

Probabilistic Predictive Principal Component Analysis for Spatially-Misaligned and High-Dimensional Air Pollution Data with Missing Observations

Phuong T. Vu¹, Timothy V. Larson² and Adam A. Szpiro¹

¹Department of Biostatistics, University of Washington

²Department of Civil & Environmental Engineering, University of
Washington

Supplemental Materials

Corresponding author: Phuong T. Vu, Department of Biostatistics, University of Washington, F-600, Health Sciences Building, Box 357232, Seattle, WA 98195-7232.
Email: phuongvu@uw.edu.

1 The ProPrPCA-Krige model and algorithm

1.1 The model

The ProPrPCA-Krige assumes that $\mathbf{X} = \sum_{l=1}^q (\mathbf{u}_l \mathbf{v}_l^\top + \mathbf{E}_l)$ and $\mathbf{u}_l = \mathbf{R}\boldsymbol{\beta}_l + \boldsymbol{\eta}_l$, where $\mathbf{u}_l \in \mathbb{R}^n$, $\mathbf{v}_l \in \mathbb{R}^p$, $\boldsymbol{\beta}_l \in \mathbb{R}^k$, $E_{ij} \sim \mathcal{N}(0, \gamma^2)$, and $\boldsymbol{\eta}_l \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}(\boldsymbol{\xi}_l))$. For notation simplification, we are going to ignore the subscript l for the following mathematical derivation. The parameter estimation is the same for all PC. For each PC, the model becomes:

$$\begin{aligned}\mathbf{X} &= \mathbf{u}\mathbf{v}^\top + \mathbf{E}, \\ \mathbf{u} &= \mathbf{R}\boldsymbol{\beta} + \boldsymbol{\eta},\end{aligned}$$

Denote $\Theta = \{\mathbf{v}, \boldsymbol{\beta}, \gamma^2, \boldsymbol{\xi}\}$ as the collection of the model parameters. To solve for Θ , we first rewrite the model in the conventional vectorized version. Denote $\mathbf{W} \in \mathbb{R}^N$, for $N = np$, as the vectorized version of \mathbf{X} , i.e.

$$\mathbf{W} = \begin{bmatrix} \mathbf{W}_1 \\ \vdots \\ \mathbf{W}_n \end{bmatrix} = \begin{bmatrix} \mathbf{X}_{1:} \\ \vdots \\ \mathbf{X}_{n:} \end{bmatrix},$$

where \mathbf{X}_i is the i -th row of \mathbf{X} . The model assumes that $\mathbf{W}_i = \mathbf{X}_i = u_i \mathbf{v} + \boldsymbol{\epsilon}_i$, for $i = 1, \dots, n$. Here \mathbf{v} represents the transformation from the latent variable space to the multi-pollutant exposure space, and $\boldsymbol{\epsilon}_i$'s are i.i.d. Gaussian noises distributed with mean zero and variance γ^2 . The full model can then be written as $\mathbf{W} = \mathbf{V}\mathbf{u} + \boldsymbol{\epsilon}$, where $\mathbf{V} = \mathbf{I}_n \otimes \mathbf{v}$ and \otimes denotes the Kronecker product. The model also assumes that the latent variables are normally distributed with a spatial mean model and covariance structure. That is, $\mathbf{u} \sim \mathcal{N}(\mathbf{R}\boldsymbol{\beta}, \boldsymbol{\Sigma}(\boldsymbol{\xi}))$. In this paper, we assume $\boldsymbol{\Sigma}(\boldsymbol{\xi})$ has an exponential structure with no nugget effect. For identifiability, we assume that $\|\mathbf{v}\|_2 = 1$. When every element of \mathbf{X} is observed, we have the following hierarchical model:

$$\begin{aligned}\mathbf{W} \mid \mathbf{u} &\sim \mathcal{N}(\mathbf{V}\mathbf{u}, \gamma^2 \mathbf{I}_N), \\ \mathbf{u} &\sim \mathcal{N}(\mathbf{R}\boldsymbol{\beta}, \boldsymbol{\Sigma}(\boldsymbol{\xi})).\end{aligned}$$

In practice, not all pollutants are measured at every monitoring location. Denote $\mathbf{W}_o \in \mathbb{R}^{N_o}$ as the collection of all observed elements of \mathbf{W} , and $\mathbf{W}_m \in \mathbb{R}^{N_m}$ as the collection of all missing entries, where $N_o + N_m = N$. Algebraically, there exists a linear transformation \mathbf{G} such that

$$\mathbf{G}\mathbf{W} = \begin{bmatrix} \mathbf{G}_o \\ \mathbf{G}_m \end{bmatrix} \mathbf{W} = \begin{bmatrix} \mathbf{W}_o \\ \mathbf{W}_m \end{bmatrix},$$

where $\mathbf{G}_o \in \mathbb{R}^{N_o \times N}$ and $\mathbf{G}_m \in \mathbb{R}^{N_m \times N}$. Each row and column of \mathbf{G} contains exactly one element of one and $(N - 1)$ zeros. Thus by construction, \mathbf{G}_o and \mathbf{G}_m are both full row rank, as well as $\mathbf{G}_o \mathbf{G}_o^\top = \mathbf{I}_{N_o}$ and $\mathbf{G}_m \mathbf{G}_m^\top = \mathbf{I}_{N_m}$. The hierarchical model for the observed

elements become:

$$\begin{aligned}\mathbf{W}_o | \mathbf{u} &\sim \mathcal{N}(\mathbf{G}_o \mathbf{V} \mathbf{u}, \gamma^2 \mathbf{I}_{N_o}), \\ \mathbf{u} &\sim \mathcal{N}(\mathbf{R} \boldsymbol{\beta}, \boldsymbol{\Sigma}(\boldsymbol{\xi})).\end{aligned}$$

1.2 Estimation of model parameters when monitoring data is complete

Our approach to estimate the model parameters is similar to the EM algorithm employed by Tipping and Bishop (1999). We consider the latent variable \mathbf{u} to be the ‘‘missing’’ portion. Thus the ‘‘complete’’ data consists of the observed data \mathbf{W} , and the latent variable \mathbf{u} . The goal is then to maximize the joint likelihood of (\mathbf{W}, \mathbf{u}) , i.e. $\mathcal{L} = f(\mathbf{W}, \mathbf{u}) = f(\mathbf{W}|\mathbf{u})f(\mathbf{u})$. The ‘‘complete’’ log-likelihood, up to a constant, is:

$$\ell_c = -\frac{N}{2} \log \gamma^2 - \frac{1}{2} \log |\boldsymbol{\Sigma}| - \frac{1}{2\gamma^2} (\mathbf{W} - \mathbf{V} \mathbf{u})^\top (\mathbf{W} - \mathbf{V} \mathbf{u}) - \frac{1}{2} (\mathbf{u} - \mathbf{R} \boldsymbol{\beta})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{u} - \mathbf{R} \boldsymbol{\beta}).$$

Because this log-likelihood involves \mathbf{u} which is unobserved, in each E step, we find the expectation of ℓ_c with respect to the conditional distribution of $\mathbf{u}|\mathbf{W}$. We first derive this distribution as follows:

$$\begin{aligned}f(\mathbf{u}|\mathbf{W}) &\propto \exp \left[-\frac{1}{2\gamma^2} (\mathbf{W} - \mathbf{V} \mathbf{u})^\top (\mathbf{W} - \mathbf{V} \mathbf{u}) \right] \times \exp \left[-\frac{1}{2} (\mathbf{u} - \mathbf{R} \boldsymbol{\beta})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{u} - \mathbf{R} \boldsymbol{\beta}) \right] \\ &\propto \exp \left\{ -\frac{1}{2} \left[\mathbf{u}^\top \left(\frac{1}{\gamma^2} \mathbf{V}^\top \mathbf{V} + \boldsymbol{\Sigma}^{-1} \right) \mathbf{u} - \mathbf{u}^\top \left(\frac{1}{\gamma^2} \mathbf{V}^\top \mathbf{W} + \boldsymbol{\Sigma}^{-1} \mathbf{R} \boldsymbol{\beta} \right) \right. \right. \\ &\quad \left. \left. - \left(\frac{1}{\gamma^2} \mathbf{W}^\top \mathbf{V} + \boldsymbol{\beta}^\top \mathbf{R}^\top \boldsymbol{\Sigma}^{-1} \right) \mathbf{u} \right] \right\}.\end{aligned}$$

