



POLITECNICO
MILANO 1863

SCUOLA DI INGEGNERIA INDUSTRIALE
E DELL'INFORMAZIONE

A novel spatio-temporal model for Bayesian Source Apportionment for the PM_{10} pollution

Giorgio Aurina, Paolo Bettanin, Alberto Calanchi, Raffaele Scorza, Claudio Strino, Filippo Volpicelli

Tutors:

Michela Frigeri
Alessandra Guglielmi

Academic year:

2024-2025

Abstract: In recent years, environmental pollution has become one of the most critical topics in scientific research, driving numerous efforts to study the sources and impacts of various pollutants. Understanding the sources of air pollution is crucial for effective environmental policies and health risk assessments. This study presents a Bayesian spatio-temporal model for source apportionment of particulate matter pollution, addressing key challenges such as temporal dependence, spatial variability, and the unknown number of pollution sources. By leveraging Bayesian inference, we incorporate prior knowledge and constraints to improve source identification. To validate our approach, we implemented a Monte Carlo Markov Chain (MCMC) method alongside a Turing.jl model, allowing us to compare results and assess model reliability. Our implementation is tested on simulated data and is designed for application to real-world pollution datasets, such as those from ARPA Lombardia.

Key-words: source apportionment, Bayesian functional warping, infinite latent factors model

<https://github.com/GioAurina...>

<https://github.com/Paolofox01...>

1. Introduction

Environmental pollutants frequently exist as complex mixtures of chemical and non-chemical constituents. Using receptors (e.g., ambient monitors), we can measure individual constituents within a mixture. However, we lack direct information on the sources that generated them, nor do we know the contribution of each constituent to the total pollution.

Particulate matter (PM) and other airborne pollutants significantly affect air quality and human health. For example, complex pollutant mixtures like PM are often formed by multiple sources (e.g., biomass burning or traffic). Identifying these sources and quantifying their contributions is crucial for environmental policy and mitigation strategies.

Traditional source apportionment methods, such as Principal Component Analysis (PCA) and Absolute Principal Component Analysis (APCA), have limitations, including violations of non-negativity constraints and unrealistic assumptions about source independence. Most models for source apportionment represent observed constituents at a receptor as a linear combination of the unknown mixture concentrations from each source and the constituent composition of each source. However, these models often lack uniqueness and require additional constraints to provide meaningful solutions.

To address these challenges, we propose a novel spatio-temporal Bayesian model for source apportionment. This Bayesian framework allows for the incorporation of prior knowledge, enforces realistic constraints on source contributions, and accounts for both spatial and temporal dependencies in pollutant concentrations. By modeling the total pollution mass emitted from each source over time and its contribution to different pollutants, our approach provides a more robust and interpretable analysis of pollution dynamics.

The primary objectives of this study are:

- To estimate the contribution of different pollution sources to observed pollutant concentrations.
- To incorporate spatial dependencies, ensuring a more accurate representation of pollution dispersion.
- To model the temporal evolution of pollution sources and identify potential shifts in source contributions over time.
- To improve inference accuracy through a Bayesian hierarchical approach, integrating expert knowledge and observational data.

By addressing these challenges, this work aims to provide a powerful tool for policymakers and environmental agencies to better understand pollution sources and develop targeted mitigation strategies.

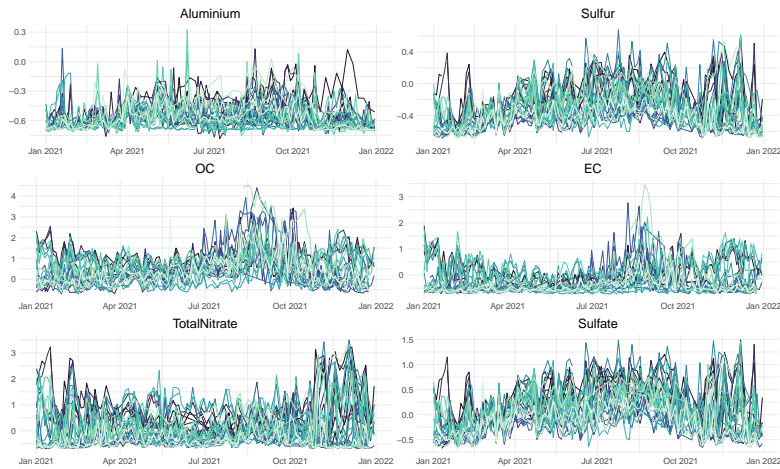


Figure 1: Different pollutants measured over time

2. Modeling

In this section, we present the detailed modeling approach used for source apportionment of PM10 pollution. The proposed model incorporates spatial and temporal dependencies through a Bayesian hierarchical framework. We introduce key components of the model step by step, providing mathematical formulations and explanations.

2.1. Notation in our hierarchical model

We define the following parameters and variables:

- $i = 1, \dots, N$: indices for the monitoring sites.
- $k = 1, \dots, K$: indices for the pollution sources.
- $c = 1, \dots, C$: indices for the different pollutants under analysis.
- $t = 1, \dots, T$: indices for the time steps (e.g., days).

The observed concentration of pollutant c at site s_i on day t is assumed to follow a Gaussian distribution [1]:

$$y^c(\mathbf{s}_i, t) \mid \mu^c(\mathbf{s}_i, t), \sigma_c^2 \stackrel{ind}{\sim} \mathcal{N}(\mu^c(\mathbf{s}_i, t), \sigma_c^2) \quad i = 1, \dots, N, \quad c = 1, \dots, C \quad (1)$$

where:

- $\mu^c(\mathbf{s}_i, t)$ represents the expected concentration, determined by the contribution of multiple sources.
- σ_c^2 accounts for the residual variance of pollutant C , with a marginal prior:

$$\sigma_c^2 \mid a, b \stackrel{iid}{\sim} \text{InvGamma}(a, b), \quad c = 1, \dots, C \quad (2)$$

2.2. Source Contribution Decomposition

The expected concentration $\mu^c(\mathbf{s}_i, t)$ is modeled as a weighted sum of the contributions from the K pollution sources [4]:

$$\mu^c(\mathbf{s}_i, t) = \sum_{k=1}^K h_k^c g_k(\mathbf{s}_i, t) \quad (3)$$

where:

- h_k^c represents the impact of source k on pollutant c , with a marginal prior:

$$\mathbf{h}_k \mid \alpha_0 \stackrel{iid}{\sim} \text{Dirichlet}(\alpha_0, \dots, \alpha_0), \quad k = 1, \dots, K \quad (4)$$

- $g_k(\mathbf{s}_i, t)$ is the spatiotemporal intensity function of source k at site \mathbf{s}_i and time t .

For $k = 1, \dots, K$ and $i = 1, \dots, N$, $t = 1, \dots, T$ we model $g_k(\mathbf{s}_i, t)$ as follows:

$$g_k(\mathbf{s}_i, t) = e^{\gamma_{ki}} f(t - \tau_i) \quad (5)$$

reflecting Bayesian functional warping [5].

