

Australasian Applied Statistics Conference 2024

# Cokrig-and-Regress for Spatially Misaligned Environmental Data

---

**Zhi Yang Tho**

Joint work with Dr. Francis Hui, Prof. Alan Welsh, Dr. Tao Zou



**Australian  
National  
University**

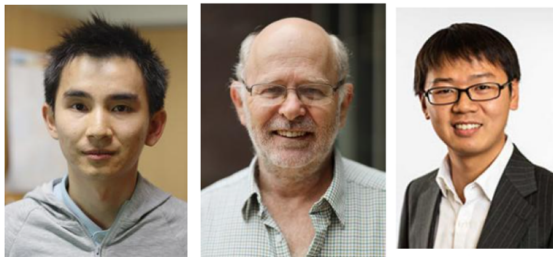


Figure: Left to right: Dr. Francis Hui, Prof. Alan Welsh, Dr. Tao Zou

1. Motivation
2. Model Set-Up
3. Cokrig-and-Regress (CNR)
4. Uncertainty Quantification
5. Simulation Studies

# Motivation

---

# Motivation

**Spatial misalignment** problem between  $PM_{2.5}$  concentration (response) and meteorological variables (covariates) such as temperature and precipitation.

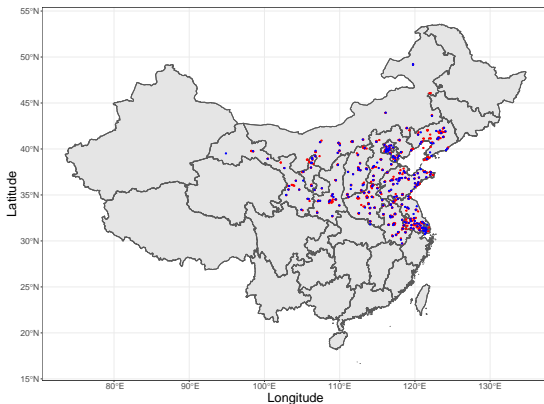


Figure: Map of China with geographic locations of **pollution monitoring stations** and **meteorological stations**.

# Motivation (Cont'd)

Existing methods to construct spatially aligned datasets:

- ⊙ **Nearest-neighbor interpolation** (Jhun et al., 2015; Greenstone et al., 2022);
- ⊙ **Kriging** for each meteorological covariate separately and treating predicted covariates as **fixed** (Reich et al., 2011; Liu et al., 2020);
- ⊙ **Krige-and-regress (KNR)** method that accounts for additional variability of the predicted covariate (Madsen et al., 2008; Szpiro et al., 2011; Pouliot, 2023)

However, KNR method only allows for

- ⊙ A single misaligned meteorological covariate;
- ⊙ Simple linear pollution-meteorological relationship.

# Model Set-Up



# Model Set-Up

Denote the pollution stations and meteorological stations as  $S = \{s_1, \dots, s_N\}$  and  $\tilde{S} = \{\tilde{s}_1, \dots, \tilde{s}_M\}$ , respectively.

Let  $\mathbf{y} = (y_1, \dots, y_N)^\top = (y(s_1), \dots, y(s_N))^\top$ ,  
 $\mathbf{x}_k = (x_{1k}, \dots, x_{Nk})^\top = (x_k(s_1), \dots, x_k(s_N))^\top$  and  
 $\tilde{\mathbf{x}}_k = (\tilde{x}_{1k}, \dots, \tilde{x}_{Mk})^\top = (x_k(\tilde{s}_1), \dots, x_k(\tilde{s}_M))^\top$ ,

$$y_i = \beta_0 + \sum_{k=1}^K \mathbf{f}_k(x_{ik})^\top \boldsymbol{\beta}_k + \rho_i + \epsilon_i, \text{ for } i = 1, \dots, N, \quad (1)$$

where **only**  $\mathbf{y}$  and  $\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_K$  are observed but not  $\mathbf{x}_1, \dots, \mathbf{x}_K$ .



## Model Set-Up (Cont'd)

Let  $\mathbf{x} = (\tilde{\mathbf{x}}_1^\top, \mathbf{x}_1^\top, \dots, \tilde{\mathbf{x}}_K^\top, \mathbf{x}_K^\top)^\top$  denote the stacked  $K(M+N)$ -vector of  $K$  covariates at locations in both  $\tilde{S}$  and  $S$ . Assume

$$\mathbf{x} \sim N(\boldsymbol{\mu} \otimes \mathbf{1}_{M+N}, \boldsymbol{\Sigma}),$$

$$\boldsymbol{\Sigma} = \text{Bdiag}(\mathbf{L}_1, \dots, \mathbf{L}_K)(\mathbf{R} \otimes \mathbf{I}_{M+N})\text{Bdiag}(\mathbf{L}_1, \dots, \mathbf{L}_K)^\top,$$

where  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_K)^\top$ ,  $\mathbf{R}$  is a  $K \times K$  cross-correlation matrix for the  $K$  meteorological covariates,  $\mathbf{L}_k$  is the lower Cholesky factor of  $\boldsymbol{\Sigma}_k$  which is the Matern spatial covariance matrix for the  $k$ -th covariate.

