde{X}_{t,ij} \beta_j + \epsilon_{t,i}, \text{ for } i = 1, \dots, n, \text{ and } t = 1, \dots, T. \quad (2.3)$$

The measurements are available as the observed concentration level similar to $\phi(\mathbf{s}_i, t)$, but with the added noise. The physical dispersion process in (2.1) is incorporated in $\tilde{X}_{t,ij} \beta_j$, and the inherent randomness of the real measurements are modeled with $\epsilon_{t,i}$.

The convenience comes from the linearity of the emission intensity β . When calculating the dispersion, the emission needs not be known. As seen above, the dispersion from a source at s_1 can be calculated with a unit vector $Q(\mathbf{s}) = (1, 0, 0)$ for the desired time interval. Once the dispersion field is obtained, the intensity can be scaled easily by multiplying it with β_1 . Similarly, dispersion from the other two coordinates can be calculated with $(0, 1, 0)$ and $(0, 0, 1)$ and scaled with β_2 and β_3 , which makes $Q(\mathbf{s})$ linear in $(\beta_1, \beta_2, \beta_3)$.

In real applications, the goal is to estimate the pollution emission surface $Q(\mathbf{s})$ over the spatial domain. Without further simplifications, estimation of the pollution emission over all of the spatial domain requires solving an infinite dimensional problem. Because of this challenge, (2.1) has been used with strong assumptions on the sources. For example, Keats et al. (2007) considered the problem of finding the location and emission intensity of a source, but assuming there exists only one.

Alternatively, we approximate the emission surface by using a basis representation

$$Q(\mathbf{s}) = \sum_{j=1}^p \beta_j \Phi(\|\mathbf{s} - \mathbf{v}_j\|; \tau) = \sum_{j=1}^p \beta_j \Phi_j, \quad (2.4)$$

where β_j is the emission intensity associated with j th component, and $\Phi(\|\mathbf{s} - \mathbf{v}_j\|; \tau) = \exp(-\|\mathbf{s} - \mathbf{v}_j\|^2/2\tau^2)/2\pi\tau^2$ with $\tau > 0$, a Gaussian density kernel centered at location \mathbf{v}_j . With a given τ , $\Phi(\|\mathbf{s} - \mathbf{v}_j\|; \tau)$ serves as pre-determined basis functions. Then, the emission surface is approximated by a sum of p smooth kernels, centered at the fixed grid points covering the domain, $\mathbf{v}_1, \dots, \mathbf{v}_p$. The number of kernels is chosen to be much smaller than the number of grids, i.e. $p \ll \tilde{p}$, so it makes the problem tractable. The amount of emission contribution from the area around \mathbf{v}_j is estimated by β_j . The computation is still conducted over \tilde{p} grids, but the emission from each kernel is assumed to behave together. Figure 1 shows two basis functions at two different locations, where the entire domain is divided into 64 cells, and each cell is represented by the basis centered at the grid in dashed-lines.

As an illustration, Figure 2 shows a series of dispersion model outputs from two sources in the form of (2.4) marked by the red circles, for three consecutive hours. The arrows represent the wind field changing over time, where the color contours depict $\phi(\mathbf{s}, t)$ due to the emission concentration. In the following section, we illustrate how (2.3) and (2.4) are paired to solve the real problem.

3. Model

In this section, we present our model extended from equations (2.3) and (2.4). As illustrated in the discussion following equation (2.4), we divide the domain of

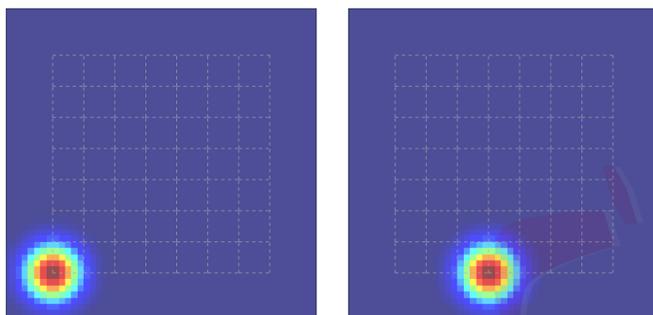


Figure 1: Illustration of our model approach. Dashed grid show the location of centers of basis functions, eight in both horizontal and vertical direction, where two panels show the two sites located at the two grid points of the domain.

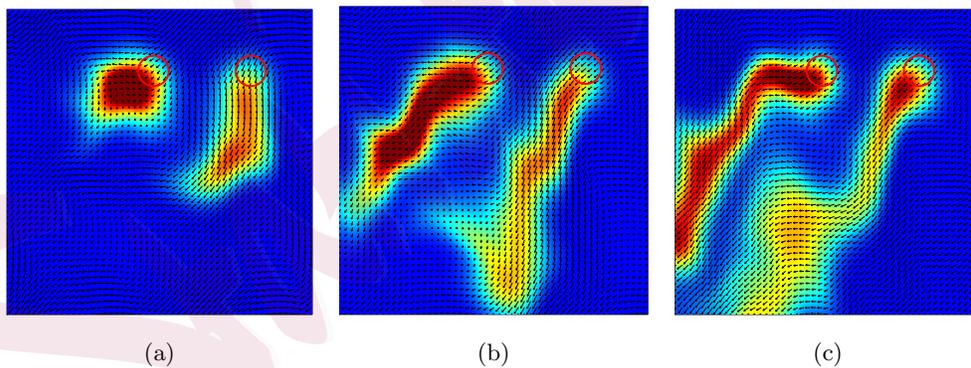


Figure 2: A series of model outputs depicting the dispersion from two sources marked by the red circles, where the snapshots are taken at hour 1 to 3 (a-c). The arrows in the background indicate the changing wind field.

interest into p sources. We further assume that the emission rate of each source has a diurnal pattern according to the 24-hour cycle of the day. Hence, for $j = 1, \dots, p$, there is a vector of intensity parameter $\boldsymbol{\beta}(j) = (\beta_{j,1}, \dots, \beta_{j,24})$. The dispersion emitted from source j and located at station i at time t , calculated using (2.1), is denoted by $X_{t,ij}$. Following these assumptions, we have

$$y_{t,i} = \sum_{j=1}^p X_{t,ij} \beta_{j,h(t)} + \epsilon_{t,i}, \quad (3.1)$$

where $h(t) = (t \bmod 24) + 1$. The equation implies that, each measurement is represented as a sum of p source locations' contribution, each of which is further decomposed into sum of 24 components. Similar idea was adopted in the classical chemical mass balance model (Christensen and Gunst, 2004), but in a simpler form. The error terms $\epsilon_{t,i}$ are assumed to be independent of each other.