Thus the distribution of $\mathbf{u}|\mathbf{W}$ is $\mathcal{N}(\mathbf{M}, \mathbf{S})$, where:

$$\begin{aligned}\mathbf{S} &= \left(\frac{1}{\gamma^2} \mathbf{V}^\top \mathbf{V} + \boldsymbol{\Sigma}^{-1} \right)^{-1}, \\ \mathbf{M} &= \mathbf{S} \left(\frac{1}{\gamma^2} \mathbf{V}^\top \mathbf{W} + \boldsymbol{\Sigma}^{-1} \mathbf{R} \boldsymbol{\beta} \right).\end{aligned}$$

We can further simplify these expressions by noticing that $\mathbf{V}^\top \mathbf{V} = \mathbf{I}_n$ and $\mathbf{V}^\top \mathbf{W} = \mathbf{X} \mathbf{v}$, using properties of Kronecker products. The conditional covariance and mean become

$$\begin{aligned}\mathbf{S} &= \left(\frac{1}{\gamma^2} \mathbf{I}_n + \boldsymbol{\Sigma}^{-1} \right)^{-1}, \\ \mathbf{M} &= \mathbf{S} \left(\frac{1}{\gamma^2} \mathbf{X} \mathbf{v} + \boldsymbol{\Sigma}^{-1} \mathbf{R} \boldsymbol{\beta} \right).\end{aligned}$$

Given a current estimate $\tilde{\Theta}$, the expectation of ℓ_c with respect to $\mathbf{u}|\mathbf{W}$ is:

$$E \left[\ell_c \mid \mathbf{W}, \tilde{\Theta} \right] = -\frac{N}{2} \log \gamma^2 - \frac{1}{2} \log |\Sigma| - \frac{1}{2\gamma^2} E \left[(\mathbf{W} - \mathbf{V}\mathbf{u})^\top (\mathbf{W} - \mathbf{V}\mathbf{u}) \mid \mathbf{W}, \tilde{\Theta} \right] - \frac{1}{2} E \left[(\mathbf{u} - \mathbf{R}\boldsymbol{\beta})^\top \Sigma^{-1} (\mathbf{u} - \mathbf{R}\boldsymbol{\beta}) \mid \mathbf{W}, \tilde{\Theta} \right] \quad (1)$$

The conditional distribution $\mathbf{u}|\mathbf{W}, \tilde{\Theta}$ is $\mathcal{N}(\tilde{\mathbf{M}}, \tilde{\mathbf{S}})$, where $\tilde{\mathbf{M}} = \mathbf{M}(\tilde{\Theta})$ and $\tilde{\mathbf{S}} = \mathbf{S}(\tilde{\Theta})$. This implies that

$$\begin{aligned} \mathbf{V}\mathbf{u} - \mathbf{W} &\sim \mathcal{N} \left(\mathbf{V}\tilde{\mathbf{M}} - \mathbf{W}, \mathbf{V}\tilde{\mathbf{S}}\mathbf{V}^\top \right), \\ \mathbf{u} - \mathbf{R}\boldsymbol{\beta} &\sim \mathcal{N} \left(\tilde{\mathbf{M}} - \mathbf{R}\boldsymbol{\beta}, \tilde{\mathbf{S}} \right). \end{aligned}$$

Thus the first expectation term of (1) is

$$\begin{aligned} E \left[(\mathbf{W} - \mathbf{V}\mathbf{u})^\top (\mathbf{W} - \mathbf{V}\mathbf{u}) \mid \mathbf{W}, \tilde{\Theta} \right] &= \text{Tr} \left(\mathbf{V}\tilde{\mathbf{S}}\mathbf{V}^\top \right) + (\mathbf{V}\tilde{\mathbf{M}} - \mathbf{W})^\top (\mathbf{V}\tilde{\mathbf{M}} - \mathbf{W}) \\ &= \text{Tr}(\tilde{\mathbf{S}}) + (\mathbf{V}\tilde{\mathbf{M}} - \mathbf{W})^\top (\mathbf{V}\tilde{\mathbf{M}} - \mathbf{W}). \end{aligned}$$

The second expectation term of (1) is

$$E \left[(\mathbf{u} - \mathbf{R}\boldsymbol{\beta})^\top \Sigma^{-1} (\mathbf{u} - \mathbf{R}\boldsymbol{\beta}) \mid \mathbf{W}, \tilde{\Theta} \right] = \text{Tr} \left(\Sigma^{-1} \tilde{\mathbf{S}} \right) + (\tilde{\mathbf{M}} - \mathbf{R}\boldsymbol{\beta})^\top \Sigma^{-1} (\tilde{\mathbf{M}} - \mathbf{R}\boldsymbol{\beta})$$

Hence (1) can be simplified as

$$\begin{aligned} E \left[\ell_c \mid \mathbf{W}, \tilde{\Theta} \right] &= -\frac{N}{2} \log \gamma^2 - \frac{1}{2} \log |\Sigma| - \frac{1}{2\gamma^2} \text{Tr}(\tilde{\mathbf{S}}) - \frac{1}{2\gamma^2} (\mathbf{V}\tilde{\mathbf{M}} - \mathbf{W})^\top (\mathbf{V}\tilde{\mathbf{M}} - \mathbf{W}) \\ &\quad - \frac{1}{2} \text{Tr} \left(\Sigma^{-1} \tilde{\mathbf{S}} \right) - \frac{1}{2} (\tilde{\mathbf{M}} - \mathbf{R}\boldsymbol{\beta})^\top \Sigma^{-1} (\tilde{\mathbf{M}} - \mathbf{R}\boldsymbol{\beta}). \end{aligned}$$

To solve for \mathbf{v} , we effectively maximize $\{-\frac{1}{2\gamma^2} (\mathbf{V}\tilde{\mathbf{M}} - \mathbf{W})^\top (\mathbf{V}\tilde{\mathbf{M}} - \mathbf{W})\}$, which can be rewrite as follows:

$$\begin{aligned} &-\frac{1}{2\gamma^2} (\mathbf{V}\tilde{\mathbf{M}} - \mathbf{W})^\top (\mathbf{V}\tilde{\mathbf{M}} - \mathbf{W}) \\ &= -\frac{1}{2\gamma^2} \left\| \mathbf{W} - \mathbf{V}\tilde{\mathbf{M}} \right\|_2^2 = -\frac{1}{2\gamma^2} \left\| \begin{bmatrix} \mathbf{X}_{1:} \\ \mathbf{X}_{2:} \\ \vdots \\ \mathbf{X}_{n:} \end{bmatrix} - \begin{bmatrix} \mathbf{v} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{v} & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{v} \end{bmatrix} \begin{bmatrix} \tilde{M}_1 \\ \tilde{M}_2 \\ \vdots \\ \tilde{M}_n \end{bmatrix} \right\|_2^2 \\ &= -\frac{1}{2\gamma^2} \sum_{i=1}^n \left\| \mathbf{X}_{i\cdot} - \tilde{M}_i \mathbf{v} \right\|_2^2 = -\frac{1}{2\gamma^2} \sum_{i=1}^n \sum_{j=1}^p \left(X_{ij} - \tilde{M}_i v_j \right)^2. \end{aligned}$$