# Cokrig-and-Regress (CNR)

---

Estimate the parameters of the **joint** distribution of the meteorological covariates, based on the **observed misaligned meteorological data**  $\mathbf{x}_{\tilde{S}} = (\tilde{\mathbf{x}}_1^\top, \dots, \tilde{\mathbf{x}}_K^\top)^\top$ .

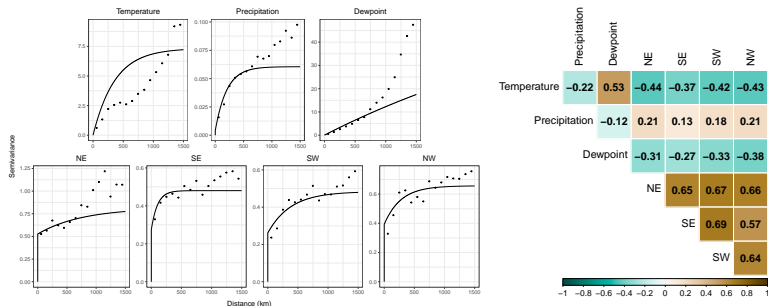


Figure: Estimated marginal Matérn covariance models (left), estimated cross-correlation matrix (right) from CNR Step 1.

**Predict** the unobserved  $\mathbf{x}_S = (\mathbf{x}_1^\top, \dots, \mathbf{x}_K^\top)^\top$  by **cokriging**, based on the observed  $\mathbf{x}_{\tilde{S}}$  and estimated parameters for the joint distribution of the meteorological covariates from CNR Step 1.

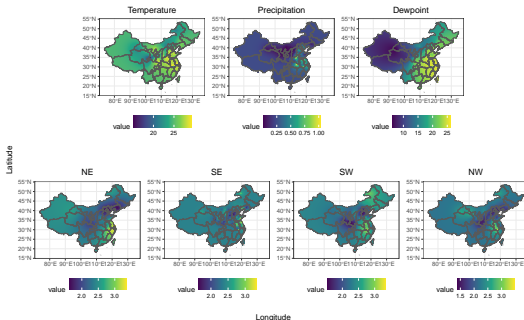


Figure: Spatial maps for the cokriging predictor of each meteorological covariate in CNR Step 2.

Replace the unobserved meteorological covariates with their cokriging prediction and **fit the spatial linear mixed model**

(1) i.e.,  $y_i = \beta_0 + \sum_{k=1}^K \mathbf{f}_k(\hat{x}_{ik})^\top \boldsymbol{\beta}_k + \rho_i + \epsilon_i$ .

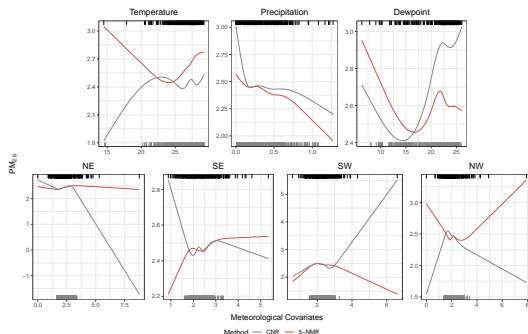


Figure: Estimated conditional smoothers for the seven meteorological covariates based on CNR and 5-NMR using natural cubic splines.

# Uncertainty Quantification

---

To estimate variance of the CNR estimates or construct confidence intervals for  $\beta_k$ :

- (i) Perform **preliminary parametric bootstrap** to bias-correct CNR spatial covariance parameter estimates (which was shown to be biased).
- (ii) Based on the bias-corrected CNR spatial covariance parameter estimates, perform **secondary parametric bootstrap** to obtain bootstrap samples of CNR estimates for  $\beta_k$ .

# Bootstrap CI compared to Naive Variance Estimator

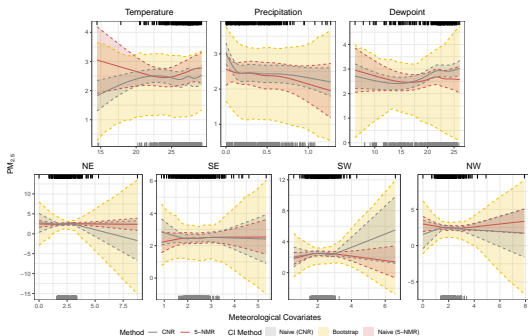


Figure: Estimated conditional smoothers for the seven meteorological covariates based on CNR (grey solid lines) and 5-NMR (red solid lines) using natural cubic splines. Also shown are 95% confidence bands based on the naive variance estimator for CNR (shaded regions between dashed lines in grey), bootstrap percentile confidence bands for CNR (shaded regions between dashed lines in yellow) and the naive variance estimator for 5-NMR (shaded regions between dashed lines in red).



