

# Spatial Prediction of Fine Particulate Matter in Taiwan

Hsin-Cheng Huang

Institute of Statistical Science, Academia Sinica

July 25, 2017

Joint work with

Guowen Huang (National Tsing Hua University, Taiwan)

Wen-Han Hwang (National Chung Hsing University, Taiwan)

# Outline

1 PM<sub>2.5</sub> Data

2 Resolution Adaptive Fixed Rank Kriging

3 Extensions

4 PM<sub>2.5</sub> Applications

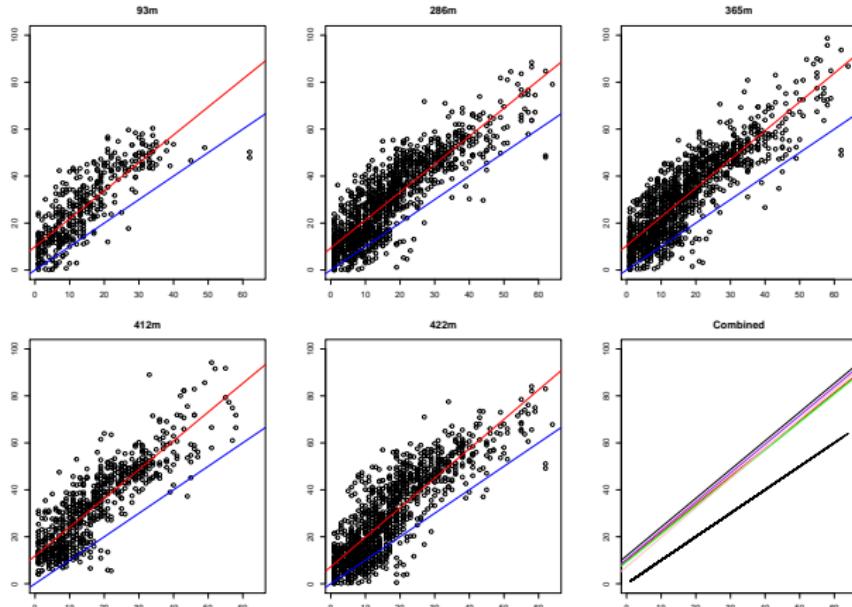
# PM<sub>2.5</sub> Data

# PM<sub>2.5</sub> Data Analysis

- Particulate matter less than 2.5 micrometers (PM<sub>2.5</sub>) can cause cardiovascular and lung diseases
- Two data sources
  - 1,147 AirBoxes (available at <https://pm25.lass-net.org/>)
  - 76 Taiwan EPA monitoring stations
- Data: Hourly measurements from Jan 1 to Feb 28, 2017
- Goals
  - Spatial prediction of ground-level PM<sub>2.5</sub> concentration at any location in Taiwan
  - Detection of abnormal AirBox measurements and emission sources

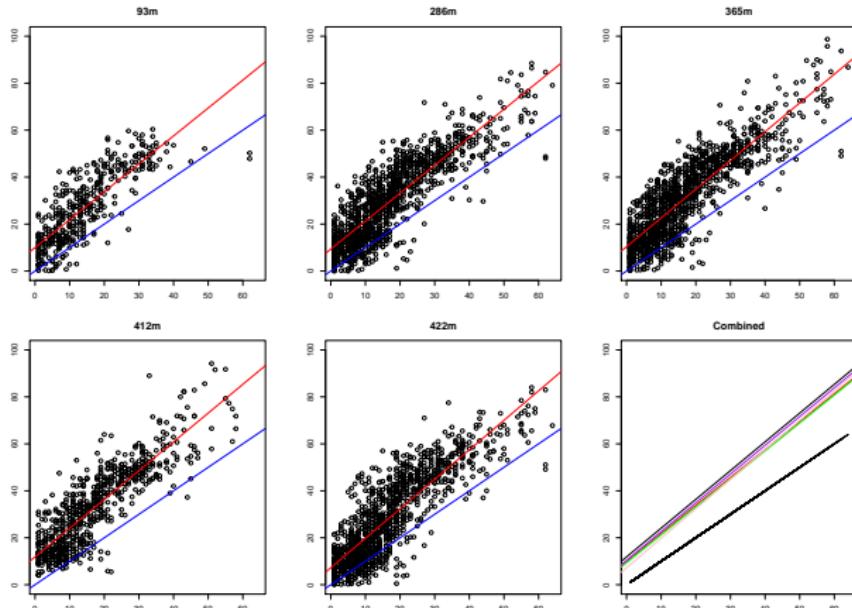
1,147 AirBoxes (Left) & 76 Monitoring Stations (Right)