Differentiate this expression with respect to each v_j , we get

$$\check{v}_j = \frac{\sum_{i=1}^n X_{ij} \tilde{M}_i}{\sum_{i=1}^n \tilde{M}_i^2} = \frac{\sum_{i=1}^n X_{ij} \tilde{M}_i}{\|\tilde{\mathbf{M}}\|_2^2}.$$

Thus, the solution for \mathbf{v} can be written in closed-form as

$$\hat{\mathbf{v}} = \frac{\check{\mathbf{v}}}{\|\check{\mathbf{v}}\|_2}, \quad \text{where } \check{\mathbf{v}} = \frac{\mathbf{X}^\top \tilde{\mathbf{M}}}{\|\tilde{\mathbf{M}}\|_2^2}$$

To solve for γ^2 , we maximize $\left\{ -\frac{N}{2} \log \gamma^2 - \frac{1}{2\gamma^2} \text{Tr}(\tilde{\mathbf{S}}) - \frac{1}{2\gamma^2} (\mathbf{V} \tilde{\mathbf{M}} - \mathbf{W})^\top (\mathbf{V} \tilde{\mathbf{M}} - \mathbf{W}) \right\}$. The closed-form solution for γ^2 is simply

$$\hat{\gamma}^2 = \frac{1}{N} [\text{Tr}(\tilde{\mathbf{S}}) + \|\mathbf{V} \tilde{\mathbf{M}} - \mathbf{W}\|_2^2].$$

The solution for $\boldsymbol{\beta}$ by maximizing $\left\{ -\frac{1}{2} (\tilde{\mathbf{M}} - \mathbf{R}\boldsymbol{\beta})^\top \boldsymbol{\Sigma}^{-1} (\tilde{\mathbf{M}} - \mathbf{R}\boldsymbol{\beta}) \right\}$ is

$$\hat{\boldsymbol{\beta}} = (\mathbf{R}^\top \boldsymbol{\Sigma}^{-1} \mathbf{R})^{-1} \mathbf{R}^\top \boldsymbol{\Sigma}^{-1} \tilde{\mathbf{M}}.$$

Finally, to solve for $\boldsymbol{\xi}$, we maximize $\left\{ -\frac{1}{2} \log |\boldsymbol{\Sigma}| - \frac{1}{2} \text{Tr}(\boldsymbol{\Sigma}^{-1} \tilde{\mathbf{S}}) - \frac{1}{2} (\tilde{\mathbf{M}} - \mathbf{R}\boldsymbol{\beta})^\top \boldsymbol{\Sigma}^{-1} (\tilde{\mathbf{M}} - \mathbf{R}\boldsymbol{\beta}) \right\}$ numerically, where $\boldsymbol{\Sigma}$ is a function of $\boldsymbol{\xi}$. In this paper, we adopt the exponential covariance structure for $\boldsymbol{\Sigma}$. For identifiability, we assume that $\boldsymbol{\Sigma}$ has no nugget effect.

Thus, parameter estimation of ProPrPCA-Krige with complete monitoring data can be summarized as:

Algorithm: ProPrPCA-Krige with complete monitoring data

Input \mathbf{X} , \mathbf{R} , q , and t_{max}
for l in $\{1, \dots, q\}$ **do**
 $\mathbf{X}_l \leftarrow \mathbf{X}_{l-1} - \hat{\mathbf{u}}_{l-1} \hat{\mathbf{v}}_{l-1}^\top$ where $\mathbf{X}_0 = \mathbf{X}$, $\hat{\mathbf{u}}_0 = \mathbf{0}$, and $\hat{\mathbf{v}}_0 = \mathbf{0}$
Initialize $\mathbf{v}_l^{(0)}$, $(\gamma_l^{(0)})^2$, $\beta_l^{(0)}$, $\xi_l^{(0)}$, and $t = 1$
 $\Sigma_l^{(0)} \leftarrow \Sigma(\xi_l^{(0)})$
while not converged **or** $t < t_{max}$ **do**
 $\tilde{\mathbf{S}}_l \leftarrow \left[(\gamma_l^{(t)})^{-2} \mathbf{I}_n + (\Sigma_l^{(t)})^{-1} \right]^{-1}$
 $\tilde{\mathbf{M}}_l \leftarrow \tilde{\mathbf{S}}_l \left[(\gamma_l^{(t)})^{-2} \mathbf{X}_l \mathbf{v}_l^{(t)} + (\Sigma_l^{(t)})^{-1} \mathbf{R} \beta_l^{(t)} \right]$
 $\tilde{\mathbf{v}}_l \leftarrow \mathbf{X}_l^\top \tilde{\mathbf{M}}_l / \|\tilde{\mathbf{M}}_l\|_2^2$
 $\mathbf{v}_l^{(t+1)} \leftarrow \tilde{\mathbf{v}}_l / \|\tilde{\mathbf{v}}_l\|_2$
 $(\gamma_l^{(t+1)})^2 \leftarrow (np)^{-1} \left[\text{Tr}(\tilde{\mathbf{S}}_l) + \|(\mathbf{I}_n \otimes \mathbf{v}_l^{(t+1)}) \tilde{\mathbf{M}}_l - \text{vec}(\mathbf{X}_l)\|_2^2 \right]$
 $\xi_l^{(t+1)} \leftarrow \arg \max_{\xi_l} \left\{ -\log |\Sigma_l| - \text{Tr} \left(\Sigma_l^{-1} \tilde{\mathbf{S}}_l \right) - (\tilde{\mathbf{M}}_l - \mathbf{R} \beta_l^{(t)})^\top \Sigma_l^{-1} (\tilde{\mathbf{M}}_l - \mathbf{R} \beta_l^{(t)}) \right\}$

 where $\Sigma_l = \Sigma(\xi_l)$
 $\beta_l^{(t+1)} \leftarrow \left(\mathbf{R}^\top \hat{\Sigma}(\xi_l^{(t+1)})^{-1} \mathbf{R} \right)^{-1} \mathbf{R}^\top \hat{\Sigma}(\xi_l^{(t+1)})^{-1} \tilde{\mathbf{M}}_l$
 $t \leftarrow t + 1$
end while
 $\hat{\mathbf{v}}_l \leftarrow \mathbf{v}_l^{(t)}$, $\hat{\gamma}_l^2 \leftarrow (\gamma_l^{(t)})^2$, $\hat{\beta}_l \leftarrow \beta_l^{(t)}$, $\hat{\xi}_l \leftarrow \xi_l^{(t)}$
 $\hat{\mathbf{u}}_l = \mathbf{X}_l \hat{\mathbf{v}}_l$
end for
Output $\{\hat{\mathbf{v}}_1, \dots, \hat{\mathbf{v}}_q\}$, $\{\hat{\mathbf{u}}_1, \dots, \hat{\mathbf{u}}_q\}$, $\{\hat{\beta}_1, \dots, \hat{\beta}_q\}$, $\{\hat{\gamma}_1^2, \dots, \hat{\gamma}_q^2\}$, $\{\hat{\xi}_1, \dots, \hat{\xi}_q\}$

1.3 Parameter estimation and model-based imputation with missing monitoring data

The hierarchical model in the case of missing monitoring data can be written as

$$\begin{aligned} \mathbf{W}_o \mid \mathbf{u} &\sim \mathcal{N}(\mathbf{G}_o \mathbf{V} \mathbf{u}, \gamma^2 \mathbf{I}_{N_o}), \\ \mathbf{u} &\sim \mathcal{N}(\mathbf{R} \boldsymbol{\beta}, \Sigma(\boldsymbol{\xi})). \end{aligned}$$

Thus the ‘‘complete’’ log-likelihood becomes

$$\begin{aligned} \ell_c &= -\frac{N_o}{2} \log \gamma^2 - \frac{1}{2} \log |\Sigma| - \frac{1}{2\gamma^2} (\mathbf{W}_o - \mathbf{G}_o \mathbf{V} \mathbf{u})^\top (\mathbf{W}_o - \mathbf{G}_o \mathbf{V} \mathbf{u}) \\ &\quad - \frac{1}{2} (\mathbf{u} - \mathbf{R} \boldsymbol{\beta})^\top \Sigma^{-1} (\mathbf{u} - \mathbf{R} \boldsymbol{\beta}). \end{aligned}$$