We stack the hourly emission coefficients for each location, $\boldsymbol{\beta}(j)$, to define the $p \times 24$ matrix $\boldsymbol{\beta}$ such that j th row of $\boldsymbol{\beta}$ is $\boldsymbol{\beta}(j)$, i.e. $\boldsymbol{\beta} = [\boldsymbol{\beta}(1)^\top, \dots, \boldsymbol{\beta}(24)^\top]^\top = (\beta_{j,k})$. Although this problem can easily be solved with multivariate, multi-response linear regression, there are issues due to high dimensionality. In most practical applications, the number of possible emission locations, p , is high, and estimation of hourly emissions requires us to fit $p \times 24$ variables. On the other hand, we expect to have a small number of samples, $n \times T$. Simulating pollution dispersion at a high resolution requires substantial computational effort, and pollution sensors are often costly, which lead to small T and n , respectively.

Thus, in a setting where the number of variables is on the same order of magnitude as the number of samples, we have to utilize prior knowledge and enforce a special structure on the estimated coefficients. This study was done with emphasis on urban spaces, where two patterns stand out (Gurjar et al., 2004; Morawska, 2006; Saarikoski et al., 2008):

1. Most locations emit negligible amount of pollution.
2. Major pollution sources are traffic (i.e. motor vehicles, airports and sea-ports) and industrial areas, which produce the pollution in different hourly patterns.

We formulate our model based on these two findings.

Based on the first pattern, $\beta(j)$ should be zero for many j since most locations do not generate significant pollution. In other words, the rows of β should be sparse. We enforce this assumption with the group lasso penalty (Yuan and Lin, 2006; Simon and Tibshirani, 2012), by adding a penalty term proportional to $\sum_{j=1}^p \|\beta(j)\|_2$, where $\|\cdot\|_2$ is the Euclidean norm. The sub-level sets of the group lasso penalty function contain coefficients for which only a few of the rows are non-zero, and hence this penalty ‘sparsifies’ the solution with respect to the rows of β . Group lasso is commonly used in high dimensional problems where coefficients are expected to be active (i.e. non-zero) in groups of variables. For more details on the theoretical properties of group lasso, see the aforementioned references.

According to the second finding, there are different sources of pollution, but the same sources tend to have similar hourly patterns. Traffic pollution occurs commonly around the main highways, usually spikes in the morning and evening during the rush hours, and are generally constant otherwise; while industrial areas often have emissions that peak around noon, and there are certain industrial areas (e.g. facilities that are not shutdown during the night) that emit a constant level of pollution throughout the day. Ideally, locations’ hourly pollution patterns are known in advance, and these can be enforced as conditions on the estimated coefficients, β . Unfortunately, in practice, it is usually difficult to know which regions might have which patterns, especially when the number of potential emission sources is high. Therefore, we assume that there are a few unknown daily patterns, and each location follows one of these patterns or a combination of them (e.g. emissions from locations that have power plants and highways would be given by the sum of two different daily patterns). These unknown daily patterns can also be thought as latent factors which generate the emission coefficients. In this setting, we expect the rows of β to be linearly dependent; that is $\text{rank}(\beta)$ should be small. We enforce this assumption with a nuclear norm penalty, which encourages sparsity in singular values (Candès and Recht, 2009). To see this, first observe that for an $p \times m$ matrix \mathbf{A} with ($p \geq m$), its nuclear norm $\|\cdot\|_*$ is given by

$$\|\mathbf{A}\|_* = \sum_{l=1}^m \sigma_l,$$

where σ_l is the l^{th} singular value of \mathbf{A} , l^{th} element of diagonal matrix of $\mathbf{\Sigma}$, defined by the singular value decomposition (SVD) of $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$. Hence this penalty controls linear dependencies among the rows of $\boldsymbol{\beta}$, $\boldsymbol{\beta}(j)$. Example 1 demonstrates how the penalty on singular values works in practice.

Example 1. We display sets of coefficients obtained by solving the nuclear norm penalized regression problem:

$$\min_{\boldsymbol{\beta}} n^{-1} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_F^2 + \lambda_{\text{nuc}} \|\boldsymbol{\beta}\|_*,$$

where $\|\cdot\|_F$ is the matrix Frobenius norm (i.e. square root of the sum of the squared entries of the matrix), $\mathbf{y} \in \mathbb{R}^{T \times 24}$, $\mathbf{X} \in \mathbb{R}^{T \times p}$ and $\boldsymbol{\beta} \in \mathbb{R}^{p \times 24}$ are the response, predictor and coefficient matrices, respectively. Note that these variables bear no relationship to other \mathbf{y} , \mathbf{X} and $\boldsymbol{\beta}$ used in the remainder of the paper. The term $\lambda_{\text{nuc}} \geq 0$ is a tuning parameter that controls the effect of the nuclear norm penalty; as λ_{nuc} increases, the minimizer is forced to have more linear dependency. We create a toy example with $p = 12$ and $T = 200$ and generate observations from equation (3.1). Other details, such as the choice of \mathbf{X} and ϵ , are relegated to the Supplementary material. The true coefficients, $\boldsymbol{\beta}$ and the estimates for three different choices of λ_{nuc} are given in Figure 3. The true coefficients are given by the set of three vectors: β_{Type1} , given by a constant vector (i.e. $\beta_{\text{Type1}} = (1, \dots, 1)^\top$), β_{Type2} is set equal to a concave quadratic function that peaks at the 12th hour, and β_{Type3} is a vector that contains zeroes except for hours 6-8 and hours 15-17. For each type, four β 's are generated which gives $\boldsymbol{\beta}$ with $p = 12$ when stacked.