# Calibration Between Two Types of Measurements



- x-axis: PM<sub>2.5</sub> measurements at an EPA monitoring station
- y-axis: PM<sub>2.5</sub> measurements at nearby AirBoxes
- Monitoring stations produce higher quality measurements but lower PM<sub>2.5</sub> values, because they are located at higher elevations

# Calibration Between Two Types of Measurements



- x-axis: PM<sub>2.5</sub> measurements at an EPA monitoring station
- y-axis: PM<sub>2.5</sub> measurements at nearby AirBoxes
- Monitoring stations produce higher quality measurements but lower PM<sub>2.5</sub> values, because they are located at higher elevations

# Resolution Adaptive Fixed Rank Kriging

# Spatial Random Effects Models

- Spatial process of interest at time  $t$ ;  $t = 1, \dots, T$

$$y_t(\mathbf{s}) = \sum_{k=1}^K \mathbf{w}_{t,k} \varphi_k(\mathbf{s}); \quad \mathbf{s} \in D$$

- $\mathbf{w}_t = (w_{t,1}, \dots, w_{t,K})' \sim N(\mathbf{0}, \mathbf{M})$
- $\varphi_1(\mathbf{s}), \dots, \varphi_K(\mathbf{s})$  are basis functions

- Spatial covariance function

$$C_y(\mathbf{s}, \mathbf{s}^*) = \varphi(\mathbf{s})' \mathbf{M} \varphi(\mathbf{s}^*)$$

- $\varphi(\mathbf{s}) = (\varphi_1(\mathbf{s}), \dots, \varphi_K(\mathbf{s}))'$
- Data:  $\mathbf{Z}_t = (z_t(\mathbf{s}_1), \dots, z_t(\mathbf{s}_n))'$ ;  $t = 1, \dots, T$

$$z_t(\mathbf{s}_i) = y_t(\mathbf{s}_i) + \varepsilon_t(\mathbf{s}_i)$$

- $\varepsilon_t(\mathbf{s}_i) \sim N(0, \sigma^2)$

# Spatial Random Effects Models

- Spatial process of interest at time  $t$ ;  $t = 1, \dots, T$

$$y_t(\mathbf{s}) = \sum_{k=1}^K \mathbf{w}_{t,k} \varphi_k(\mathbf{s}); \quad \mathbf{s} \in D$$

- $\mathbf{w}_t = (w_{t,1}, \dots, w_{t,K})' \sim N(\mathbf{0}, \mathbf{M})$
- $\varphi_1(\mathbf{s}), \dots, \varphi_K(\mathbf{s})$  are basis functions

- Spatial covariance function

$$C_y(\mathbf{s}, \mathbf{s}^*) = \varphi(\mathbf{s})' \mathbf{M} \varphi(\mathbf{s}^*)$$

- $\varphi(\mathbf{s}) = (\varphi_1(\mathbf{s}), \dots, \varphi_K(\mathbf{s}))'$

- Data:  $\mathbf{Z}_t = (z_t(\mathbf{s}_1), \dots, z_t(\mathbf{s}_n))'$ ;  $t = 1, \dots, T$

$$z_t(\mathbf{s}_i) = y_t(\mathbf{s}_i) + \varepsilon_t(\mathbf{s}_i)$$

- $\varepsilon_t(\mathbf{s}_i) \sim N(0, \sigma^2)$

# Spatial Random Effects Models

- Spatial process of interest at time  $t$ ;  $t = 1, \dots, T$

$$y_t(\mathbf{s}) = \sum_{k=1}^K \mathbf{w}_{t,k} \varphi_k(\mathbf{s}); \quad \mathbf{s} \in D$$

- $\mathbf{w}_t = (w_{t,1}, \dots, w_{t,K})' \sim N(\mathbf{0}, \mathbf{M})$
- $\varphi_1(\mathbf{s}), \dots, \varphi_K(\mathbf{s})$  are basis functions

- Spatial covariance function

$$C_y(\mathbf{s}, \mathbf{s}^*) = \varphi(\mathbf{s})' \mathbf{M} \varphi(\mathbf{s}^*)$$