2.3. Temporal Dependence

To capture the temporal dynamics of the pollution sources, we introduce temporal dependencies through a Gaussian process [5]. Specifically, the temporal dependence of the pollution sources is modeled as follows:

$$\mathbf{f}_k \mid \rho_k \stackrel{iid}{\sim} \mathcal{N}_T(\mathbf{0}, \mathbf{\Sigma}_{f_k}) \quad (6)$$

where the temporal covariance is defined by:

$$\mathbf{\Sigma}_{f_k}[m, n] = \exp \left\{ -\frac{\rho_k^2}{2} (t_m - t_n)^2 \right\} \quad (7)$$

and the correlation between the source k and site i over time is modeled by:

$$\mathbf{\Sigma}_{ki}[m, n] = \exp \left\{ -\frac{\rho_k^2}{2} (t_m - \tau_i - t_n)^2 \right\} \quad (8)$$

where ρ_k is a parameter that controls the degree of temporal correlation, and τ_i models the temporal shift for site \mathbf{s}_i . The marginal prior on ρ_k is given by:

$$\rho_k \mid a_\rho, b_\rho \stackrel{iid}{\sim} \text{Gamma}(a_\rho, b_\rho) \quad \text{for } k = 1, \dots, K \quad (9)$$

The temporal shifts τ_i are modeled as:

$$\boldsymbol{\tau} \sim \mathcal{N}_N \left(0, \sigma_\tau^2 \left(\mathbf{I}_N - \frac{1}{N+1} \mathbf{1}\mathbf{1}' \right) \right), \quad (10)$$

where the marginal variance σ_τ^2 is chosen according to the plausible range of latency of our problem

2.4. Spatial Dependence

To account for spatial correlation, we introduce a Gaussian Process prior over the site-specific scaling factors. The spatial dependence of the sources at different monitoring sites is modeled as follows [2]:

$$\boldsymbol{\gamma}_k \mid \boldsymbol{\beta}_k, \phi_k \stackrel{iid}{\sim} \mathcal{N}_N(\mathbf{X}\boldsymbol{\beta}_k, \mathbf{\Sigma}_{\gamma_k}) \quad (11)$$

where $\mathbf{\Sigma}_{\gamma_k}[i, j]$ represents the spatial covariance between sites \mathbf{s}_i and \mathbf{s}_j , which decays with the euclidean distance between the sites:

$$\mathbf{\Sigma}_{\gamma_k}[i, j] = \exp \left\{ -\frac{\phi_k^2}{2} \|\mathbf{s}_i - \mathbf{s}_j\|^2 \right\} \quad (12)$$

with ϕ_k controlling the spatial correlation length, modeled as:

$$\phi_k \mid a_\phi, b_\phi \stackrel{iid}{\sim} \text{Gamma}(a_\phi, b_\phi) \quad \text{for } k = 1, \dots, K \quad (13)$$

The vector \mathbf{X}_i contains the site-specific covariates, and $\boldsymbol{\beta}_k$ represents the regression coefficients, with a marginal prior:

$$\boldsymbol{\beta}_k \stackrel{iid}{\sim} \mathcal{N}_p(\mathbf{0}, \mathbf{I}_p) \quad (14)$$

3. Dataset

We use a dataset that primarily contains categorical information about the sites where pollutant measurements are taken. These categories include site type, location, and other relevant attributes that help contextualize the pollutant levels recorded at each site. This information is crucial for analyzing patterns and variations in air quality across different site characteristics.

FinalStations

Latitude	Longitude	Elevation	Local Site Name	County Name	City Name	LandUse	LandSetting
39.76168	-121.84047	1	Chico-East Avenue	Butte	Chico	IndustrialTraffic	SUBURBAN
41.560952	-124.083964	235	Redwood NP	Del Norte	Redwood National Park	LowPopulation	RURAL
38.9248	-119.97	1935	Lake Tahoe Community College	El Dorado	South Lake Tahoe	LowPopulation	RURAL
36.78538	-119.77321	96	Fresno - Garland	Fresno	Fresno	IndustrialTraffic	URBAN AND CENTER CITY
37.22064	-119.155557	2598	Kaiser Wilderness	Fresno	Not in a City	LowPopulation	RURAL
32.67618	-115.48307	1	Calexico-Ethel Street	Imperial	Calexico	ResidentialAgri	SUBURBAN
37.360684	-118.330783	1257	White Mountain Research Center - Owens Valley Lab	Inyo	Not in a City	ResidentialAgri	RURAL
35.356615	-119.062613	0	Bakersfield-California	Kern	Bakersfield	IndustrialTraffic	URBAN AND CENTER CITY
34.821922	-118.887598	1	Lebec-Peace Valley/Frazier Park Roads	Kern	Lebec	IndustrialTraffic	RURAL
34.06659	-118.22688	87	Los Angeles-North Main Street	Los Angeles	Los Angeles	ResidentialAgri	URBAN AND CENTER CITY
34.813034	-118.884819	1	Lebec	Los Angeles	Not in a City	IndustrialTraffic	RURAL
36.122979	-122.90944	76	Point Reyes NS Ranger Station	Marin	Point Reyes National Seashore	LowPopulation	RURAL
37.71325	-119.7062	1599	Yosemite NP - Turtleback Dome	Mariposa	Not in a City	LowPopulation	RURAL
38.088023	-119.178069	2561	Hoover Wilderness	Mono	Not in a City	LowPopulation	RURAL
33.83062	-117.93845	10	Anaheim	Orange	Anaheim	ResidentialAgri	SUBURBAN
39.81336	-120.47069	8	Portola	Plumas	Portola	ResidentialAgri	SUBURBAN
33.99958	-117.41601	250	Rubidoux	Riverside	Rubidoux	ResidentialAgri	SUBURBAN
33.463644	-116.971457	508	Aqua Tibia Wilderness	Riverside	Not in a City	LowPopulation	RURAL
38.613779	-121.368014	8	Sacramento-Del Paso Manor	Sacramento	Arden-Arcade	ResidentialAgri	SUBURBAN

Figure 2: California’s locations

3.1. Simulated Data

To evaluate the proposed approach, we generate a synthetic dataset that mimics real-world PM10 pollution trends. The dataset is created by simulating pollution levels across multiple monitoring sites, taking into account spatial and temporal variations.

The generation process starts by assigning each monitoring station a set of attributes, including location coordinates, elevation, and site classification (e.g., urban or suburban). These characteristics influence the expected pollution levels and their variation over time.

For each site, pollution levels are generated based on a combination of site-specific conditions and general trends observed in air quality data. The simulated pollution values reflect daily and seasonal fluctuations, ensuring that the dataset captures realistic patterns Figure 3.

This synthetic dataset serves as a controlled environment for testing and validating the source apportionment model before applying it to actual observations. By doing so, we ensure that our approach is robust and capable of identifying pollution sources accurately.