Similar to the case with complete data, we first derive the conditional distribution of

$\mathbf{u}|\mathbf{W}_o$ as follows

$$\begin{aligned} f(\mathbf{u}|\mathbf{W}_o) &\propto \exp \left[-\frac{1}{2\gamma^2} (\mathbf{W}_o - \mathbf{G}_o \mathbf{V} \mathbf{u})^\top (\mathbf{W}_o - \mathbf{G}_o \mathbf{V} \mathbf{u}) \right] \times \exp \left[-\frac{1}{2} (\mathbf{u} - \mathbf{R}\boldsymbol{\beta})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{u} - \mathbf{R}\boldsymbol{\beta}) \right] \\ &\propto \exp \left\{ -\frac{1}{2} \left[\mathbf{u}^\top \left(\frac{1}{\gamma^2} \mathbf{V}^\top \mathbf{G}_o^\top \mathbf{G}_o \mathbf{V} + \boldsymbol{\Sigma}^{-1} \right) \mathbf{u} - \mathbf{u}^\top \left(\frac{1}{\gamma^2} \mathbf{V}^\top \mathbf{G}_o^\top \mathbf{W}_o + \boldsymbol{\Sigma}^{-1} \mathbf{R}\boldsymbol{\beta} \right) \right. \right. \\ &\quad \left. \left. - \left(\frac{1}{\gamma^2} \mathbf{W}_o^\top \mathbf{G}_o \mathbf{V} + \boldsymbol{\beta}^\top \mathbf{R}^\top \boldsymbol{\Sigma}^{-1} \right) \mathbf{u} \right] \right\}. \end{aligned}$$

Thus the distribution of $\mathbf{u}|\mathbf{W}_o$ is $\mathcal{N}(\mathbf{M}, \mathbf{S})$, where:

$$\begin{aligned} \mathbf{S} &= \left(\frac{1}{\gamma^2} \mathbf{V}^\top \mathbf{G}_o^\top \mathbf{G}_o \mathbf{V} + \boldsymbol{\Sigma}^{-1} \right)^{-1}, \\ \mathbf{M} &= \mathbf{S} \left(\frac{1}{\gamma^2} \mathbf{V}^\top \mathbf{G}_o^\top \mathbf{W}_o + \boldsymbol{\Sigma}^{-1} \mathbf{R}\boldsymbol{\beta} \right). \end{aligned}$$

Hence the expectation of ℓ_c with respect to the distribution of $\mathbf{u}|\mathbf{W}_o, \tilde{\Theta}$ is:

$$\begin{aligned} E \left[\ell_c \mid \mathbf{W}_o, \tilde{\Theta} \right] &= -\frac{N_o}{2} \log \gamma^2 - \frac{1}{2} \log |\boldsymbol{\Sigma}| - \frac{1}{2\gamma^2} \text{Tr}(\mathbf{V}^\top \mathbf{G}_o^\top \mathbf{G}_o \mathbf{V} \tilde{\mathbf{S}}) \\ &\quad - \frac{1}{2\gamma^2} (\mathbf{G}_o \mathbf{V} \tilde{\mathbf{M}} - \mathbf{W}_o)^\top (\mathbf{G}_o \mathbf{V} \tilde{\mathbf{M}} - \mathbf{W}_o) \\ &\quad - \frac{1}{2} \text{Tr} \left(\boldsymbol{\Sigma}^{-1} \tilde{\mathbf{S}} \right) - \frac{1}{2} (\tilde{\mathbf{M}} - \mathbf{R}\boldsymbol{\beta})^\top \boldsymbol{\Sigma}^{-1} (\tilde{\mathbf{M}} - \mathbf{R}\boldsymbol{\beta}) \end{aligned} \quad (2)$$

The solutions for γ^2 , $\boldsymbol{\beta}$, and $\boldsymbol{\xi}$ that maximize (2), given current estimates, are relatively similar to the complete case:

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= (\mathbf{R}^\top \boldsymbol{\Sigma}^{-1} \mathbf{R})^{-1} \mathbf{R}^\top \boldsymbol{\Sigma}^{-1} \tilde{\mathbf{M}}, \\ \hat{\boldsymbol{\xi}} &= \arg \max_{\boldsymbol{\xi}} \left\{ -\frac{1}{2} \log |\boldsymbol{\Sigma}| - \frac{1}{2} \text{Tr} \left(\boldsymbol{\Sigma}^{-1} \tilde{\mathbf{S}} \right) - \frac{1}{2} (\tilde{\mathbf{M}} - \mathbf{R}\boldsymbol{\beta})^\top \boldsymbol{\Sigma}^{-1} (\tilde{\mathbf{M}} - \mathbf{R}\boldsymbol{\beta}) \right\}, \\ \hat{\gamma}^2 &= \frac{1}{N_o} \left[\text{Tr}(\mathbf{V}^\top \mathbf{G}_o^\top \mathbf{G}_o \mathbf{V} \tilde{\mathbf{S}}) + \|\mathbf{G}_o \mathbf{V} \tilde{\mathbf{M}} - \mathbf{W}_o\|_2^2 \right]. \end{aligned}$$

To solve for \mathbf{v} , we maximize a slightly different function,

$$\begin{aligned}
 h(\mathbf{v}) &= -\frac{1}{2\gamma^2} \text{Tr}(\mathbf{V}^\top \mathbf{G}_o^\top \mathbf{G}_o^\top \mathbf{V} \tilde{\mathbf{S}}) - \frac{1}{2\gamma^2} (\mathbf{G}_o \mathbf{V} \tilde{\mathbf{M}} - \mathbf{W}_o)^\top (\mathbf{G}_o \mathbf{V} \tilde{\mathbf{M}} - \mathbf{W}_o) \\
 &= -\frac{1}{2\gamma^2} \left[\text{Tr} \left(\begin{bmatrix} \sum_{j \in \Omega_1} v_j^2 & 0 & \dots & 0 \\ 0 & \sum_{j \in \Omega_2} v_j^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sum_{j \in \Omega_n} v_j^2 \end{bmatrix} \tilde{\mathbf{S}} \right) + \sum_{i=1}^n \sum_{j \in \Omega_i} (X_{ij} - \tilde{M}_i v_j)^2 \right] \\
 &= -\frac{1}{2\gamma^2} \left[\sum_{i=1}^n \sum_{j \in \Omega_i} \tilde{S}_{ii} v_j^2 + \sum_{i=1}^n \sum_{j \in \Omega_i} (X_{ij} - \tilde{M}_i v_j)^2 \right] \\
 &= -\frac{1}{2\gamma^2} \sum_{i=1}^n \sum_{j \in \Omega_i} \left[\tilde{S}_{ii} v_j^2 + (X_{ij} - \tilde{M}_i v_j)^2 \right] \\
 &= -\frac{1}{2\gamma^2} \sum_{i=1}^n \left[\tilde{S}_{ii} v_j^2 + (X_{ij} - \tilde{M}_i v_j)^2 \right] \mathbf{1}_{[j \in \Omega_i]}.
 \end{aligned}$$

Here Ω_i denotes the set of observed elements across the i -th row of \mathbf{X} and $\mathbf{1}_{[\cdot]}$ denotes the indicator function. Taking derivative of $h(\mathbf{v})$ with respect to each v_j and setting it equal to zero, we can find the closed-form unscaled solution

$$\check{v}_j = \frac{\sum_{i=1}^n X_{ij} \tilde{M}_i \mathbf{1}_{[j \in \Omega_i]}}{\sum_{i=1}^n (\tilde{S}_{ii} + \tilde{M}_i^2) \mathbf{1}_{[j \in \Omega_i]}}.$$

The final solution for \mathbf{v} is then $\hat{\mathbf{v}} = \frac{\check{\mathbf{v}}}{\|\check{\mathbf{v}}\|_2}$, where $\check{\mathbf{v}} = (\check{v}_1, \dots, \check{v}_p)$.