The effects of the nuclear norm penalty on the fitted coefficients are depicted in Figure 3. Coefficients obtained with penalization tend to be more similar, as we force the rows of $\boldsymbol{\beta}$ to be more linearly dependent. The statistical benefits of the penalty are obvious; the estimate with no penalization has a high variance, and a high estimation error as in Figure 3b. When we introduce the nuclear norm penalty term, the estimates have lower matrix rank and since the true coefficients also have a low-rank, this results in estimates with low variance and estimation error as in Figure 3c. Note that, we also introduce some bias; with $\lambda_{\text{nuc}} = 0.25$, Type 2 locations' coefficients are underestimated, but the estimates are much

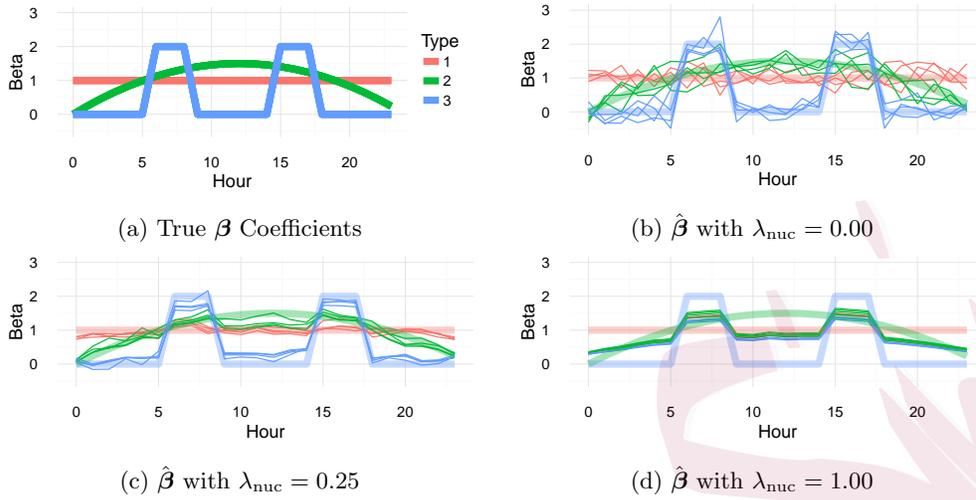


Figure 3: Effect of the nuclear norm penalty. Figure 3a displays the values of β coefficients for each type. Figures 3b,3c and 3d display the fitted β obtained using different λ_{nuc} , where each line is one row of $\hat{\beta}$ and the colors display the type of the coefficient. As the weight of the nuclear norm penalty increases, so does the linear dependence between the estimates.

closer to the truth than the results obtained without any penalties. When the penalty term is very large, $\lambda_{\text{nuc}} = 1.00$, $\hat{\beta}$ has rank one; all of the coefficients are given by the same vector multiplied by a scalar as seen in Figure 3d, where rows of $\hat{\beta}$ have the same shape. This is the case where we introduce massive bias by significantly shrinking the variance of the estimates.

Finally, last assumption in our model is that each source only can *add* emission to the ambient air. Any decay or deposit is assumed to be negligible and hence is absorbed in the error term. Thus we force all of the emission coefficients to be non-negative.

We combine these regularization penalties with a least-squares loss to obtain the following objective function:

$$\min_{\beta_{jk} \geq 0} \frac{1}{nT} \sum_{t=1}^T \sum_{i=1}^n \left(y_{t,i} - \sum_{j=1}^p X_{t,ij} \beta_{j,h(t)} \right)^2 + \lambda_{\text{gl}} \sum_{j=1}^p \|\beta(j)\|_2 + \lambda_{\text{nuc}} \|\beta\|_*, \quad (3.2)$$

where $\lambda_{\text{gl}} \geq 0$ and $\lambda_{\text{nuc}} \geq 0$ are tuning parameters for the group lasso and nuclear norm penalties, respectively. As before, $h(t) = (t \bmod 24) + 1$. The objective

function is a sum of convex functions, and hence is convex. The non-negativity constraint is also convex, resulting in a convex problem. In the next section, we discuss how we can efficiently minimize this constrained non-differentiable objective function.

4. Estimation

Multiple non-differentiable components in (3.2) make our problem challenging. For efficient optimization, we propose an *Alternating Direction Method of Multipliers* (ADMM) based approach (Boyd et al., 2010; Parikh and Boyd, 2014). ADMM is built to minimize objective functions with separable components. The algorithm works in a distributed and iterative manner.

As aforementioned, there are four main components in (3.2); (i) sum of squared errors, (ii) group Lasso penalty, (iii) nuclear norm penalty, (iv) nonnegative constraints. In this regard, the regularized least square problem in (3.2) is rewritten as

$$\text{minimize} \quad f_{\text{mse}}(\boldsymbol{\beta}) + f_{\text{gl}}(\boldsymbol{\beta}) + f_{\text{nuc}}(\boldsymbol{\beta}) + f_{\text{nn}}(\boldsymbol{\beta}), \quad (4.1)$$

where

$$\begin{aligned} f_{\text{mse}}(\boldsymbol{\beta}) &= (nT)^{-1} \sum_{t=1}^T \sum_{i=1}^n \left(y_{t,i} - \sum_{j=1}^p X_{t,ij} \beta_{j,h(t)} \right)^2, \\ f_{\text{gl}}(\boldsymbol{\beta}) &= \lambda_{\text{gl}} \sum_{j=1}^p \|\boldsymbol{\beta}(j)\|_2, \\ f_{\text{nuc}}(\boldsymbol{\beta}) &= \lambda_{\text{nuc}} \|\boldsymbol{\beta}\|_*, \end{aligned}$$

and $f_{\text{nn}}(\boldsymbol{\beta})$ is the indicator function for the non-negative orthant. That is, $f_{\text{nn}}(\boldsymbol{\beta}) = \sum_{j=1}^p \sum_{k=1}^{24} \delta_{\mathbb{R}_+}(\beta_{j,k})$, where

$$\delta_{\mathbb{R}_+}(z) = \begin{cases} 0 & \text{if } z \geq 0, \\ +\infty & \text{otherwise.} \end{cases}$$