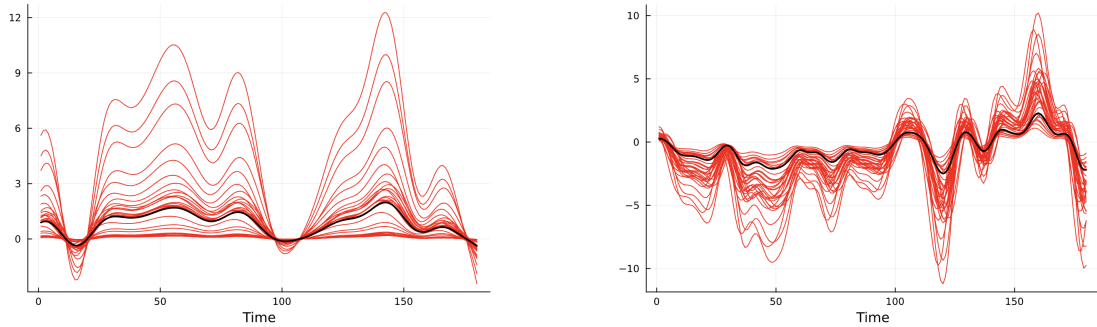


Figure 3: Global (black) and local contributions in each site i of source k , $i=1...32$, $k=1,2$

3.2. Generating data in Julia

```
1 function simulate_data(df, seed, K, n, n_time)
2     Random.seed!(seed)
3
4     # Parameters
5     # K = 2 #number of sources
6     # n = 32 #number of sites
7     # n_time = 180
8
9     intercept = ones(n)
10    x1 = zeros(n)
11    x2 = zeros(n)
12
13    for cont in 1:n
14        if df[cont, 8] == "SUBURBAN"
15            x1[cont] = 1
16        end
17
18        if df[cont, 8] == "URBAN_AND_CENTER_CITY"
19            x2[cont] = 1
20        end
21    end
22
23
24    sites = Matrix(DataFrame(Latitude = df[:, 1], Longitude = df[:, 2],
25        intercept = intercept, x1 = x1, x2=x2, Elevation = df[:, 3]))
26    maximum(euclid_dist(sites[:, 1], sites[:, 2], n))
27
28    theta = Dict{Int64, Dict{Any,Any}}{ }
29    theta[1] = Dict(
30        :rho => 0.1,
31        :phi => 1/300.0,
32        :gamma => zeros(32),
33        :beta => [-0.5, 0.5, 0.8, 0.1],
34        :tau => rand(Normal(0, 1), 32)
35    )
36
37    theta[2] = Dict(
38        :rho => 0.2,
39        :phi => 1/400.0,
40        :gamma => zeros(32),
41        :beta => [0.3, -0.4, -0.7, -0.1],
42        :tau => rand(Normal(0, 1), 32)
43    )
44
45    dat = generate_data(sites, n, K, n_time, theta)
46
47    theta_true = copy(theta)
48    df_new = Dict{Int64, DataFrame}{}
49    dat_trials = Dict{Int64, DataFrame}{}
50    for k in 1:K
51        theta_true[k][:gamma] = dat[:gamma][:,k]
52        df_new[k] = DataFrame(dat[:g][:, k, :], :auto)
53        dat_trials[k] = stack(df_new[k], variable_name = "time",
54            value_name = "value")
55        dat_trials[k].trial = repeat(1:n,n_time)
56    end
57 end
```

```

1  #' Generate data.
2  #'
3  #' @param n Number of trials.
4  #' @param n_time Number of time points.
5  #' @param theta Named list of parameter values.
6  function generate_data(sites, n, K, n_time, theta)
7
8      X = sites[:, 3:end] # Design matrix
9      coords = sites[:, 1:2]
10
11     t = range(1, stop=n_time, length=n_time)
12     f = zeros(Float64, K, n_time)
13     g = zeros(Float64, n, K, n_time)
14     gamma = zeros(Float64, n, K)
15
16     dist = euclid_dist(coords[:, 1], coords[:, 2], n)
17
18     for k in 1:K
19         theta_k = theta[k]
20         Sigma_f = sq_exp_kernel(t, theta_k[:rho]; nugget=1e-6)
21         Sigma_f_inv = inv(Sigma_f)
22         f[k, :] = rand(MvNormal(zeros(n_time), Sigma_f))
23
24         Sigma_gamma = get_Sigma_gamma(dist, theta_k[:phi])
25         gamma[:, k] = rand(MvNormal(X*theta_k[:beta], Sigma_gamma))
26         theta_k[:gamma] = gamma[:, k]
27         for i in 1:n
28             Sigma_i = get_Sigma_i(i, t, theta_k)
29
30             g[i, k, :] = get_mu_g(i, t, f[k,:], theta_k, Sigma_f_inv)
31         end
32     end
33
34     return Dict(:g => g, :f => f, :gamma => gamma)
35 end
36

```

4. Sample from the Joint Posterior Distribution

4.1. Full Conditional of g_k

We have the following joint distribution:

$$\mathcal{L}(\mathbf{f}_k, \mathbf{g}_k(\mathbf{s}_i)) = \begin{pmatrix} f_k(1) \\ \vdots \\ f_k(T) \\ a_i f_k(1 - \tau_i) \\ \vdots \\ a_i f_k(T - \tau_i) \end{pmatrix} \sim \mathcal{N}_{2T} \left(\begin{pmatrix} \mathbf{0} \\ \mathbf{0} \end{pmatrix}, \begin{bmatrix} \Sigma_{f_k} & a_i \Sigma_{k_i} \\ a_i \Sigma_{k_i}^T & a_i^2 \Sigma_{f_k} \end{bmatrix} \right),$$

where $a_i = e^{\gamma_{ki}}$.

This representation shows that $\mathbf{g}_k(\mathbf{s}_i)$ is the same Gaussian process of \mathbf{f}_k but evaluated at new points, effectively performing temporal kriging over τ_i , $i = 1, \dots, 32$. This allows us to derive the conditional distribution:

$$\mathbf{g}_k(\mathbf{s}_i) | \mathbf{f}_k, \gamma_{ki}, \tau_i \stackrel{ind}{\sim} \mathcal{N}_T(\boldsymbol{\mu}_{g_{ki}}, \boldsymbol{\Sigma}_{g_{ki}}), \quad i = 1, \dots, N, \quad k = 1, 2$$

where the mean and covariance are given by:

$$\boldsymbol{\mu}_{g_{ki}} = e^{\gamma_{ki}} \Sigma_{k_i} \Sigma_{f_k}^{-1} \mathbf{f}_k$$

$$\Sigma_{g_{ki}} = e^{2\gamma_{ki}} (\Sigma_{f_k} - \Sigma_{k_i} \Sigma_{f_k}^{-1} \Sigma_{k_i}^T + \epsilon I)$$

This result follows from standard Gaussian conditioning, where $\mathbf{g}_k(\mathbf{s}_i)$ inherits the properties of the underlying process \mathbf{f}_k , but is adjusted according to the transformation $a_i f(t - \tau_i)$.

Note about the Matrices Σ_{f_k} and Σ_{k_i} :

The temporal covariance Σ_{f_k} is defined by:

$$\Sigma_{f_k}[m, n] = \exp \left\{ -\frac{\rho_k^2}{2} (t_m - t_n)^2 \right\}$$

where t_m and t_n are time indices, and ρ_k is the length-scale parameter for source k .

The correlation between the source k and site i over time is modeled by:

$$\Sigma_{k_i}[m, n] = \exp \left\{ -\frac{\rho_k^2}{2} (t_m - \tau_i - t_n)^2 \right\}$$

where τ_i represents the shift for site i .

On the Nugget ϵ :

If $\tau_i = 0$, the variance of the posterior would approach zero, making $\mathbf{g}_k(\mathbf{s}_i)$ deterministic rather than stochastic. To preserve the stochastic nature of $\mathbf{g}_k(\mathbf{s}_i)$ even when $\tau_i = 0$, we introduce a nugget term ϵ . This nugget is modeled as a small noise term, drawn from a normal distribution with mean zero and a small variance, i.e., $\epsilon \sim \mathcal{N}(0, b_\epsilon)$, where b_ϵ is a fixed parameter.