- $\varphi(\mathbf{s}) = (\varphi_1(\mathbf{s}), \dots, \varphi_K(\mathbf{s}))'$
- Data:  $\mathbf{Z}_t = (z_t(\mathbf{s}_1), \dots, z_t(\mathbf{s}_n))'$ ;  $t = 1, \dots, T$

$$z_t(\mathbf{s}_i) = y_t(\mathbf{s}_i) + \varepsilon_t(\mathbf{s}_i)$$

- $\varepsilon_t(\mathbf{s}_i) \sim N(0, \sigma^2)$

# Spatial Random Effects Models

- Minimum mean-squared error predictor of  $y_t(\mathbf{s})$

$$\hat{y}_t(\mathbf{s}) = \varphi(\mathbf{s})' \mathbf{M} \Phi' (\Phi \mathbf{M} \Phi' + \sigma^2 \mathbf{I}_n)^{-1} \mathbf{Z}_t$$

- $\Phi = (\varphi_j(\mathbf{s}_i))_{n \times K}$
- $\varphi(\mathbf{s}) = (\varphi_1(\mathbf{s}), \dots, \varphi_K(\mathbf{s}))'$

- Mean squared prediction error (prediction variance)

$$E(\hat{y}_t(\mathbf{s}) - y_t(\mathbf{s}))^2 = \varphi(\mathbf{s})' \{ \mathbf{M} - \mathbf{M} \Phi' (\Phi \mathbf{M} \Phi' + \sigma^2 \mathbf{I})^{-1} \Phi \mathbf{M} \} \varphi(\mathbf{s})$$

- Sherman-Morrison-Woodbury formula

$$(\Phi \mathbf{M} \Phi' + \sigma^2 \mathbf{I}_n)^{-1} = \frac{1}{\sigma^2} \mathbf{I}_n - \frac{1}{\sigma^4} \Phi \left\{ \mathbf{M}^{-1} + \frac{1}{\sigma^2} \Phi' \Phi \right\}^{-1} \Phi'$$

- Models for  $\mathbf{M}$

- LatticeKrig (Nychka *et al.*, 2015):  $\mathbf{M}(\theta)$  depends on a small number of parameters  $\theta$
- Fixed rank kriging (Cressie and Johannesson, 2008): No structure is imposed on  $\mathbf{M}$

# Spatial Random Effects Models

- Minimum mean-squared error predictor of  $y_t(\mathbf{s})$

$$\hat{y}_t(\mathbf{s}) = \varphi(\mathbf{s})' \mathbf{M} \Phi' (\Phi \mathbf{M} \Phi' + \sigma^2 \mathbf{I}_n)^{-1} \mathbf{Z}_t$$

- $\Phi = (\varphi_j(\mathbf{s}_i))_{n \times K}$
- $\varphi(\mathbf{s}) = (\varphi_1(\mathbf{s}), \dots, \varphi_K(\mathbf{s}))'$

- Mean squared prediction error (prediction variance)

$$E(\hat{y}_t(\mathbf{s}) - y_t(\mathbf{s}))^2 = \varphi(\mathbf{s})' \{ \mathbf{M} - \mathbf{M} \Phi' (\Phi \mathbf{M} \Phi' + \sigma^2 \mathbf{I})^{-1} \Phi \mathbf{M} \} \varphi(\mathbf{s})$$

- Sherman-Morrison-Woodbury formula

$$(\Phi \mathbf{M} \Phi' + \sigma^2 \mathbf{I}_n)^{-1} = \frac{1}{\sigma^2} \mathbf{I}_n - \frac{1}{\sigma^4} \Phi \left\{ \mathbf{M}^{-1} + \frac{1}{\sigma^2} \Phi' \Phi \right\}^{-1} \Phi'$$

- Models for  $\mathbf{M}$

- LatticeKrig (Nychka *et al.*, 2015):  $\mathbf{M}(\theta)$  depends on a small number of parameters  $\theta$
- Fixed rank kriging (Cressie and Johannesson, 2008): No structure is imposed on  $\mathbf{M}$

# Fixed Rank Kriging (FRK)

$$y_t(\mathbf{s}) = \sum_{k=1}^K w_{t,k} \varphi_k(\mathbf{s}), \quad \mathbf{w}_t = (w_{t,1}, \dots, w_{t,K})' \sim N(\mathbf{0}, \mathbf{M})$$

- No structure is assumed on  $\mathbf{M}$  except that it is non-negative definite
- Parameters:  $\mathbf{M}$  and  $\sigma^2$  (with  $K(K + 1)/2 + 1$  parameters)
- Advantages
  - Flexible nonstationary spatial covariance models
  - Handle massive amounts of data
- Issues
  1. Maximum likelihood estimation: EM algorithm (Katzfuss and Cressie, 2009)
  2. How to allocate basis functions?
  3. How many basis functions?