Intuitively, our algorithm separately updates the solution to improve in

Algorithm 1 ADMM algorithm to obtain the estimator.

Set the initial estimates $\beta_{\text{mse}}^{(1)}, \beta_{\text{gl}}^{(1)}, \beta_{\text{nuc}}^{(1)}, \beta_{\text{nn}}^{(1)}, \bar{\beta}^{(1)}$;
Set the initial differences $\mathbf{u}_{\text{mse}}^{(1)}, \mathbf{u}_{\text{gl}}^{(1)}, \mathbf{u}_{\text{nuc}}^{(1)}, \mathbf{u}_{\text{nn}}^{(1)}$.
for $k = 1, \dots, 24$ **do** ▷ Reshape data
 for $i = 1, \dots, n$ **do**
 $\mathbf{X}(i, k) \leftarrow \{X_{t,ij} : (t \bmod 24) = k - 1\}$;
 $\mathbf{y}(i, k) \leftarrow \{y_{t,i} : (t \bmod 24) = k - 1\}$;
 end for
 $\mathbf{X}(k)^\top \leftarrow [\mathbf{X}(1, k)^\top, \dots, \mathbf{X}(n, k)^\top]$;
 $\mathbf{y}(k)^\top \leftarrow [\mathbf{y}(1, k)^\top, \dots, \mathbf{y}(n, k)^\top]$;
end for
for $m = 1, \dots, M$ **do** ▷ ADMM Iterations
 for $k = 1, \dots, 24$ **do**
 $\beta_{\text{mse},(:,k)}^{(m+1)} \leftarrow \left(\mathbf{X}(k)^\top \mathbf{X}(k) + \frac{1}{2\rho} \mathbb{I}_{p \times p} \right)^{-1} \mathbf{X}(k)^\top \mathbf{y}(k) + \frac{1}{2\rho} \left(\bar{\beta}_{:,k}^{(m)} - \rho \mathbf{u}_{\text{mse}}^{(m)}(k) \right)$; ▷ MSE
 end for
 $\beta_{\text{gl}}^{(m+1)} \leftarrow \text{sign}(\bar{\beta}^{(m)} - \mathbf{u}_{\text{gl}}^{(m)}) \left(1 - \frac{\lambda_{\text{gl}}}{\|\bar{\beta}^{(m)} - \mathbf{u}_{\text{gl}}^{(m)}\|_2} \right) \mathbf{1}(\bar{\beta}^{(m)} - \mathbf{u}_{\text{gl}}^{(m)} > \lambda_{\text{gl}}\rho)$; ▷ Group Lasso
 $\mathbf{U}^{(m)}, \mathbf{\Sigma}^{(m)}, \mathbf{V}^{(m)} \leftarrow \text{SVD}(\bar{\beta}^{(m)} - \mathbf{u}_{\text{nuc}}^{(m)})$;
 $\tilde{\mathbf{\Sigma}}^{(m)} \leftarrow (\mathbf{\Sigma}^{(m)} - \rho \lambda_{\text{nuc}} \mathbb{I}_{p \times p})_+$;
 $\beta_{\text{nuc}}^{(m+1)} \leftarrow \mathbf{U}^{(m)} \tilde{\mathbf{\Sigma}}^{(m)} \mathbf{V}^{(m)\top}$; ▷ Nuclear Norm
 $\beta_{\text{nn}}^{(m+1)} \leftarrow (\bar{\beta}^{(m)} - \mathbf{u}_{\text{nn}}^{(m)})_+$; ▷ Nonnegative Projection
 $\bar{\beta}^{(m+1)} \leftarrow \left(\beta_{\text{mse}}^{(m+1)} + \beta_{\text{gl}}^{(m+1)} + \beta_{\text{nuc}}^{(m+1)} + \beta_{\text{nn}}^{(m+1)} \right) / 4$; ▷ Consensus
 for $g = \{\text{mse}, \text{gl}, \text{nuc}, \text{nn}\}$ **do**
 $\mathbf{u}_g^{(m+1)} \leftarrow \mathbf{u}_g^{(m)} + \left(\beta_g^{(m+1)} - \bar{\beta}^{(m+1)} \right)$ ▷ Dual Variables
 end for
end for

$f_{\text{mse}}, \dots, f_{\text{nn}}$, and then gradually pulls the solutions toward their average while improving in each direction. The optimization is summarized in Algorithm 1. The detailed derivation of each step in Algorithm 1 is deferred to the Supplement.

The original problem in (3.2) is complex and requires a semidefinite program. Decomposition by (4.1), however, makes the minimization trivial. All of the proximal steps can be calculated in an expeditious manner, which leads to a very efficient and fast algorithm. This is a considerable benefit because the entire algorithm must be executed repeatedly for finding the appropriate tuning parameters. Furthermore, the memory requirement is only linear in the number of variables as the algorithm only tracks the parameters themselves.