4.2. Full conditional of β_k

We assume the following model:

$$\gamma_k | \beta_k, \phi_k \stackrel{ind}{\sim} \mathcal{N}_N(\mathbf{X} \beta_k, \Sigma_{\gamma_k})$$

$$\Sigma_{\gamma_k}[i, j] = \exp \left(-\frac{\phi_k^2}{2} \|\mathbf{s}_i - \mathbf{s}_j\|^2 \right)$$

The marginal prior on β_k is:

$$\beta_k \stackrel{iid}{\sim} \mathcal{N}_P(\mathbf{0}, \mathbf{I}_p)$$

where P represents the number of covariate

Applying Bayes' theorem, the full conditional is proportional to:

$$\beta_k \mid \gamma_k, \phi_k \propto \mathcal{L}(\gamma_k \mid \beta_k, \phi_k) \pi(\beta_k)$$

Expanding the law of γ_k given the rest:

$$\mathcal{L}(\gamma_k \mid \beta_k, \phi_k) \propto \exp \left(-\frac{1}{2} (\gamma_k - \mathbf{X} \beta_k)^T \Sigma_{\gamma_k}^{-1} (\gamma_k - \mathbf{X} \beta_k) \right)$$

Expanding the marginal prior:

$$\pi(\beta_k) \propto \exp \left(-\frac{1}{2} \beta_k^T \mathbf{I}_p \beta_k \right)$$

Rearranging terms:

$$\begin{aligned} & \propto \exp \left[-\frac{1}{2} ((\gamma_k - \mathbf{X} \beta_k)^T \Sigma_{\gamma_k}^{-1} (\gamma_k - \mathbf{X} \beta_k) + \beta_k^T \beta_k) \right] \\ & \propto \exp \left[-\frac{1}{2} (-2 \gamma_k^T \Sigma_{\gamma_k}^{-1} \mathbf{X} \beta_k + \beta_k^T \mathbf{X}^T \Sigma_{\gamma_k}^{-1} \mathbf{X} \beta_k + \beta_k^T \beta_k) \right] \\ & \propto \exp \left[-\frac{1}{2} (-2 \gamma_k^T \Sigma_{\gamma_k}^{-1} \mathbf{X} \beta_k + \beta_k^T (\mathbf{X}^T \Sigma_{\gamma_k}^{-1} \mathbf{X} + \mathbf{I}_p) \beta_k) \right] \end{aligned}$$

Defining:

$$\begin{aligned} \mathbf{S} &= \mathbf{X}^T \Sigma_{\gamma_k}^{-1} \mathbf{X} + \mathbf{I}_p, \\ & \propto \exp \left[-\frac{1}{2} (-2 \gamma_k^T \Sigma_{\gamma_k}^{-1} \mathbf{X} \mathbf{S}^{-1} \mathbf{S} \beta_k + \beta_k^T \mathbf{S} \beta_k) \right] \end{aligned}$$

Defining:

$$\begin{aligned} \mathbf{m}^T &= \gamma_k^T \Sigma_{\gamma_k}^{-1} \mathbf{X} \mathbf{S}^{-1} \\ & \propto \exp \left[-\frac{1}{2} (-2 \mathbf{m}^T \mathbf{S} \beta_k + \beta_k^T \mathbf{S} \beta_k) \right] \end{aligned}$$

Now we look for a quadratic form, getting rid of terms which do not depend on β_k . We retrieve:

$$\propto \exp \left[-\frac{1}{2} (\beta_k - \mathbf{m})^T \mathbf{S} (\beta_k - \mathbf{m}) \right]$$

Thus, the full conditional follows in close form:

$$\beta_k \mid \gamma_k, \phi_k \stackrel{ind}{\sim} \mathcal{N}_P(\mathbf{m}, \mathbf{S}^{-1})$$

and the final Expression follows:

$$\beta_k \mid \gamma_k, \phi_k \stackrel{ind}{\sim} \mathcal{N}_P((\mathbf{X}^T \Sigma_{\gamma_k}^{-1} \mathbf{X} + \mathbf{I}_p)^{-1} \mathbf{X}^T \Sigma_{\gamma_k}^{-1} \gamma_k, (\mathbf{X}^T \Sigma_{\gamma_k}^{-1} \mathbf{X} + \mathbf{I}_p)^{-1}) \quad (15)$$

4.3. Full conditional of ϕ_k

Since in the Metropolis Hasting steps the proposal must be positive, we work with a transformed variable:

$$\theta_k = \log(\phi_k) \quad (16)$$

s.t $\phi_k = e^{\theta_k}$

The Gaussian proposal is then:

$$q(\theta^* | \theta^{\text{old}}) = \mathcal{N}(\theta^{\text{old}}, \varepsilon) \quad (17)$$

and the marginal prior transforms as:

$$\pi_{\theta_k}(\theta_k) = \pi_{\phi_k}(e^{\theta_k})e^{\theta_k} \quad (18)$$

Thus, the MH target function is:

$$\mathcal{L}(\gamma_k | \beta_k, e^{\theta_k}) \pi_{\phi_k}(e^{\theta_k}) e^{\theta_k} \quad (19)$$

4.4. Full conditional of τ

The MH target function is:

$$\mathcal{L}(\tau | \mathbf{g}_k(\mathbf{s}_i), \gamma_k, \mathbf{f}_k) \propto \prod_{i=1}^N \mathcal{L}(\mathbf{g}_k(\mathbf{s}_i) | \gamma_{ki}, \tau_i, \mathbf{f}_k) \pi(\tau) \quad (20)$$

4.5. Full conditional of ρ_k

Defining $\theta_k = \log(\rho_k)$, we get:

$$\mathcal{L}(\rho_k | \mathbf{g}_k(\mathbf{s}_i), \gamma_k, \tau, \mathbf{f}_k) \propto \prod_{i=1}^N \mathcal{L}(\mathbf{g}_k(\mathbf{s}_i) | \gamma_{ki}, \tau_i, \mathbf{f}_k) \mathcal{L}(\mathbf{f}_k | \rho_k) \pi_{\rho_k}(e^{\theta_k}) e^{\theta_k} \quad (21)$$

4.6. Full conditional of \mathbf{f}_k

It is trivial to show that it follows a closed-form Gaussian distribution [5]:

$$\mathcal{L}(\mathbf{f}_k | \mathbf{G}_k, \gamma_k, \tau, \rho_k) \propto \mathcal{L}(\mathbf{G}_k | \gamma_k, \tau, \mathbf{f}_k) \mathcal{L}(\mathbf{f}_k | \rho_k) \quad (22)$$

and so we get:

$$\mathbf{f}_k | \mathbf{G}_k, \gamma_k, \tau, \rho_k \sim \mathcal{N}_T(\boldsymbol{\mu}_{f_k\text{-post}}, \boldsymbol{\Sigma}_{f_k\text{-post}}) \quad (23)$$

$$\boldsymbol{\Sigma}_{f_k\text{-post}} = (\boldsymbol{\Sigma}_{f_k}^{-1} + \sum_{i=1}^n (\boldsymbol{\Sigma}_{k_i} \boldsymbol{\Sigma}_{f_k}^{-1})^T \boldsymbol{\Sigma}_{g_{ki}}^{-1} \boldsymbol{\Sigma}_{k_i} \boldsymbol{\Sigma}_{f_k}^{-1})^{-1} \quad (24)$$

$$\boldsymbol{\mu}_{f_k\text{-post}} = \left(\boldsymbol{\Sigma}_{f_k}^{-1} + \sum_{i=1}^n (\boldsymbol{\Sigma}_{k_i} \boldsymbol{\Sigma}_{f_k}^{-1})^T \boldsymbol{\Sigma}_{g_{ki}}^{-1} \boldsymbol{\Sigma}_{k_i} \boldsymbol{\Sigma}_{f_k}^{-1} \right)^{-1} \sum_{i=1}^n \mathbf{G}_k[i, :]^T \boldsymbol{\Sigma}_{g_{ki}}^{-1} \boldsymbol{\Sigma}_{k_i} \boldsymbol{\Sigma}_{f_k}^{-1} \quad (25)$$

This result follows from standard Gaussian conditioning, a more general case is shown in Vannucci's paper. To understand the updating formulas for $\boldsymbol{\mu}_{f\text{-post}}, \boldsymbol{\Sigma}_{f\text{-post}}$ we invite the reader to check the code on Github.