When some elements of the exposure data are missing, parameter estimation for each PC is based only on the observed elements \mathbf{W}_o . Estimate for PC score can then be made by projecting the model-based imputed exposure data onto the direction of \mathbf{v} . The joint distribution of \mathbf{W}_o and \mathbf{W}_m can be written as

$$\begin{bmatrix} \mathbf{W}_o \\ \mathbf{W}_m \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mathbf{m}_o \\ \mathbf{m}_m \end{bmatrix}, \begin{bmatrix} \mathbf{C}_{oo} & \mathbf{C}_{om} \\ \mathbf{C}_{mo} & \mathbf{C}_{mm} \end{bmatrix} \right) = \mathcal{N}(\mathbf{M}, \mathbf{C})$$

where $\mathbf{M} = \mathbf{GVR}\boldsymbol{\beta}$ and $\mathbf{C} = \gamma^2 \mathbf{G}\mathbf{G}^\top + \mathbf{GV}\boldsymbol{\Sigma}\mathbf{V}^\top \mathbf{G}^\top$. The missing elements, \mathbf{W}_m , can then be imputed by the conditional mean,

$$E(\mathbf{W}_m \mid \mathbf{W}_o) = \mathbf{m}_m + \mathbf{C}_{mo} \mathbf{C}_{oo}^{-1} (\mathbf{W}_o - \mathbf{m}_o).$$

Thus, the parameter estimation of ProPrPCA-Krige with missing monitoring data can be summarized as:

Algorithm: ProPrPCA-Krige with missing monitoring data

Input \mathbf{X} , \mathbf{R} , \mathbf{G}_o q , and t_{max}
for l in $\{1, \dots, q\}$ **do**
 $\mathbf{X}_l \leftarrow \mathbf{X}_{l-1}^{zero} - \hat{\mathbf{u}}_{l-1}^{zero} \hat{\mathbf{v}}_{l-1}^\top$ where $\mathbf{X}_0^{zero} = \mathbf{X}$ imputed with zeros, $\hat{\mathbf{u}}_0^{zero} = \mathbf{0}$, and $\hat{\mathbf{v}}_0 = \mathbf{0}$
 $\mathbf{W}_o \leftarrow \mathbf{G}_o \text{vec}(\mathbf{X}_l)$
Initialize $\mathbf{v}_l^{(0)}$, $(\gamma_l^{(0)})^2$, $\beta_l^{(0)}$, $\xi_l^{(0)}$, and $t = 1$
 $\Sigma_l^{(0)} \leftarrow \Sigma(\xi_l^{(0)})$
while not converged **or** $t < t_{max}$ **do**
 $\tilde{\mathbf{S}}_l \leftarrow \left[(\gamma_l^{(t)})^{-2} \mathbf{V}_l^{(t)\top} \mathbf{G}_o^\top \mathbf{G}_o \mathbf{V}_l^{(t)} + (\Sigma_l^{(t)})^{-1} \right]^{-1}$ where $\mathbf{V}_l^{(t)} = \mathbf{I}_n \otimes \mathbf{v}_l^{(t)}$
 $\tilde{\mathbf{M}}_l \leftarrow \tilde{\mathbf{S}}_l \left[(\gamma_l^{(t)})^{-2} \mathbf{V}_l^{(t)\top} \mathbf{G}_o^\top \mathbf{W}_o + (\Sigma_l^{(t)})^{-1} \mathbf{R} \beta_l^{(t)} \right]$
 $\mathbf{v}_l^{(t+1)} \leftarrow \tilde{\mathbf{v}}_l / \|\tilde{\mathbf{v}}_l\|_2$ where the j -th element of $\tilde{\mathbf{v}}_l$ (for $j = 1, \dots, p$) is calculated as:

$$\frac{\sum_{i=1}^n (\mathbf{X}_l)_{ij} (\tilde{\mathbf{M}}_l)_i \mathbf{1}_{[j \in \Omega_i.]}}{\sum_{i=1}^n \left((\tilde{\mathbf{S}}_l)_{ii} + (\tilde{\mathbf{M}}_l)_i^2 \right) \mathbf{1}_{[j \in \Omega_i.]}}$$

 $(\gamma_l^{(t+1)})^2 \leftarrow (N_o)^{-1} \left[\text{Tr}(\mathbf{V}_l^{(t+1)\top} \mathbf{G}_o^\top \mathbf{G}_o \mathbf{V}_l^{(t+1)} \tilde{\mathbf{S}}_l) + \|\mathbf{V}_l^{(t+1)\top} \mathbf{G}_o^\top \tilde{\mathbf{M}}_l - \mathbf{W}_o\|_2^2 \right]$

 where $\mathbf{V}_l^{(t+1)} = \mathbf{I}_n \otimes \mathbf{v}_l^{(t+1)}$
 $\xi_l^{(t+1)} \leftarrow \arg \max_{\xi_l} \left\{ -\log |\Sigma_l| - \text{Tr} \left(\Sigma_l^{-1} \tilde{\mathbf{S}}_l \right) - (\tilde{\mathbf{M}}_l - \mathbf{R} \beta_l^{(t)})^\top \Sigma_l^{-1} (\tilde{\mathbf{M}}_l - \mathbf{R} \beta_l^{(t)}) \right\}$

 where $\Sigma_l = \Sigma(\xi_l)$
 $\beta_l^{(t+1)} \leftarrow \left(\mathbf{R}^\top \hat{\Sigma}(\xi_l^{(t+1)})^{-1} \mathbf{R} \right)^{-1} \mathbf{R}^\top \hat{\Sigma}(\xi_l^{(t+1)})^{-1} \tilde{\mathbf{M}}_l$
 $t \leftarrow t + 1$
end while
 $\hat{\mathbf{v}}_l \leftarrow \mathbf{v}_l^{(t)}$, $\hat{\gamma}_l^2 \leftarrow (\gamma_l^{(t)})^2$, $\hat{\beta}_l \leftarrow \beta_l^{(t)}$, $\hat{\xi}_l \leftarrow \xi_l^{(t)}$
 $\mathbf{X}_l^{zero} \leftarrow \mathbf{X}_l$ with elements at missing indices replaced with zero

 $\mathbf{X}_l^{imp} \leftarrow \mathbf{X}_l$ with elements at missing indices replaced with conditional means

 $\hat{\mathbf{u}}_l^{zero} = \mathbf{X}_l^{zero} \hat{\mathbf{v}}_l$
 $\hat{\mathbf{u}}_l = \mathbf{X}_l^{imp} \hat{\mathbf{v}}_l$
end for
Output $\{\hat{\mathbf{v}}_1, \dots, \hat{\mathbf{v}}_q\}$, $\{\hat{\mathbf{u}}_1, \dots, \hat{\mathbf{u}}_q\}$, $\{\hat{\beta}_1, \dots, \hat{\beta}_q\}$, $\{\hat{\gamma}_1^2, \dots, \hat{\gamma}_q^2\}$, $\{\hat{\xi}_1, \dots, \hat{\xi}_q\}$