Unfortunately, ADMM algorithms do not generally have convergence guar-

antees (Tran-Dinh and Cevher, 2015). By choosing a proper step size, however, those issues can be mitigated. For our problem, ρ needs to be chosen proportional to the minimum eigenvalue of the Hessian of f_{mse} . We can achieve that by setting $\rho \geq (nT)^{-1} \sum_{k=1}^{24} \lambda_{\min}(\mathbf{X}(k)^\top \mathbf{X}(k))$ where $\lambda_{\min}(\cdot)$ is the minimum eigenvalue and $\mathbf{X}(k)$ is the observations in \mathbf{X} that correspond to hour k , whose definition can be found in Algorithm 1.

This condition is required for convergence. However, the speed of convergence with this choice of ρ can be excruciatingly slow. In our analysis, we have observed that by adding a step-size selection procedure, we can empirically ensure convergence and a fast speed of convergence. In order to avoid an extra computational burden during the step-size selection, we first create a list of step-size candidates from a geometric sequence. In the first 20 steps of the algorithm, all candidates for ρ are tried and the one that gives the largest reduction for the cost function is chosen. In the following steps, only ρ that was chosen in the last round, and 4 other ρ candidates that are the closest in value to the previously chosen ρ are tested in the step-size search. By limiting the search, we can avoid extensive computation, which results in faster convergence.

There are two tuning parameters in the objective function λ_{gl} and λ_{nuc} . We employ a brute force search over a grid of possible values, and use cross validation to estimate the sample error for each parameter choice. Choosing large values for λ_{gl} or λ_{nuc} will force all the variables to be zero. Using Karush-Kuhn-Tucker (KKT) conditions, it can be shown that the variables are reduced to 0 when λ_{gl} is larger than $\max_k \bar{\lambda}_{\text{gl},k}$, where $\bar{\lambda}_{\text{gl},k} = \|\mathbf{X}(k)^\top \mathbf{y}(k)\|_2$. Similarly, the variables are shrunk to 0 when λ_{nuc} is larger than $\sqrt{\lambda_{\max}(\mathbf{Z}^\top \mathbf{Z})}$ where $\mathbf{Z}_{:,k} = \mathbf{X}(k)^\top \mathbf{y}(k)$.

Using these two facts, we choose two sequences $\mathbf{d}_{\lambda_{\text{gl}}}$, $\mathbf{d}_{\lambda_{\text{nuc}}}$ for λ_{gl} and λ_{nuc} , whose ranges are $\exp(-5, \log(\max_k \bar{\lambda}_{\text{gl},k}))$ and $\exp(-5, \log(\sqrt{\lambda_{\max}(\mathbf{Z}^\top \mathbf{Z})}))$, respectively. Then, the grid for λ_{gl} , λ_{nuc} candidates is given by $\mathbf{d}_{\lambda_{\text{gl}}} \times \mathbf{d}_{\lambda_{\text{nuc}}}$. In applications where a brute force grid search cannot be afforded due to time constraints, sequential Kriging optimization or global Bayesian optimization methods can be used to find the best tuning parameters (Huang et al., 2006; Snoek et al., 2012).

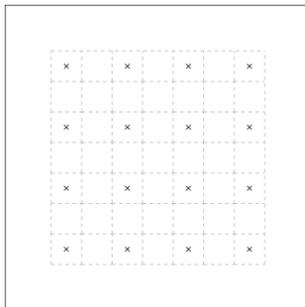


Figure 4: Station locations. The gridded lines indicate the location of the 64 source locations, where \times signs indicate the 16 monitoring locations.

5. Case Study

In this section, we illustrate our methodology and evaluate the performance of our proposed estimator with synthetic data. We considered a spatial domain of size 40 km by 40 km, divided into 64 potential source locations (8×8), which gives each grid a 5 km \times 5 km resolution. We generated 14 day long meteorological conditions by using a stochastic Fourier series to simulate atmospheric turbulence. In real data setting, one can use numerical weather prediction model (e.g., Skamarock et al., 2008). The details on simulated atmospheric turbulence are deferred to the Supplementary material.

Dispersion model calculates the transport from each pollution source transported to monitoring stations by using the simulated wind condition. We considered 16 monitoring stations, whose locations are depicted in Figure 4. Also, we considered 24 levels for diurnal cycle, based on the time of the day. This setup gives $n \times T = 16 \times 336 = 5376$ observations, each of which is associated with $1536 = 64 \times 24 = p \times 24$ variables and hence $\beta \in \mathbb{R}^{64 \times 24}$. The number of variables is comparable to the sample size, and any method that does not employ appropriate regularization is expected to overfit to the data.

We assumed that 16 sources are active among the 64 candidates, which resulted in the simulated $\beta^* \in \mathbb{R}^{64 \times 24}$, in which only 16 rows of β^* have non-zero elements. We divided the 64 source candidates into three different groups \mathcal{A} , \mathcal{B} and \mathcal{C} . Half of the active pollution sources were classified as group \mathcal{A} , and

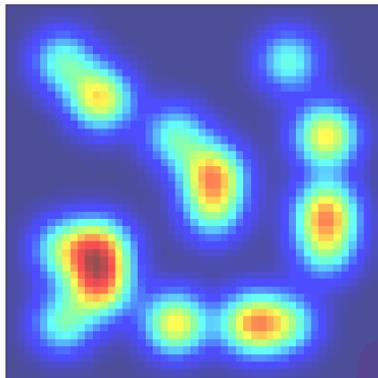


Figure 5: Surface of daily average emission in the simulation.

were assumed to be constantly active for all 24 hours, $\beta^*(j)^\top = \{1, \dots, 1\}$ for all $j \in \mathcal{A}$. The remaining eight active sources were assigned to group \mathcal{B} , and were fixed to have a diurnal cycle which steadily increases from the morning (6 AM) to noon (12 PM) and then steadily decreases until 6 PM. Sites in group \mathcal{C} , did not produce any pollutants, that is $\beta^*(j) = 0$ for $j \in \mathcal{C}$. The locations of active 16 pollution sources were fixed for all simulations. They were randomly placed, and their daily average emission is depicted in Figure 5. Note that we also tested the robustness of our simulation studies by changing the pollution sources, and the results were very similar.