4.7. Full conditional of γ_k

The full conditional distribution of γ_k is not available in closed form. Moreover, Metropolis-Hastings methods tend to perform poorly in high-dimensional or highly correlated multivariate settings. To mitigate this issue, we propose an alternative approach: instead of sampling γ_k jointly, we marginalize over γ_k and work directly with its conditional components γ_{ki} .

To illustrate this, we start by considering the distribution of γ_k , that here will be denoted as γ to make the notation easier:

$$\gamma = \begin{pmatrix} \gamma_1 \\ \vdots \\ \gamma_N \end{pmatrix} = \begin{pmatrix} \gamma_{-N} \\ \gamma_N \end{pmatrix} \sim \mathcal{N}_N \left(\begin{pmatrix} X_{[-N,:]} \beta \\ X_{[N,:]} \beta \end{pmatrix}, \begin{bmatrix} \Sigma_{[-N,-N]} & \Sigma_{[-N,N]} \\ \Sigma_{[N,-N]} & \Sigma_{[N,N]} \end{bmatrix} \right).$$

From this, we derive the marginal distribution of γ_N given γ_{-N} :

$$\gamma_N \mid \gamma_{-N} \stackrel{ind}{\sim} \mathcal{N}(\mu_{N\text{marginal}}, \Sigma_{N\text{marginal}}).$$

We now focus on the conditional structure of γ_k . Specifically, we consider:

$$\gamma_{ki} \mid \gamma_k^{(-i)}, \tau, \mathbf{f}_k \propto \mathcal{L}(\mathbf{g}_k(s_i) \mid \gamma_{ki}, \tau, \mathbf{f}_k) \mathcal{L}(\gamma_{ki} \mid \gamma_k^{(-i)}). \quad (26)$$

Using the properties of the multivariate normal distribution, the conditional distribution of γ_{ki} given $\gamma_k^{(-i)}$ follows:

$$\gamma_{ki} \mid \gamma_k^{(-i)} \stackrel{ind}{\sim} \mathcal{N} \left(\mathbf{X}_i^\top \beta_k + \Sigma_{\gamma_k}[i, -i] (\Sigma_{\gamma_k}[-i, -i])^{-1} (\gamma_k^{(-i)} - \mathbf{X}_{[-i,:]} \beta_k), \right. \\ \left. \Sigma_{\gamma_k}[i, i] - \Sigma_{\gamma_k}[i, -i] (\Sigma_{\gamma_k}[-i, -i])^{-1} \Sigma_{\gamma_k}[-i, i] \right). \quad (27)$$

$$\Sigma_{\gamma_k} = \begin{bmatrix} \boxed{\Sigma_{[-i, -i]}} & \boxed{\Sigma_{[-i, i]}} \\ \boxed{\Sigma_{[i, -i]}} & \boxed{\Sigma_{[i, i]}} \end{bmatrix}$$

NB : Following the R language alike notation, $\Sigma_{\gamma_k}[-i, -i]$ indicates the matrix Σ_{γ_k} without the i -th row and the i -th column (Dim = (number of sites -1) x (number of sites -1)), $\Sigma_{\gamma_k}[i, -i]$ indicates only the i -th row of Σ_{γ_k} , without the only element in the i -th column of the matrix (Dim = 1 x (number of sites -1)), $\Sigma_{\gamma_k}[-i, i]$ indicates only the i -th column of Σ_{γ_k} , without the only element in the i -th row of the matrix (Dim = (number of sites -1) x 1), and $\Sigma_{\gamma_k}[i, i]$ is the element in the i -th row and i -th column of Σ_{γ_k} .

Thus, the Metropolis-Hastings target function is:

$$\mathcal{L}(\mathbf{g}_k(s_i) \mid \gamma_{ki}, \tau_i, \mathbf{f}_k) \mathcal{L}(\gamma_{ki} \mid \gamma_k^{(-i)}). \quad (28)$$

Computational Advantages:

This approach simplifies the sampling process by reducing the dimensionality of the sampling space. Instead of handling the full joint posterior of γ_k , we iteratively sample its components γ_{ki} . This improves computational efficiency and convergence properties.

5. MCMC Algorithm

Given the full conditionals derived in the previous section, we can now proceed to implement the following MCMC algorithm in order to sample from the joint posterior distribution, with a Gibbs sampler for parameters for which we have a known closed form, and Metropolis Hasting steps when the closed form is not available.

```

Initialize  $\gamma_k^{(0)}, f_k^{(0)}, \beta_k^{(0)}, \Phi_k^{(0)}, \tau^{(0)}, \rho_k^{(0)}$ 
for  $iter \in 1 : num.iterations$  do
    Update  $f_k^{(iter)} | \gamma_k^{(iter-1)}, \tau^{(iter-1)}, \rho_k^{(iter-1)}, G_k^{(iter-1)}$  (closed-form)
    Update  $\beta_k^{(iter)} | \gamma_k^{(iter-1)}, \Phi_k^{(iter-1)}$  (closed-form)
    Update  $\Phi_k^{(iter)} | \gamma_k^{(iter-1)}, \beta_k^{(iter)}, g_k^{(iter-1)}$ 
    for  $i \in 1 : N$  do
        Update  $\gamma_{ki}^{(iter)} | \gamma_{k(-i)}^{(iter)}, \tau_i^{(iter-1)}, f_k^{(iter)}, \beta_k^{(iter)}, \Phi_k^{(iter)}$  (MH on single marginal)
    end
    Update  $\tau^{(iter)} | \gamma_k^{(iter)}, f_k^{(iter)}, g_k^{(iter-1)}$ 
    Update  $\rho_k^{(iter)} | \gamma_k^{(iter)}, \tau^{(iter)}, f_k^{(iter)}, g_k^{(iter-1)}$ 
end

```

6. MCMC structure in Julia

```

1 function fit_model(sites, g, n_iter, theta0, hyperparam)
2     chain = Vector{Any}(undef, n_iter)
3     chain_f = Vector{Any}(undef, n_iter)
4     chain_g = Vector{Any}(undef, n_iter)
5     chain_z = Vector{Any}(undef, n_iter)
6     chain_beta = Vector{Any}(undef, n_iter)
7     chain_gamma = Vector{Any}(undef, n_iter)
8     chain_tau = Vector{Any}(undef, n_iter)
9     chain_rho = Vector{Any}(undef, n_iter)
10    chain_phi = Vector{Any}(undef, n_iter)
11    n_time = size(g, 2)
12    n = size(g, 1)
13
14    t = range(1, stop=n_time, length=n_time)
15    dist = euclid_dist(sites[:, 1], sites[:, 2], n)
16    X = sites[:, 3:6]
17
18    chain[1] = copy(theta0)
19
20    chain_beta[1] = theta0[:beta]
21    chain_tau[1] = theta0[:tau]
22    chain_gamma[1] = theta0[:gamma]
23    chain_rho[1] = theta0[:rho]
24    chain_phi[1] = theta0[:phi]
25    chain_f[1] = sample_f(g, chain[1], 1)
26    Sigma_f = sq_exp_kernel(t, chain[1][:rho], nugget = 1e-9)
27    Sigma_f_inv = inv(Sigma_f)
28
29    chain_g[1] = get_mu_g_matrix(g, chain_f[1], t, chain[1], Sigma_f_inv)
30    chain_z[1] = g - chain_g[1]
31