# Improvement 1: Maximum Likelihood Estimation

- Parameters:  $\mathbf{M}$  and  $\sigma^2$  (with  $K(K+1)/2 + 1$  parameters)

Theorem (Tzeng and Huang. 2017+, Technometrics)

$$\hat{\mathbf{M}} = (\Phi' \Phi)^{-1/2} \mathbf{P} \text{diag}(\hat{d}_1, \dots, \hat{d}_K) \mathbf{P}' (\Phi' \Phi)^{-1/2},$$

$$\hat{\sigma}^2 = \max \left\{ \frac{1}{n - K^*} \left( \text{tr}(\mathbf{S}) - \sum_{k=0}^{K^*} d_k \right), 0 \right\}$$

- $\mathbf{P} \text{diag}(d_1, \dots, d_K) \mathbf{P}'$  is the eigen-decomposition of

$$(\Phi' \Phi)^{-1/2} \Phi' \mathbf{S} \Phi (\Phi' \Phi)^{-1/2}$$

- $\mathbf{S} = \frac{1}{T} \sum_{t=1}^T \mathbf{Z}_t \mathbf{Z}_t'$

- $d_0 \equiv 0$  and  $\hat{d}_k = \max(d_k - \hat{\sigma}^2, 0); k = 1, \dots, K$

- $K^* = \max \left( \{0\} \cup \left\{ L : d_L > \frac{1}{n - L} \left( \text{tr}(\mathbf{S}) - \sum_{k=0}^L d_k \right) \right\} \right)$

# Improvement 1: Maximum Likelihood Estimation

- Parameters:  $\mathbf{M}$  and  $\sigma^2$  (with  $K(K+1)/2 + 1$  parameters)

Theorem (Tzeng and Huang. 2017+, Technometrics)

$$\hat{\mathbf{M}} = (\Phi' \Phi)^{-1/2} \mathbf{P} \operatorname{diag}(\hat{d}_1, \dots, \hat{d}_K) \mathbf{P}' (\Phi' \Phi)^{-1/2},$$

$$\hat{\sigma}^2 = \max \left\{ \frac{1}{n - K^*} \left( \operatorname{tr}(\mathbf{S}) - \sum_{k=0}^{K^*} d_k \right), 0 \right\}$$

- $\mathbf{P} \operatorname{diag}(d_1, \dots, d_K) \mathbf{P}'$  is the eigen-decomposition of

$$(\Phi' \Phi)^{-1/2} \Phi' \mathbf{S} \Phi (\Phi' \Phi)^{-1/2}$$

- $\mathbf{S} = \frac{1}{T} \sum_{t=1}^T \mathbf{Z}_t \mathbf{Z}_t'$

- $d_0 \equiv 0$  and  $\hat{d}_k = \max(d_k - \hat{\sigma}^2, 0); k = 1, \dots, K$

- $K^* = \max \left( \{0\} \cup \left\{ L : d_L > \frac{1}{n - L} \left( \operatorname{tr}(\mathbf{S}) - \sum_{k=0}^L d_k \right) \right\} \right)$

# Improvement 2: Multi-Resolution Basis Functions

- $J(f) = \int_{\mathbb{R}^d} \sum_{\nu_1 + \dots + \nu_d = 2} \frac{2}{\nu_1! \dots \nu_d!} \left( \frac{\partial^2 f(\mathbf{s})}{\partial x_1^{\nu_1} \dots \partial x_d^{\nu_d}} \right)^2 d\mathbf{s}; \quad \mathbf{s} = (x_1, \dots, x_d)'$

Definition (Multiresolution spline basis functions)