2 The ProPrPCA-Spline model and algorithm

2.1 The model

For each PC, the ProPrPCA-Spline algorithm assumes the following model

$$\mathbf{X} = \mathbf{Z}\boldsymbol{\beta}\mathbf{v}^\top + \mathbf{E}$$

Here \mathbf{Z} contains both geographic covariates and the thin-plate spline basis functions. The collection of model parameters Θ now includes $\{\mathbf{v}, \boldsymbol{\beta}, \gamma^2\}$. Using the same vectorization established in previous section, the model assumes $\mathbf{W}_i = \mathbf{X}_i = (\mathbf{Z}\boldsymbol{\beta})_i \mathbf{v} + \boldsymbol{\epsilon}_i$, for $i = 1, \dots, n$. We can then rewrite this model in the vectorized form as

$$\mathbf{W} \mid \Theta \sim \mathcal{N}(\mathbf{V}\mathbf{Z}\boldsymbol{\beta}, \gamma^2 \mathbf{I}_N),$$

where $\mathbf{V} = \mathbf{I}_n \otimes \mathbf{v}$. When there are missing data, the distribution of interest becomes

$$\mathbf{W}_o \mid \Theta \sim \mathcal{N}(\mathbf{G}_o \mathbf{V}\mathbf{Z}\boldsymbol{\beta}, \gamma^2 \mathbf{I}_{N_o}).$$

2.2 Estimation of model parameters when monitoring data is complete

To solve for the parameters, we maximize the log-likelihood (up to a constant) directly:

$$\ell(\Theta \mid \mathbf{W}) = -\frac{1}{N} \log \gamma^2 - \frac{1}{2\gamma^2} (\mathbf{W} - \mathbf{V}\mathbf{Z}\boldsymbol{\beta})^\top (\mathbf{W} - \mathbf{V}\mathbf{Z}\boldsymbol{\beta})$$

To solve for \mathbf{v} , we effectively maximize the following function:

$$-\frac{1}{2\gamma^2} (\mathbf{W} - \mathbf{V}\mathbf{Z}\boldsymbol{\beta})^\top (\mathbf{W} - \mathbf{V}\mathbf{Z}\boldsymbol{\beta})$$

Denote $\mathbf{K} = \mathbf{Z}\boldsymbol{\beta} \in \mathbb{R}^{n \times 1}$, we can rewrite this function similarly to the function involved \mathbf{v} with complete data for ProPrPCA-Krige. Thus, the solution for \mathbf{v} can be written in closed-form as

$$\hat{\mathbf{v}} = \frac{\check{\mathbf{v}}}{\|\check{\mathbf{v}}\|_2}, \quad \text{where } \check{\mathbf{v}} = \frac{\mathbf{X}^\top \mathbf{K}}{\|\mathbf{K}\|_2^2} = \frac{\mathbf{X}^\top \mathbf{Z}\boldsymbol{\beta}}{\|\mathbf{Z}\boldsymbol{\beta}\|_2^2}.$$

The closed-form solution for $\boldsymbol{\beta}$ is straightforwardly a result of ordinary least squares, $\hat{\boldsymbol{\beta}} = [(\mathbf{V}\mathbf{Z})^\top (\mathbf{V}\mathbf{Z})]^{-1} (\mathbf{V}\mathbf{Z})^\top \mathbf{W}$. This can be further simplified thanks to the constraint on \mathbf{v} and noticing that $\mathbf{V}^\top \mathbf{V} = \mathbf{I}_n$. Thus we have, $\hat{\boldsymbol{\beta}} = (\mathbf{Z}^\top \mathbf{Z})^{-1} (\mathbf{Z} \otimes \mathbf{v})^\top \mathbf{W}$. Finally, the solution for γ^2 is simply $\hat{\gamma}^2 = N^{-1} \|\mathbf{W} - \mathbf{V}\mathbf{Z}\boldsymbol{\beta}\|_2^2$. Thus parameter estimation of ProPrPCA-Spline with complete monitoring data can be summarized as:

Algorithm: ProPrPCA-Spline with complete monitoring data

Input \mathbf{X} , \mathbf{Z} , q , and t_{max}
for l in $\{1, \dots, q\}$ **do**
 $\mathbf{X}_l \leftarrow \mathbf{X}_{l-1} - \hat{\mathbf{u}}_{l-1} \hat{\mathbf{v}}_{l-1}^\top$ where $\mathbf{X}_0 = \mathbf{X}$, $\hat{\mathbf{u}}_0 = \mathbf{0}$, and $\hat{\mathbf{v}}_0 = \mathbf{0}$
Initialize $\mathbf{v}_l^{(0)}$, $(\gamma_l^{(0)})^2$, $\beta_l^{(0)}$, and $t = 1$
while not converged **or** $t < t_{max}$ **do**
 $\mathbf{v}_l^{(t+1)} \leftarrow \tilde{\mathbf{v}}_l / \|\tilde{\mathbf{v}}_l\|_2$ where $\tilde{\mathbf{v}}_l \leftarrow \mathbf{X}_l^\top \mathbf{Z} \beta_l^{(t)} / \|\mathbf{Z} \beta_l^{(t)}\|_2^2$
 $\beta_l^{(t+1)} \leftarrow (\mathbf{Z}^\top \mathbf{Z})^{-1} \left(\mathbf{Z} \otimes \mathbf{v}_l^{(t+1)} \right)^\top \text{vec}(\mathbf{X}_l)$
 $(\gamma_l^{(t+1)})^2 \leftarrow (np)^{-1} \|\text{vec}(\mathbf{X}_l) - (\mathbf{I}_n \otimes \mathbf{v}_l^{(t+1)}) \mathbf{Z} \beta_l^{(t+1)}\|_2^2$
 $t \leftarrow t + 1$
end while
 $\hat{\mathbf{v}}_l \leftarrow \mathbf{v}_l^{(t)}$, $\hat{\gamma}_l^2 \leftarrow (\gamma_l^{(t)})^2$, $\hat{\beta}_l \leftarrow \beta_l^{(t)}$
 $\hat{\mathbf{u}}_l = \mathbf{X}_l \hat{\mathbf{v}}_l$
end for
Output $\{\hat{\mathbf{v}}_1, \dots, \hat{\mathbf{v}}_q\}$, $\{\hat{\mathbf{u}}_1, \dots, \hat{\mathbf{u}}_q\}$, $\{\hat{\beta}_1, \dots, \hat{\beta}_q\}$, $\{\hat{\gamma}_1^2, \dots, \hat{\gamma}_q^2\}$

2.3 Parameter estimation and model-based imputation with missing monitoring data

The observed log-likelihood with missing monitoring data is

$$\ell(\Theta \mid \mathbf{W}_o) = -\frac{1}{N_o} \log \gamma^2 - \frac{1}{2\gamma^2} (\mathbf{W}_o - \mathbf{G}_o \mathbf{V} \mathbf{Z} \boldsymbol{\beta})^\top (\mathbf{W}_o - \mathbf{G}_o \mathbf{V} \mathbf{Z} \boldsymbol{\beta}).$$

The solutions for $\boldsymbol{\beta}$ and γ^2 are trivial and fairly similar to those with complete data. To solve for \mathbf{v} , we maximize $\{-\frac{1}{2\gamma^2} (\mathbf{W}_o - \mathbf{G}_o \mathbf{V} \mathbf{Z} \boldsymbol{\beta})^\top (\mathbf{W}_o - \mathbf{G}_o \mathbf{V} \mathbf{Z} \boldsymbol{\beta})\}$. We can rewrite the function of \mathbf{v} as

$$\begin{aligned} & -\frac{1}{2\gamma^2} (\mathbf{W}_o - \mathbf{G}_o \mathbf{V} \mathbf{Z} \boldsymbol{\beta})^\top (\mathbf{W}_o - \mathbf{G}_o \mathbf{V} \mathbf{Z} \boldsymbol{\beta}) \\ &= -\frac{1}{2\gamma^2} (\mathbf{W}_o - \mathbf{G}_o \mathbf{V} \mathbf{K})^\top (\mathbf{W}_o - \mathbf{G}_o \mathbf{V} \mathbf{K}) \\ &= \sum_{i=1}^n \sum_{j \in \Omega_i} (X_{ij} - K_i v_j)^2 = \sum_{i=1}^n (X_{ij} - K_i v_j)^2 \mathbf{1}_{[j \in \Omega_i]}, \end{aligned}$$

Taking derivative with respect to each v_j and setting it equal to zero, we can find the closed-form unscaled solution

$$\check{v} = \frac{\sum_{i=1}^n X_{ij} K_i \mathbf{1}_{[j \in \Omega_i.]}}{\sum_{i=1}^n K_i^2 \mathbf{1}_{[j \in \Omega_i.]}}.$$

The final solution for \mathbf{v} is then $\hat{\mathbf{v}} = \frac{\check{\mathbf{v}}}{\|\check{\mathbf{v}}\|_2}$, where $\check{\mathbf{v}} = (\check{v}_1, \dots, \check{v}_p)$.