```

```

32     start = time()
33     Sigma_gamma = get_Sigma_gamma(dist, chain[1][:phi])
34
35     for iter in 2:n_iter
36         if (iter/n_iter*100) % 10 == 0.0
37             println("...", floor(Int, (iter / n_iter) * 100), "%")
38         end
39
40         curr = copy(chain[iter - 1])
41         f = sample_f(g, curr, 1)
42
43         curr[:beta] = sample_beta(curr, hyperparam, Sigma_gamma, X)
44         curr[:phi], Sigma_gamma = sample_phi(dist, X, curr, hyperparam)
45         curr[:gamma] = sample_gamma(t, g, f, curr, Sigma_f, Sigma_f_inv,
46             Sigma_gamma, X, hyperparam)
47         curr[:tau] = sample_tau(t, g, f, curr, hyperparam, Sigma_f,
48             Sigma_f_inv)
49         curr[:rho] = sample_rho(t, g, f, curr, hyperparam)
50
51         Sigma_f = sq_exp_kernel(t, curr[:rho], nugget = 1e-9)
52         Sigma_f_inv = inv(Sigma_f)
53
54         g_hat = get_mu_g_matrix(g, f, t, curr, Sigma_f_inv)
55         z = g - g_hat
56
57         chain_f[iter] = copy(f)
58         chain[iter] = copy(curr)
59         chain_g[iter] = copy(g_hat)
60         chain_z[iter] = copy(z)
61         chain_beta[iter] = copy(curr[:beta])
62         chain_gamma[iter] = copy(curr[:gamma])
63         chain_tau[iter] = copy(curr[:tau])
64         chain_rho[iter] = copy(curr[:rho])
65         chain_phi[iter] = copy(curr[:phi])
66     end
67
68     fine = time()
69     runtime = fine - start
70     println("Execution_Time:", runtime)
71
72     return Dict(
73         :chain => chain,
74         :chain_f => chain_f,
75         :chain_g => chain_g,
76         :chain_z => chain_z,
77         :chain_beta => chain_beta,
78         :chain_gamma => chain_gamma,
79         :chain_tau => chain_tau,
80         :chain_phi => chain_phi,
81         :chain_rho => chain_rho,
82         :runtime => runtime
83     )
84 end

```

7. Turing: Model Implementation in Julia

In this section, we present a complementary approach explored during the development of the MCMC. Specifically, we implemented a probabilistic model using Turing.jl to obtain preliminary results more efficiently, without manually constructing a custom MCMC algorithm.

7.1. Overview of Turing.jl

Turing.jl is a general-purpose probabilistic programming language (PPL) written entirely in Julia. It provides a streamlined approach to Bayesian inference, requiring the user to define only the probabilistic model—including priors and likelihood—within a function. Once the model is specified, Turing.jl automatically performs MCMC sampling through a dedicated command. This not only simplifies the implementation process but also enhances the readability and interpretability of the analysis.

7.2. Model Construction

The model follows a structure similar to the custom MCMC approach, iterating through the definition of priors and likelihoods up to the retrieval of the latent variables g . The final component of the Turing.jl model introduces an additional step for modeling the observed data y , as outlined below.

Observed data y is modeled as follows:

$$h_k \sim \text{Dirichlet}(\mathbf{1}) \quad (29)$$

$$\sigma_y \sim \text{InverseGamma}(3, 2) \quad (30)$$

$$\mu_{y,i,t,c} = \sum_{k=1}^K h_{k,c} g_{i,t,k} \quad (31)$$

$$y_{i,t,c} \sim \mathcal{N}(\mu_{y,i,t,c}, \sigma_y) \quad (32)$$

Here, y follows a normal distribution with mean μ_y and standard deviation σ_y , where μ_y is computed as a weighted sum of the latent variables g using the Dirichlet-distributed weights h .

8. Convergence Diagnostics and Posterior Inference

In this section, we begin by discussing the convergence diagnostics used to validate the MCMC chains, followed by detailed posterior inferences for the parameters \mathbf{g}_k , \mathbf{f}_k , γ_{ki} , and β_k . Finally, we revisit convergence diagnostics by examining the autocorrelation behavior.

(*Note:* all results concern source $K=2$ and because of computational issues, we run a MCMC simulation with $T=180$.)

8.1. Convergence Diagnostic

Model validation and performance metrics.

- *Trace plots:* The MCMC trajectories should eventually look like white noise without a discernible trend, indicating stationarity. Starting chains at different initial values is a practical way to confirm that all chains converge to the same posterior region.
- *Markov chain standard error:* We computed the chain-specific standard error for each parameter. As the number of iterations increases, the chains mix well, and the standard errors stabilize.

8.2. Posterior Inference: $\mathbf{g}_k(\mathbf{s}_i, t)$

Retrieving the simulated data. Our MCMC of 11,000 iterations successfully recovers the pollutant levels $\mathbf{g}_k(\mathbf{s}_i, t)$ across different sites and time points. Figure 4 illustrates how the posterior mean (and associated credible intervals) align closely with the simulated data.

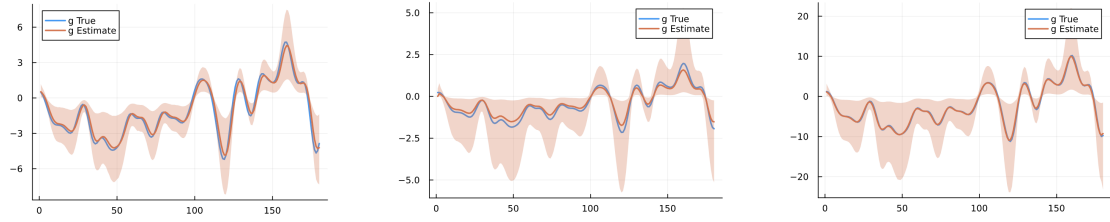


Figure 4: Retrieved $\mathbf{g}_k(\mathbf{s}_i, t)$ for sites $\{1, 2, 3\}$ with different scale of magnitude. Comparison with the "true" \mathbf{g}_k along with its 0.95 credible interval

Across all sites, the model captures the temporal trends and magnitude of \mathbf{g}_k . The slight variations around the true values reflect the model's uncertainty quantification.

8.3. Posterior Inference: \mathbf{f}_k

Recovering the second hierarchical level. We partially succeed in retrieving \mathbf{f}_k , subject to the normalization constraints. Up to a constant, we successfully recover the true function f_{true} (Figure 5). By standardizing both f_{true} and f_{estimate} , we observe a clear alignment between the two (Figure 6).

This agreement indicates that the Bayesian framework appropriately infers the underlying source composition while accounting for potential scaling ambiguities inherent in hierarchical models.