# Simulation Studies

---

# Simulation Settings

- ⊙  $\tilde{S}$  and  $S$  are locations of  $M = 243$  meteorological stations and  $N = 796$  pollution monitoring stations, respectively.
- ⊙  $K = 5$  covariates with  $\mathbf{f}_k(x_{ik})^\top \boldsymbol{\beta}_k = x_{ik}\beta_k$  for  $k = 1, \dots, 5$ , and  $\boldsymbol{\beta} = (2, 1, 0.5, 1, 0.5, 1)^\top$ .
- ⊙ A total of 400 simulated datasets.

# Simulation Settings

- ⊙ Bias and RMSE for  $\beta_k$ :
  - CNR;
  - 5-nearest-matching-and-regress (5-NMR).
- ⊙ Ratio of average estimated standard error to empirical standard deviation (ASE/ESD), empirical coverage probability of 95% CIs for  $\beta_k$ :
  - Naive variance estimator (Naive) by ignoring the cokriging prediction uncertainty;
  - Proposed bootstrap approach.

# Simulation Results

Table: Bias, RMSE, ASE/ESD and empirical coverage of  $\beta_k$  for  $k = 1, \dots, 5$ .

|        |                  | Point Estimation   |                 |               |                 |               |
|--------|------------------|--------------------|-----------------|---------------|-----------------|---------------|
|        | Method           | $\beta_1 = 1$      | $\beta_2 = 0.5$ | $\beta_3 = 1$ | $\beta_4 = 0.5$ | $\beta_5 = 1$ |
| Bias   | CNR              | -0.0068            | 0.0021          | -0.0070       | 0.0046          | -0.0111       |
|        | 5-NMR            | -0.1529            | -0.0752         | -0.1591       | -0.0665         | -0.1580       |
| RMSE   | CNR              | 0.1733             | 0.1670          | 0.2068        | 0.2137          | 0.1948        |
|        | 5-NMR            | 0.2516             | 0.2113          | 0.2705        | 0.2535          | 0.2872        |
|        |                  | ASE/ESD            |                 |               |                 |               |
| Method | Inference Method | $\beta_1 = 1$      | $\beta_2 = 0.5$ | $\beta_3 = 1$ | $\beta_4 = 0.5$ | $\beta_5 = 1$ |
| CNR    | Naive            | 0.7495             | 0.7749          | 0.7236        | 0.7811          | 0.7750        |
|        | Bootstrap        | 1.0896             | 1.0789          | 1.0300        | 1.0628          | 1.0928        |
| 5-NMR  | Naive            | 0.8096             | 0.8159          | 0.8571        | 0.8541          | 0.7827        |
|        |                  | Empirical Coverage |                 |               |                 |               |
| Method | Inference Method | $\beta_1 = 1$      | $\beta_2 = 0.5$ | $\beta_3 = 1$ | $\beta_4 = 0.5$ | $\beta_5 = 1$ |
| CNR    | Naive            | 0.8400             | 0.8650          | 0.8275        | 0.8875          | 0.8750        |
|        | Bootstrap        | 0.9525             | 0.9625          | 0.9600        | 0.9600          | 0.9550        |
| 5-NMR  | Naive            | 0.7950             | 0.8500          | 0.8175        | 0.8750          | 0.8125        |

- Greenstone, M., He, G., Jia, R., and Liu, T. (2022). Can technology solve the principal-agent problem? Evidence from China's war on air pollution. *American Economic Review: Insights*, 4:54–70.
- Jhun, I., Coull, B. A., Schwartz, J., Hubbell, B., and Koutrakis, P. (2015). The impact of weather changes on air quality and health in the United States in 1994–2012. *Environmental Research Letters*, 10:084009.
- Liu, Y., Zhou, Y., and Lu, J. (2020). Exploring the relationship between air pollution and meteorological conditions in China under environmental governance. *Scientific Reports*, 10:14518.
- Madsen, L., Ruppert, D., and Altman, N. S. (2008). Regression with spatially misaligned data. *Environmetrics*, 19:453–467.

- Pouliot, G. A. (2023). Spatial econometrics for misaligned data. *Journal of Econometrics*, 232:168–190.
- Reich, B. J., Eidsvik, J., Guindani, M., Nail, A. J., and Schmidt, A. M. (2011). A class of covariate-dependent spatiotemporal covariance functions for the analysis of daily ozone concentration. *The Annals of Applied Statistics*, 5:2425–2447.
- Szpiro, A. A., Sheppard, L., and Lumley, T. (2011). Efficient measurement error correction with spatially misaligned data. *Biostatistics*, 12:610–623.

R codes available at <https://github.com/Zy1225/CNR>.



THE  
END

THANKS!