When some elements of the exposure data are missing, parameter estimation for each PC is based only on the observed elements \mathbf{W}_o . Estimate for PC score can then be made by projecting the model-based imputed exposure data onto the direction of \mathbf{v} . The missing elements, \mathbf{W}_m , can then be imputed by its estimate $\mathbf{G}_m \hat{\mathbf{V}} \mathbf{Z} \hat{\boldsymbol{\beta}}$. Thus the parameter estimation of ProPrPCA-Spline with missing monitoring data can be summarized as:

Algorithm: ProPrPCA-Spline with missing monitoring data

Input \mathbf{X} , \mathbf{G}_o , \mathbf{Z} , q , and t_{max}
for l in $\{1, \dots, q\}$ **do**
 $\mathbf{X}_l \leftarrow \mathbf{X}_{l-1}^{zero} - \hat{\mathbf{u}}_{l-1}^{zero} \hat{\mathbf{v}}_{l-1}^\top$ where $\mathbf{X}_0^{zero} = \mathbf{X}$ imputed with zeros, $\hat{\mathbf{u}}_0^{zero} = \mathbf{0}$, and $\hat{\mathbf{v}}_0 = \mathbf{0}$
 $\mathbf{W}_o \leftarrow \mathbf{G}_o \text{vec}(\mathbf{X}_l)$
Initialize $\mathbf{v}_l^{(0)}$, $(\gamma_l^{(0)})^2$, $\beta_l^{(0)}$, and $t = 1$
while not converged **or** $t < t_{max}$ **do**
 $\mathbf{v}_l^{(t+1)} \leftarrow \tilde{\mathbf{v}}_l / \|\tilde{\mathbf{v}}_l\|_2$ where the j -element of $\tilde{\mathbf{v}}_l$ for $j = 1, \dots, p$ is calculated as:

$$\frac{\sum_{i=1}^n (\mathbf{X}_l)_{ij} (\mathbf{K}_l)_i \mathbf{1}_{[j \in \Omega_i.]}}{\sum_{i=1}^n (\mathbf{K}_l)_i^2 \mathbf{1}_{[j \in \Omega_i.]}} \text{, and } \mathbf{K}_l = \mathbf{Z} \beta_l^{(t)}$$

$$\beta_l^{(t+1)} \leftarrow \left[\left(\mathbf{G}_o \mathbf{V}_l^{(t+1)} \mathbf{Z} \right)^\top \left(\mathbf{G}_o \mathbf{V}_l^{(t+1)} \mathbf{Z} \right) \right]^{-1} \left(\mathbf{G}_o \mathbf{V}_l^{(t+1)} \mathbf{Z} \right)^\top \mathbf{W}_o$$

 where $\mathbf{V}_l^{(t+1)} = \mathbf{I}_n \otimes \mathbf{v}_l^{(t+1)}$

$$(\gamma_l^{(t+1)})^2 \leftarrow N_o^{-1} \|\mathbf{W}_o - \mathbf{G}_o \mathbf{V}_l^{(t+1)} \mathbf{Z} \beta_l^{(t+1)}\|_2^2$$

 $t \leftarrow t + 1$
end while
 $\hat{\mathbf{v}}_l \leftarrow \mathbf{v}_l^{(t)}$, $\hat{\gamma}_l^2 \leftarrow (\gamma_l^{(t)})^2$, $\hat{\beta}_l \leftarrow \beta_l^{(t)}$
 $\mathbf{X}_l^{zero} \leftarrow \mathbf{X}_l$ with elements at missing indices replaced with zero

 $\mathbf{X}_l^{imp} \leftarrow \mathbf{X}_l$ with elements at missing indices replaced with conditional means

 $\hat{\mathbf{u}}_l^{zero} = \mathbf{X}_l^{zero} \hat{\mathbf{v}}_l$
 $\hat{\mathbf{u}}_l = \mathbf{X}_l^{imp} \hat{\mathbf{v}}_l$
end for
Output $\{\hat{\mathbf{v}}_1, \dots, \hat{\mathbf{v}}_q\}$, $\{\hat{\mathbf{u}}_1, \dots, \hat{\mathbf{u}}_q\}$, $\{\hat{\beta}_1, \dots, \hat{\beta}_q\}$, $\{\hat{\gamma}_1^2, \dots, \hat{\gamma}_q^2\}$

3 Simulations

3.1 Data generating mechanism for high-dimensional simulations

To further demonstrate the performance of ProPrPCA, we simulate multi-pollutant exposure surfaces with $p = 15$. We first generate three underlying PC scores on the 100×100 grid ($N = 10,000$), such that

$$\begin{aligned} \mathbf{u}_j &\sim \mathcal{N}(\mathbf{R}_j \mathbf{b}_j, \mathbf{S}_j), \quad \text{where } j = 1, 2, 3, \\ \mathbf{R}_1 &= [\mathbf{r}_{1o} \quad \mathbf{r}_{1u}], \quad \text{where } \mathbf{r}_{1o}, \mathbf{r}_{1u} \sim \mathcal{N}(\mathbf{0}, \mathbf{1}), \quad \mathbf{b}_1^\top = [5 \quad 1], \\ \mathbf{R}_2 &= [\mathbf{r}_{2o} \quad \mathbf{r}_{2u}], \quad \text{where } \mathbf{r}_{2o}, \mathbf{r}_{2u} \sim \mathcal{N}(\mathbf{0}, \mathbf{1}), \quad \mathbf{b}_2^\top = [5 \quad 2], \\ \mathbf{R}_3 &= [\mathbf{r}_{3u}], \quad \text{where } \mathbf{r}_{3u} \sim \mathcal{N}(\mathbf{0}, \mathbf{1}), \quad \mathbf{b}_3 = 1. \end{aligned}$$

In this setting, \mathbf{r}_{jo} 's are GIS covariates observed for the model, while \mathbf{r}_{ju} 's are unobserved covariates, and used primarily to generate the scores themselves. That is, only $\mathbf{R} = [\mathbf{r}_{1o} \quad \mathbf{r}_{2o}]$ is used in the spatial prediction model. Here \mathbf{S}_1 has exponential structure with no nugget effect, partial sill of 5, and range of 50. Meanwhile, $\mathbf{S}_2 = 7.5\mathbf{I}_N$ and $\mathbf{S}_3 = 2\mathbf{I}_N$. This setup is created so that \mathbf{u}_1 is the most spatially predictable, \mathbf{u}_2 is moderately predictable in space, and \mathbf{u}_3 is not spatially predictable. Here spatial predictability refers to how well the quantity can be predicted at new locations using relevant and available covariates.