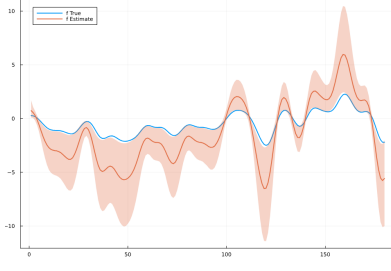


Figure 5: Retrieved original f

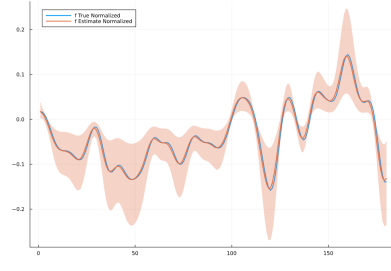


Figure 6: Retrieved adjusted f

8.4. Posterior Inference: γ_k

Normalization adjustment. Due to the normalization applied to \mathbf{f}_k , we adjust the retrieved γ ((Figure 7) accordingly. Specifically, we use:

$$\gamma_{\text{adjusted}} = \ln(\|\mathbf{f}\|) + \gamma,$$

which ensures consistency between the model’s inference and the generative process (Figure 8).

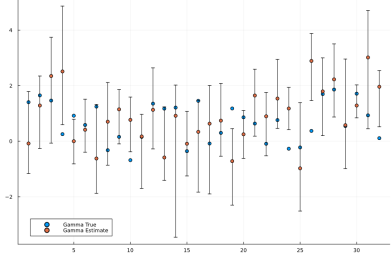


Figure 7: Retrieved γ_{ki}

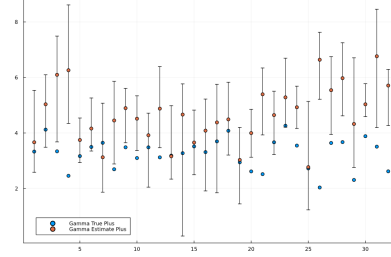


Figure 8: Adjusted γ_{ki}

The adjusted γ values align more closely with the physical interpretation of the model, correcting for the shift introduced by normalization.

8.5. Convergence Diagnostic: β_k

Trace plots and credible intervals. For the regression coefficients β_k , we have a closed-form solution in each MCMC iteration, leading to well-mixed “fat caterpillar” trace plots (Figures 9). This visual pattern suggests efficient sampling. Moreover, the recovered values of β_k lie within the 0.95 credible intervals of their true values (except for β_0), indicating that the model is capturing these parameters accurately (Figure 10).

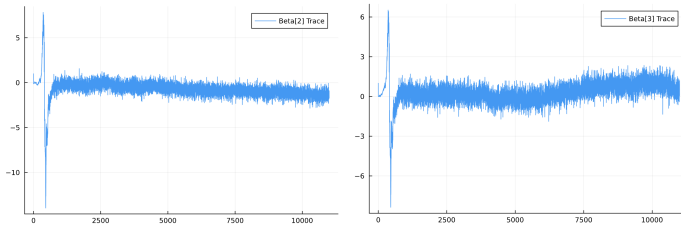


Figure 9: Traceplot of $\beta_k[1]$ and $\beta_k[2]$

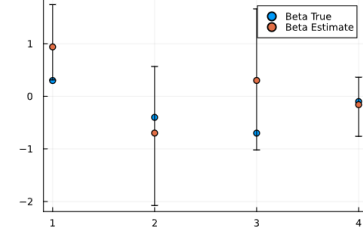


Figure 10: Estimated β_k ’s

The diagnostics confirm that the sampling scheme for β_k converges and yields estimates that reflect the true underlying parameter values.

8.6. Convergence Diagnostic: Autocorrelation

Assessing mixing quality. To further evaluate chain mixing, we examine the autocorrelation plots of key parameters. Ideally, these plots should show a decreasing pattern as the lag increases, confirming that each subsequent sample is becoming less dependent on its predecessor.

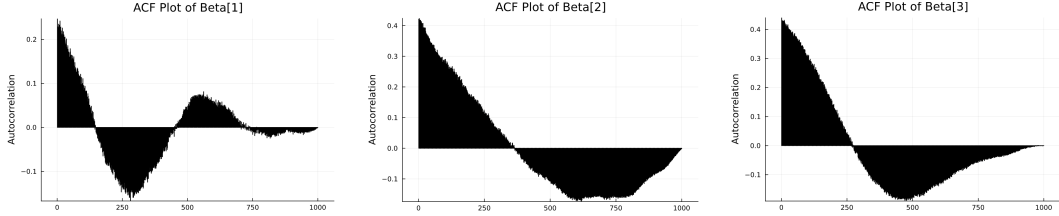


Figure 11: Autocorrelation of $\beta_0, \beta_1, \beta_2$

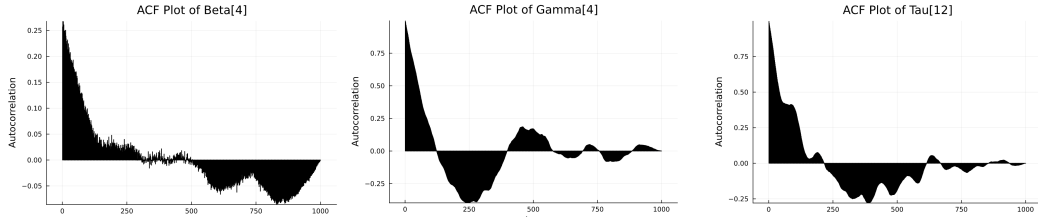


Figure 12: Autocorrelation of $\beta_3, \gamma_4, \tau_{23}$

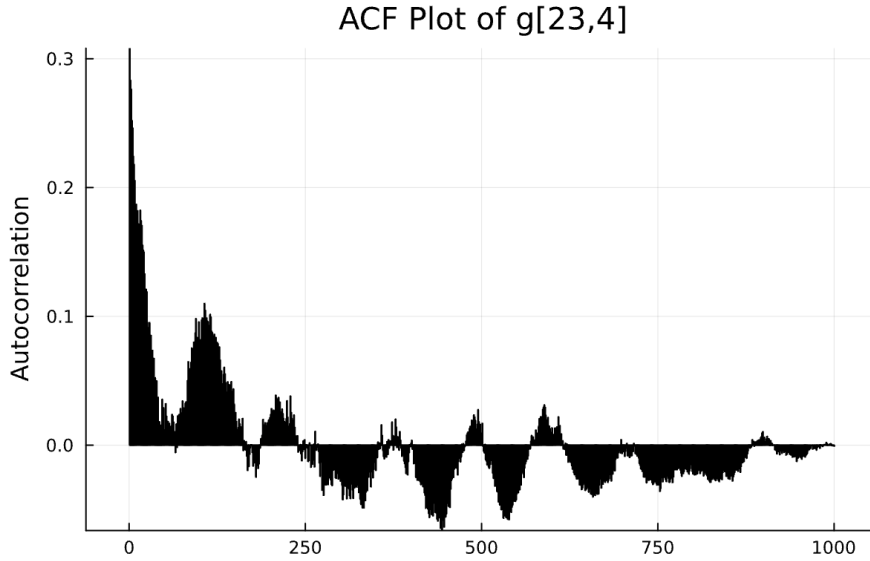


Figure 13: Autocorrelation of $g_{23,4}$

In practice, we apply thinning since the autocorrelation is high, in order to reduce bias and improve effective sample size. Here, the autocorrelation decreases rapidly with increasing lag, confirming that the chains mix efficiently.