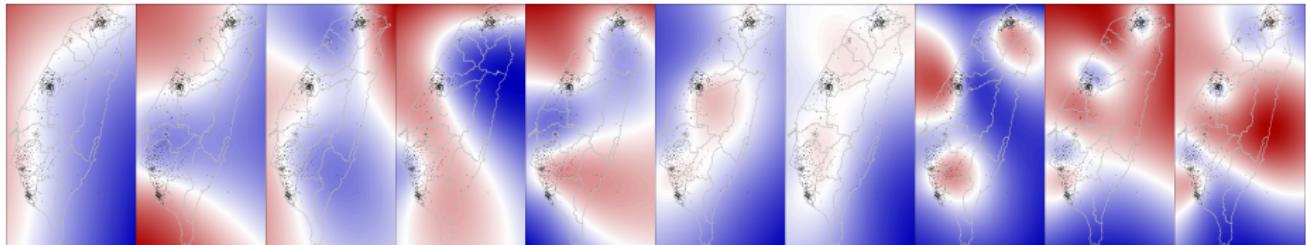
$$\varphi_1(\mathbf{s}) = 1$$

$$\varphi_k(\mathbf{s}) = x_k; \quad k = 2, \dots, d+1$$

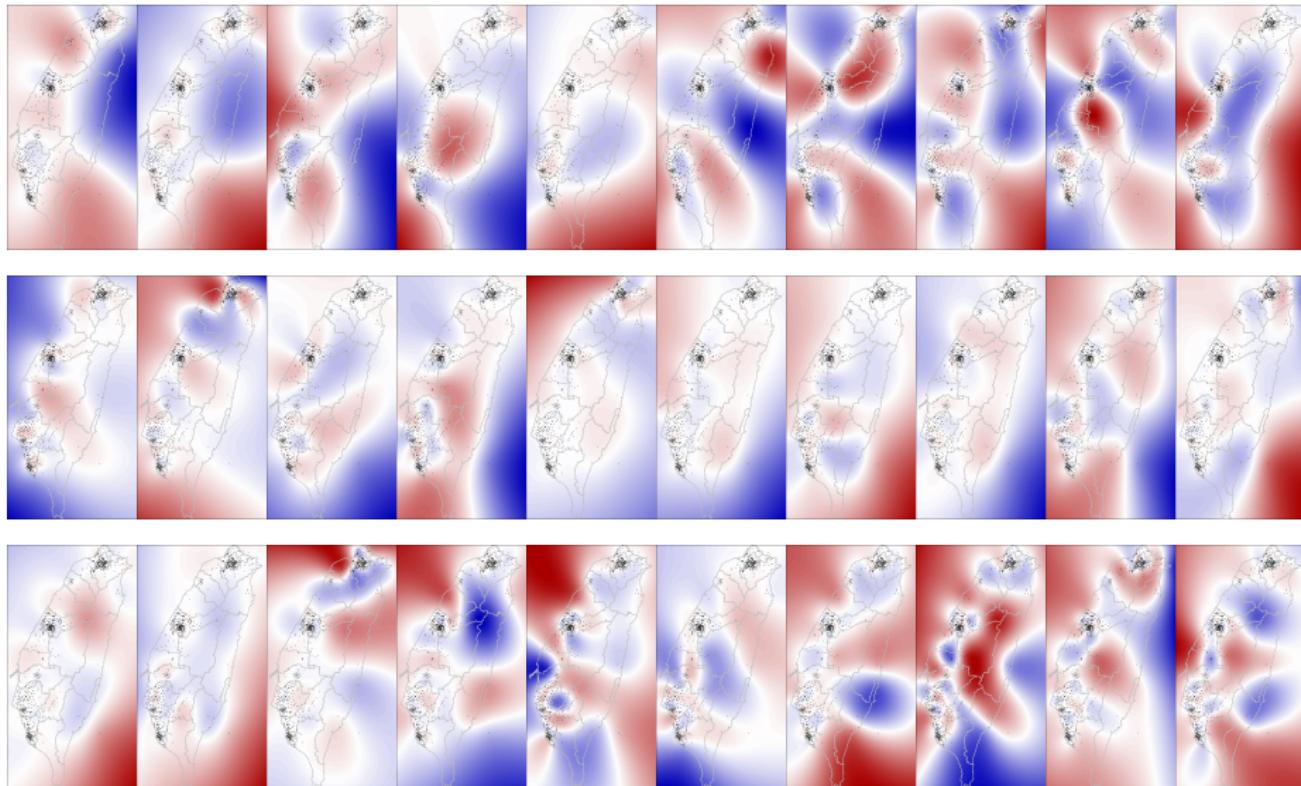
$$\varphi_k(\mathbf{s}) = \arg \min_{g(\cdot)} \{ J(g) : \mathbf{g}' \boldsymbol{\phi}_1 = \dots = \mathbf{g}' \boldsymbol{\phi}_{k-1} = 0, \|\mathbf{g}\| = 1 \};$$