We then create two scenarios in which we scale the variance of \mathbf{u}_j 's differently,

$$\begin{aligned} \text{Scenario 1: } &Var(\mathbf{u}_1) = 10, Var(\mathbf{u}_2) = 7.5, Var(\mathbf{u}_3) = 5, \\ \text{Scenario 2: } &Var(\mathbf{u}_1) = 7.5, Var(\mathbf{u}_2) = 5, Var(\mathbf{u}_3) = 10. \end{aligned}$$

In both scenarios, the multi-pollutant exposure surface is generated as

$$\begin{aligned} \mathbf{X} &= \mathbf{UV} + \mathbf{E}, \quad \text{where } E_{ij} \sim \mathcal{N}(0, 1) \\ \mathbf{V} &= [\mathbf{v}_1 \quad \mathbf{v}_2 \quad \mathbf{v}_3], \quad \text{where } \mathbf{v}_j = \frac{\check{\mathbf{v}}_j}{\|\check{\mathbf{v}}_j\|_2}, \quad \text{for } j = 1, 2, 3 \\ \check{\mathbf{v}}_1^\top &= [1 \quad 1 \quad 1 \quad 1 \quad 1 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0] \\ \check{\mathbf{v}}_2^\top &= [0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 1 \quad 1 \quad 1 \quad 1 \quad 1 \quad 0 \quad 0 \quad 0 \quad 0] \\ \check{\mathbf{v}}_3^\top &= [0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 1 \quad 1 \quad 1 \quad 1] \end{aligned}$$

The use of such sparse loadings is to clearly identify the behavior of each dimension reduction method. Because of the variance contribution setup, in scenario A we expect all three methods to pick \mathbf{u}_1 as PC1, \mathbf{u}_2 as PC2, and \mathbf{u}_3 as PC3. In scenario B, however, we expect TradPCA to pick \mathbf{u}_3 as PC1, \mathbf{u}_1 as PC2, and \mathbf{u}_2 as PC3, as \mathbf{u}_3 has the largest variance contribution. Meanwhile, PredPCA and ProPrPCA-Spline will still pick \mathbf{u}_1 as PC1, \mathbf{u}_2 as PC2.

For these high-dimensional simulations, we consider three MCAR scenarios (30%, 35%, and 40%), and one MAR scenario. In the MAR scenario, we identify training locations with

r_{10} value larger than its sample 60th percentile, and among x_1 through x_5 , 75% of these training locations become missing data. For the rest of the pollutants, from x_6 to x_{15} , each has 25% of its locations missing completely at random. This setup guarantees a mild spatial pattern in the missing data, as x_1 to x_5 are generated entirely by u_1 , which is the most predictable score based on r_{10} .

3.2 Evaluation of computational burden

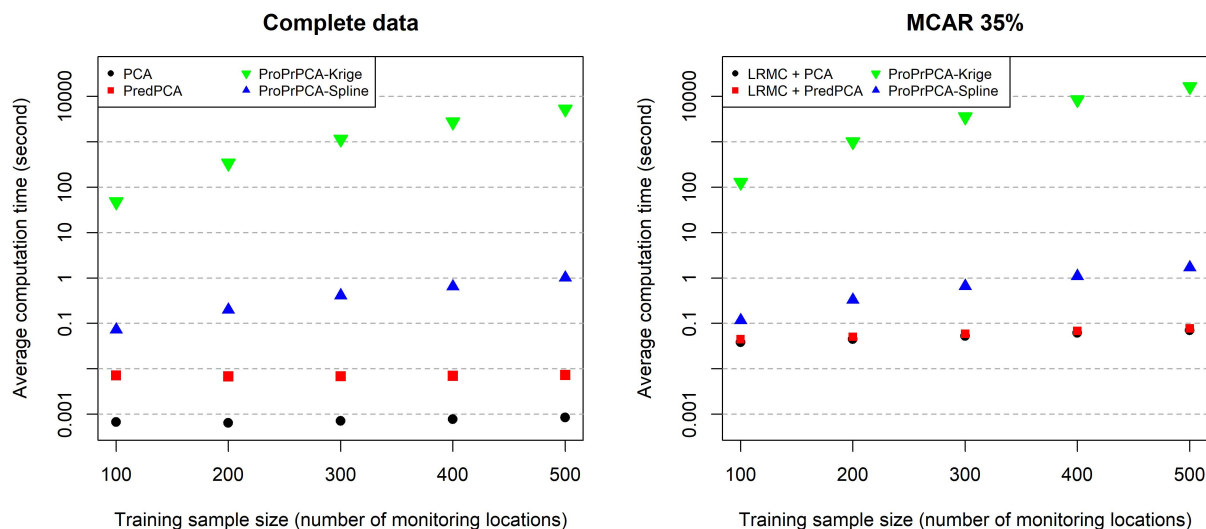


Figure 1: Computational time (average over 1,000 simulations under high-dimensional scenario 1) of PCA, PredPCA, ProPrPCA-Krige and ProPrPCA-Spline, with complete and MCAR 35% missing data by training sample size.

Figure 1 compares the computation burden among PCA, PredPCA, ProPrPCA-Krige, and ProPrPCA-Spline as the sample size increases. The results were averaged over 1,000 simulations under high-dimensional scenario 1, for complete data and MCAR 35% scenarios. The results were mostly under 1 second, on average, for PCA, PredPCA, and the Spline model. The computational burden of the Krige model is exponentially larger than the rest. While there maybe computational and programming tricks to alleviate the time cost, the Krige model would be likely to still take longer than other methods, given the nature of difficult optimization and EM algorithm.

4 Data Application

Table 1: Prediction R^2 's from leave-one-site-out cross-validation on 2010 CSN data. Sites with complete $PM_{2.5}$ component data are used as test data. Training data may include only complete sites (complete training data), or all available sites (full training data). The PCs are defined in the order of which they are obtained, i.e. PC1, PC2, PC3 are the first, second, and third PCs returned by the methods, respectively.

	PC1	PC2	PC3
PCA (complete training data)	0.24	0.51	0.51
PredPCA (complete training data)	0.52	0.44	0.62
LRMC + PredPCA (full training data)	0.54	0.53	0.45
ProPrPCA-Spline (full training data)	0.57	0.35	0.69

In this section, we evaluate the predictive performance in leave-one-site-out cross-validations. Table 1 shows the results when the PCs are simply defined in the order of which they are obtained. As discussed in the main text of the manuscript, while having decent performance for PC2 and PC3, using PCA applied to the complete training data yields a poor result for PC1. PredPCA has similar performances for PC1 with either complete or full training data. However, there is a trade-off in performances between PC2 and PC3, which can potentially be explained by the switching between PC2 and PC3 observed in the pollutant profile. ProPrPCA-Spline applied on the full training data shows the highest predictive performance for PC1 and PC3, but suffers from a decrease in the ability to predict PC2 well.

Table 2: Prediction R^2 's from leave-one-site-out cross-validation on 2010 CSN data. Sites with complete $PM_{2.5}$ component data are used as test data. Training data may include only complete sites (complete training data), or all available sites (full training data). The PCs are defined in the order of variance explained in the training data, i.e. PC1 is the component with the largest variance contributed, and so on.

	PC1	PC2	PC3
PCA (complete training data)	0.24	0.51	0.51
PredPCA (complete training data)	0.52	0.44	0.62
LRMC + PredPCA (full training data)	0.54	0.45	0.53
ProPrPCA-Spline (full training data)	0.57	0.41	0.65

Similar to Shen and Huang (2008), we also order the PCs by the variance explained in the training data, as given in Table 2. That is, in each round of the cross-validation procedure, out of the three PCs obtained by each PCA method, PC1 is defined as the component with the largest variance explained in the training data, and so on for PC2 and PC3. For spatially predictive methods, especially PredPCA, the order of variance explain was not necessary the same as the order by which the PCs were produced. Thus, while the results for PCA (and PredPCA) with complete training data were the same in both

Tables, the results with full training data were slightly different for the last two PCs. With this reordering approach, the results were more similar for PC2, and there was an apparent advantage of using ProPrPCA-Spline compared to PredPCA after imputation.

References

- Shen, H. and Huang, J. Z. (2008). Sparse principal component analysis via regularized low rank matrix approximation. *Journal of Multivariate Analysis*, 99(6):1015–1034.
- Tipping, M. E. and Bishop, C. M. (1999). Probabilistic principal component analysis. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 61(3):611–622.