The design matrix \mathbf{X} and coefficients β were determined through the calculated weather conditions and the fixed pollution sources. We then simulated data from the linear regression model

$$\mathbf{y}_i = \mathbf{x}_i^\top \beta^* + \varepsilon_i, \quad \varepsilon_i \sim \mathcal{N}(0, \mathbf{I}_T),$$

where \mathbf{x}_i^\top are computed from the dispersion model in (2.1). The design matrix was also scaled to have average variance of 1 for its columns, this gives a signal to noise ratio close to 1.

We compared three methods including ours:

1. *Non-negative Group Lasso (GL)*, the penalized least-squares estimator with a group L_2 -penalty and a non-negativity condition on the coefficients.

2. *Non-negative Nuclear-Norm (NN)*, the penalized least-squares estimator with a nuclear norm penalty and a non-negativity condition on the coefficients.
3. *Proposed method (GL+NN)*.

Note that GL method is simply our estimator with $\lambda_{\text{nuc}} = 0$; similarly the estimator for the NN method can be obtained by setting $\lambda_{\text{gl}} = 0$. All of the estimators can be fit using ADMM. To determine the tuning parameters for all estimators $\lambda = \{\lambda_{\text{gl}}, \lambda_{\text{nuc}}\}$, we ran 100 separate simulations and stored the λ chosen by 10-fold cross validation. The sequence of candidate λ were obtained from the suggested grid given in Section 3. In the simulations, λ_{gl} was fixed as the average of the 100 λ_{gl} that were chosen in the separate simulations; we repeated the same procedure for choosing λ_{nuc} .

The following two performance metrics were compared:

1. L_2 loss, which is defined as $\|\beta^* - \hat{\beta}\|_F^2$.
2. L_1 loss, which is defined as $\sum_j \|\beta^*(j) - \hat{\beta}(j)\|_1$.

We have also calculated these performance measures for subsets of β with respect to their groups.

Note that the design matrix used in the simulations, \mathbf{X} , was heavily correlated. This is common in our methodology, because dispersion from two nearby sites to one monitoring location are driven by the similar wind field, and hence behave similarly. Possibly as a result of this, in all of the 100 simulations, all estimators chose more conservative (i.e. smaller) λ values as these fits tended to have lower prediction error.

For each setting, we present the average of the performance measure based on 100 simulations. The results are depicted in Tables 1 and 2. We denote the estimation error for source j as $\Delta(j) = \beta^*(j) - \hat{\beta}(j)$. A boxplot of the L_1 and L_2 losses for different estimators is given in Figures 6 and 7.

The simulation results show that the proposed estimator that uses both nuclear norm and group Lasso penalty overperformed other estimators that make use of only one of those penalties. We see in Tables 1 and 2 that the proposed estimator had lower L_1 and L_2 loss uniformly across all groups of coefficients.

Method	$\sum_j \Delta(j) $	$\sum_{j \in \mathcal{A}} \Delta(j) $	$\sum_{j \in \mathcal{B}} \Delta(j) $	$\sum_{j \in \mathcal{C}} \Delta(j) $
GL	36.89 (0.19)	16.41 (0.09)	9.61 (0.07)	10.86 (0.12)
NN	41.51 (0.25)	13.39 (0.11)	11.74 (0.09)	16.37 (0.14)
GL+NN	30.65 (0.19)	12.43 (0.08)	8.42 (0.06)	9.79 (0.12)

Table 1: Performances of the Estimators Based on L_1 Loss. Standard errors are given in parentheses.

Method	$\sum_j \Delta(j)^2$	$\sum_{j \in \mathcal{A}} \Delta(j)^2$	$\sum_{j \in \mathcal{B}} \Delta(j)^2$	$\sum_{j \in \mathcal{C}} \Delta(j)^2$
GL	3.94 (0.03)	2.20 (0.02)	1.17 (0.01)	0.57 (0.01)
NN	3.27 (0.04)	1.21 (0.02)	1.27 (0.02)	0.79 (0.01)
GL+NN	2.43 (0.02)	1.23 (0.02)	0.81 (0.01)	0.39 (0.01)

Table 2: Performances of the Estimators Based on L_2 Loss. Standard errors are given in parentheses.

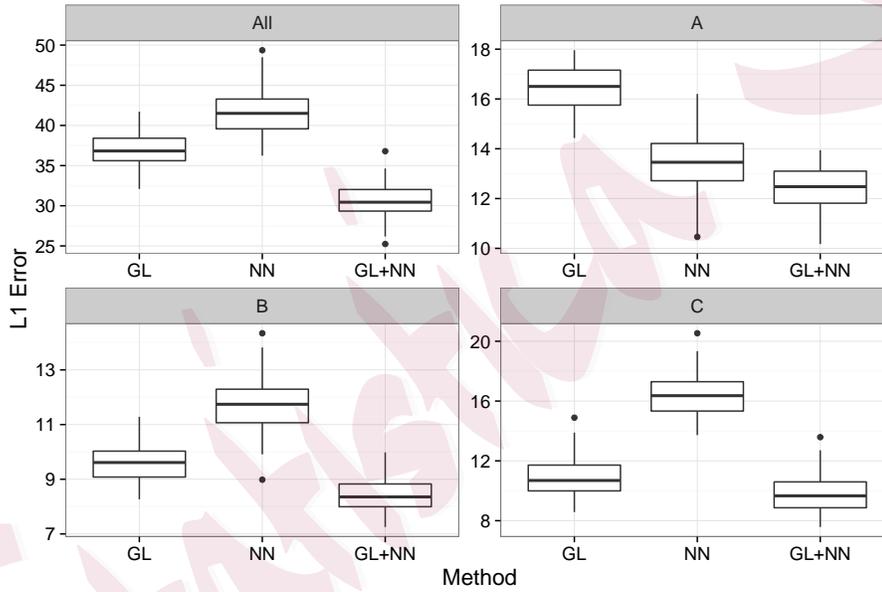


Figure 6: Boxplot of L_1 estimation errors of various estimators

This is an expected result because the true coefficients were sparse and had low rank.