Note: to see more convergence diagnostic insights, along with few issues that often appear with infinite latent factors model, see Appendix A and B.

9. Conclusions

In this study, we have successfully implemented a comprehensive Bayesian hierarchical framework to address complex spatial-temporal modeling challenges. The key findings of our work can be summarized as follows:

- **Implementation Pipeline:** Starting from Vannucci’s model for random phase amplitude Gaussian processes in R, we developed a complex data generation pipeline in Julia. This pipeline is capable of producing synthetic data that closely mimics real observations by leveraging a rigorous Bayesian hierarchical model.
- **Customized MCMC Algorithm:** After extensive derivations of full conditional distributions—including efforts to identify closed-form solutions—we devised a customized Markov Chain Monte Carlo (MCMC) algorithm.
- **Parameter Retrieval:** Our approach successfully retrieves the target functions g_k that represent the first level of the model. In addition, key second-level parameters, such as β_k and f , are accurately estimated. These results demonstrate the model’s capability to capture both the primary and secondary structures of the data.

Overall, the methodology proves to be effective in efficiently estimating parameters within intricate spatial-temporal frameworks, highlighting its potential for broader applications in environmental and other complex hierarchical models.

10. Possible Improvements

10.1. Extending MCMC to the Variable y

The next step is to extend the MCMC framework to incorporate the observed pollutant levels y , incorporating the positivity constraints.

10.2. Shrinkage in the RPAGP Model

In the initial formulation of the model, the number of pollutant sources k is fixed. However, in the context of the proposed shrinkage method [3], we introduce η_k , a parameter that controls the degree of shrinkage for each pollutant source. Specifically, the formulation becomes:

$$g_k(\mathbf{s}_i, t) = \gamma_{k,i} \cdot f_k(t - \tau_i)$$

for $k = 1, \dots, K$, $i = 1, \dots, N$, and $t = 1, \dots, T$, where:

- $f_k \stackrel{\text{ind}}{\sim} \mathcal{N}_T^+(\mathbf{0}, \frac{1}{\eta_k} \Sigma_f)$, representing a Gaussian process with positive values, where the covariance structure is governed by $\Sigma_f(t', t) = \exp\left(-\frac{\rho^2}{2}(t' - t)^2\right)$, the squared exponential kernel.
- The parameter η_k is defined as:

$$\eta_k = \prod_{l=1}^k \delta_l$$

where $\delta_1 \sim \text{Gamma}(a_1, 1)$ and $\delta_l \sim \text{Gamma}(a_2, 1)$ for $l \geq 2$, introducing a hierarchical shrinkage effect.

- The length-scale parameter ρ governs the smoothness of the source functions and is assumed to follow a Gamma distribution:

$$\rho \stackrel{\text{iid}}{\sim} \text{Gamma}(a_\rho, b_\rho)$$

where a_ρ and b_ρ are hyperparameters selected based on a plausible range of values.

The introduction of η_k ensures that only those pollutant sources whose variance is sufficiently large are kept in the model, while redundant or insignificant sources are pruned away. This mechanism enables the model to adaptively select the appropriate number of sources.

10.3. Testing the Shrinkage Model on Real Data from ARPA

The next phase involves applying the model to real-world pollution data collected by ARPA in Lombardy. This data provides actual measurements of PM10 across various locations and times. Testing the model on this real data will allow us to assess its ability to identify the relevant pollutant sources in a practical, real-world setting. The real-world validation is crucial for confirming that the dynamic source selection and shrinkage mechanism can adapt to complex, real-world pollution scenarios.

References

- [1] Oliver Baerenbold et al. *A dependent Bayesian Dirichlet process model for source apportionment of particle number size distribution*. 2022.
- [2] Sudipto Banerjee, Bradley P Carlin, and Alan E Gelfand. *Hierarchical modeling and analysis for spatial data*. Chapman and Hall/CRC, 2003.
- [3] A. Bhattacharya and D. B. Dunson. *Sparse Bayesian infinite factor models*. 2011.
- [4] Jenna R. Krall and Howard H. Chang. *Statistical methods for source apportionment*. 2019.
- [5] Dustin Pluta and Marina Vannucci. *Improved data quality and statistical power of trial-level event-related potentials with Bayesian random-shift Gaussian processes*. 2024.

Appendix A: Trace plots, nonidentifiability issues for latent parameters

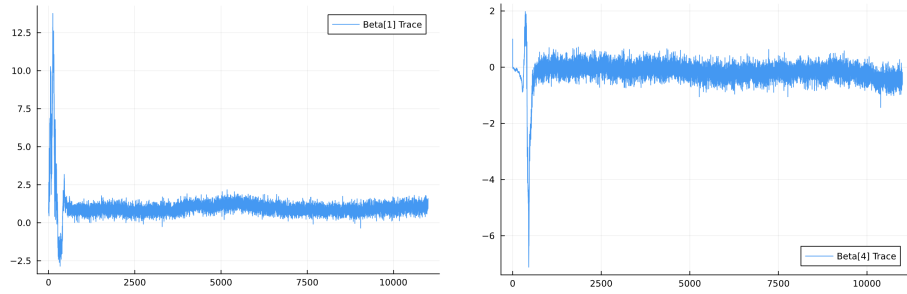


Figure 14: Trace plots of β_1, β_2

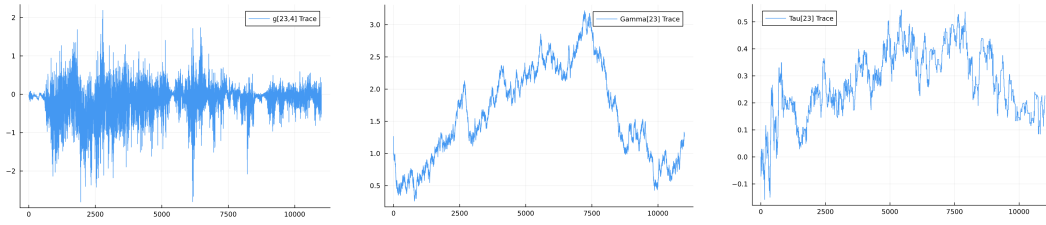


Figure 15: Trace plots of $g_{23,4}, \gamma_{23}, \tau_{23}$

Appendix B: Autocorrelation plots

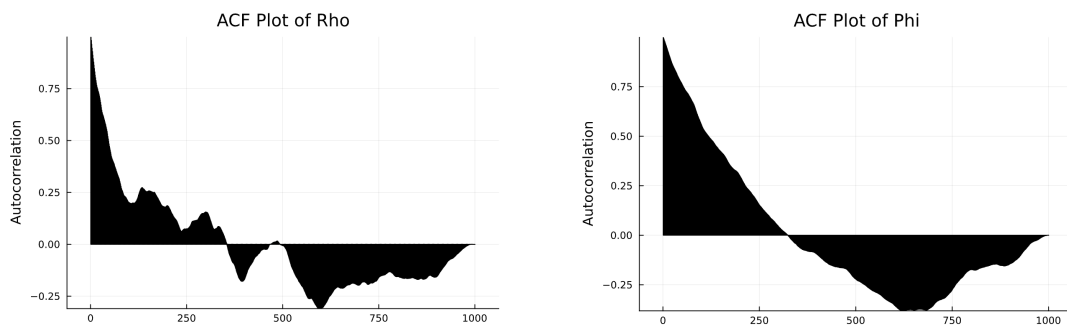


Figure 16: Autocorrelation of ρ_k and ϕ_k