Considering the performances across each coefficient group we see that our estimators had comparable errors in group \mathcal{C} to that of GL. Upon further analysis, we discovered that ours tends to pick extra sources since coefficients are forced to

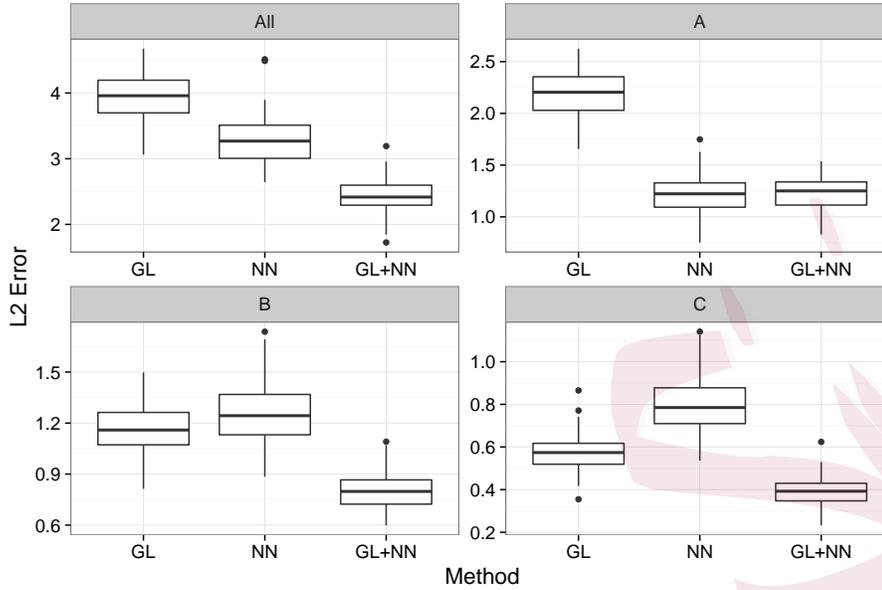


Figure 7: Boxplot of L_2 estimation errors of various estimators

be similar via the nuclear norm. Accordingly, it was common to see cases where a polluted location’s coefficients “bleed onto” nearby grids. From the observation stations’ point of view, two nearby sources had similar levels of dispersion. As a result, it became difficult to detect which one of nearby sources is the real polluter.

Comparing the nuclear norm based estimator to GL, we see that NN produced less-sparse estimates and wrongly selected some of the inactive regions. This can be seen by contrasting the proportion of estimation errors in each group over the total error; 24% of NN’s total L_2 loss is from the inactive variables (group \mathcal{C}), where as that number is 14% for GL. As a side note, although these small errors caused NN to have larger L_1 loss compared to GL, since the magnitudes of these estimates are very small, NN had less L_2 loss in total.

With regards to groups with active coefficients, GL had lower error in group \mathcal{B} but higher in group \mathcal{A} . This is because NN can generalize what it learns from one region to another, unlike GL.

We also note that as a consequence of the non-negativity constraint on the coefficients, the NN estimator returned sparse solutions (Slawski and Hein, 2013).

The median number of non-zero coefficients was 56, compared to 51 for group Lasso and 53 for our estimator. This also explains why the NN errors are not significantly worse compared to that of GL or our estimator.

Furthermore, all of the active emission sources were estimated to be active; hence the number of false negatives was zero for all methods across all simulations.

6. Discussion

In this paper, we propose a hybrid method to integrate the physical knowledge and statistical methodology to solve a challenging problem. Instead of using the typical approach from physics trying to incorporate more complex equations and parameters, we simplified the physical assumption to utilize the available resources and data. This simplification leads to a customized physics model to make a better pair with statistical method. A statistical model to incorporate the prior domain knowledge and information is proposed. An efficient algorithm is proposed to solve the difficult optimization problem.

We briefly remark on the potential future research. First, when the interest is in exploring the uncertainties associated with the similar physics, a natural extension of our work would be a Bayesian methodology incorporating the prior physical knowledge and previous emission inventory data. Second, there exist more uncertainties related to the weather model output and other physical parameters. Although it is known that such uncertainties can affect the accuracy of the final inference, it is challenging to quantify their impact. A thorough analysis to incorporate such uncertainties can be useful and informative. Finally, an extension of our method to the case where the independence assumption is violated can be considered. Research effort is needed to incorporate the dependency structure through spatio-temporal methods. Such methodology needs to address the challenges arising due to estimation of the covariance matrix of the residuals.

References

Biegler, L., Biros, G., Ghattas, O., Heinkenschloss, M., Keyes, D., Mallick, B., Marzouk, Y., Tenorio, L., van Bloemen Waanders, B., and Willcox, K. (2011),

Large-Scale Inverse Problems and Quantification of Uncertainty, West Sussex, United Kingdom: John Wiley & Sons.

Boyd, S., Parikh, N., Chu, E., Peleato, B., and Eckstein, J. (2010), “Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers,” *Foundations and Trends in Machine Learning*, 3, 1–122.

Byun, D. and Schere, K. L. (2006), “Review of the Governing Equations, Computational Algorithms, and Other Components of the Models-3 Community Multiscale Air Quality (CMAQ) Modeling System,” *Applied Mechanics Reviews*, 59, 51–77.

Candès, E. and Recht, B. (2009), “Exact Matrix Completion via Convex Optimization,” *Foundations of Computational Mathematics*, 9, 717–772.

Carroll, R., Chen, E., Li, T., Newton, H., Schmiediche, H., and Wang, N. (1997), “Ozone Exposure and Population Density in Harris County, Texas,” *Journal of the American Statistical Association*, 92, 392–404.

Christensen, W. F. and Gunst, R. F. (2004), “Measurement Error Models in Chemical Mass Balance Analysis of Air Quality Data,” *Atmospheric Environment*, 38, 733 – 744.

Fast, J. D., Gustafson, W. I., Easter, R. C., Zaveri, R. A., Barnard, J. C., Chapman, E. G., Grell, G. A., and Peckham, S. E. (2006), “Evolution of Ozone, Particulates, and Aerosol Direct Radiative Forcing in the Vicinity of Houston Using a Fully Coupled Meteorology-chemistry-aerosol Model,” *Journal of Geophysical Research: Atmospheres*, 111.

Fu, J. S., Hsu, N. C., Gao, Y., Huang, K., Li, C., Lin, N.-H., and Tsay, S.-C. (2012), “Evaluating the Influences of Biomass Burning During 2006 BASE-ASIA: a Regional Chemical Transport Modeling,” *Atmospheric Chemistry and Physics*, 12, 3837–3855.

Ghosh, S., Bhawe, P., Davis, J., and Lee, H. (2010), “Spatio-Temporal Analysis of Total Nitrate Concentrations Using Dynamic Statistical Models,” *Journal of the American Statistical Association*, 105, 538–551.

- Gurjar, B., Van Aardenne, J., Lelieveld, J., and Mohan, M. (2004), “Emission estimates and trends (1990–2000) for megacity Delhi and implications,” *Atmospheric Environment*, 38, 5663–5681.
- Higdon, D., Gattiker, J., Williams, B., and Rightley, M. (2008), “Computer Model Calibration Using High-Dimensional Output,” *Journal of the American Statistical Association*, 103, 570–583.
- Huang, D., Allen, T., Notz, W., and Zeng, N. (2006), “Global Optimization of Stochastic Black-Box Systems via Sequential Kriging Meta-Models,” *Journal of Global Optimization*, 34, 441–466.
- Keats, A., Yee, E., and Lien, F.-S. (2007), “Bayesian inference for source determination with applications to a complex urban environment,” *Atmospheric Environment*, 41, 465 – 479.
- Kennedy, M. C. and O’Hagan, A. (2001), “Bayesian Calibration of Computer Models,” *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, 63, pp. 425–464.
- Liu, F. and West, M. (2009), “A Dynamic Modelling Strategy for Bayesian Computer Model Emulation,” *Bayesian Analysis*, 4, 393–411.
- Malmberg, A., Arellano, A., Edwards, D. P., Flyer, N., Nychka, D., and Wikle, C. (2008), “Interpolating Fields of Carbon Monoxide Data Using a Hybrid Statistical-physical Model,” *Annals of Applied Statistics*, 2, 1143–1580.
- Moin, P. (2010), *Fundamentals of Engineering Numerical Analysis*, New York, NY: Cambridge University Press, 2nd ed.
- Morawska, L. (2006), “Motor Vehicle Emissions as a Source of Indoor Particles,” *Indoor Environment: Airborne Particles and Settled Dust*.
- Paciorek, C. J., Yanosky, J. D., Puett, R. C., Laden, F., and Suh, H. H. (2009), “Practical Large-scale Spatio-temporal Modeling of Particulate Matter Concentrations,” *The Annals of Applied Statistics*, 3, 370–397.
- Parikh, N. and Boyd, S. (2014), “Proximal Algorithms,” *Foundations and Trends in Optimization*, 1, 123–231.

- Saarikoski, S., Timonen, H., Saarnio, K., Aurela, M., Järvi, L., Keronen, P., Kerminen, V., and Hillamo, R. (2008), “Sources of organic carbon in fine particulate matter in northern European urban air,” *Atmos. Chem. Phys.*, 8, 6281–6295.
- Simon, N. and Tibshirani, R. (2012), “Standardization and the Group Lasso Penalty,” *Statistica Sinica*, 22, 983–1001.
- Skamarock, W. C., Klemp, J. B., Dudhia, J., Gill, D. O., Barker, D. M., Duda, M. G., Huang, X.-Y., Wang, W., and Powers, J. G. (2008), *A Description of the Advanced Research WRF Version 3*, Boulder, Colorado, USA: National Center for Atmospheric Research.
- Slawski, M. and Hein, M. (2013), “Non-Negative Least Squares for High-Dimensional Linear Models: Consistency and Sparse Recovery Without Regularization,” *Electronic Journal of Statistics*, 7, 3004–3056.
- Snoek, J., Larochelle, H., and Adams, R. P. (2012), “Practical bayesian optimization of machine learning algorithms,” in *Advances in neural information processing systems*, pp. 2951–2959.
- Tran-Dinh, Q. and Cevher, V. (2015), “Splitting the Smoothed Primal-Dual Gap: Optimal Alternating Direction Methods,” *arXiv preprint arXiv:1507.03734*.
- Williams, B., Christensen, W. F., and Reese, C. S. (2011), “Pollution Source Direction Identification: Embedding Dispersion Models to Solve an Inverse Problem,” *Environmetrics*, 22, 962–974.
- World Health Organization (2005), “WHO Air Quality Guidelines for Particulate Matter, Ozone, Nitrogen Dioxide and Sulfur Dioxide—Global Update (WHO/SDE/PHE/OEH/06.02), World Health Organization,” in *Air quality guidelines - global update 2005*.
- Yuan, M. and Lin, Y. (2006), “Model Selection and Estimation in Regression with Grouped Variables,” *Journal of the Royal Statistical Society. Series B (Methodological)*, 68, 49–67.

IBM Thomas J. Watson Research Center, Yorktown Heights, NY 10598, U.S.A.

E-mail: yhwang@us.ibm.com

kyeo@us.ibm.com

Department of Statistics, George Washington University, Washington, DC 20052,
U.S.A.

E-mail: barut@email.gwu.edu

Statistica Sinica