$$k = d+2, \dots, n$$

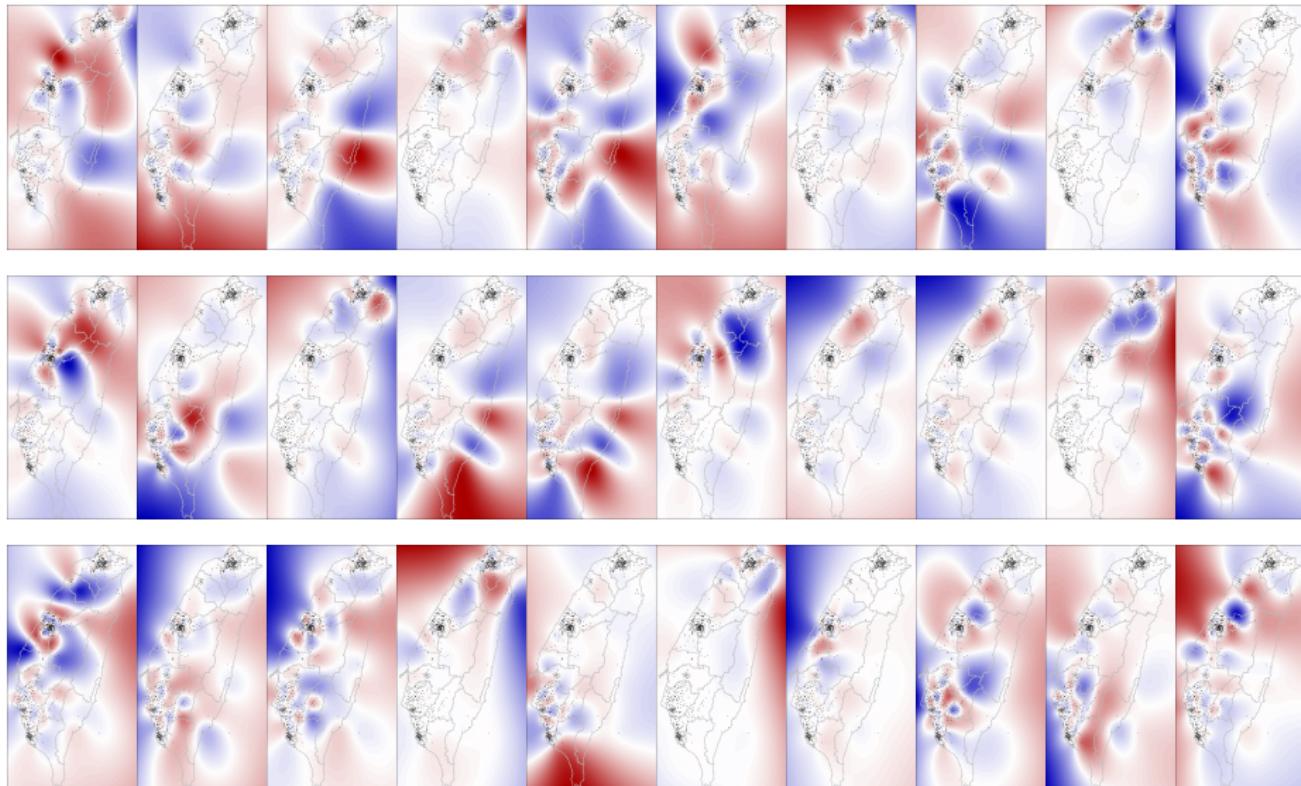
- $\boldsymbol{\phi}_j = (\varphi_j(\mathbf{s}_1), \dots, \varphi_j(\mathbf{s}_n))'$  and  $\mathbf{g} = (g(\mathbf{s}_1), \dots, g(\mathbf{s}_n))'$
- $0 = J(\varphi_1) = \dots = J(\varphi_{d+1}) < J(\varphi_{d+2}) \leq \dots \leq J(\varphi_n)$



# Basis Functions: $\varphi_{14}(\cdot), \dots, \varphi_{43}(\cdot)$



# Basis Functions: $\varphi_{44}(\cdot), \dots, \varphi_{73}(\cdot)$



# Improvement 2: Multi-Resolution Basis Functions

## Proposition

$$\varphi_k(\mathbf{s}) = \begin{cases} 1; & k=1 \\ x_k; & k=2, \dots, d+1 \\ \frac{1}{\lambda_{k-d-1}} \{ \psi(\mathbf{s}) - \Psi' \mathbf{X} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{x} \}' \mathbf{v}_{k-d-1}; & k=d+2, \dots, n \end{cases}$$

- $\mathbf{x} = (1, x_1, \dots, x_d)'$
- $\mathbf{X}$  is an  $n \times (d + 1)$  matrix with the  $i$ -th row  $(1, x_{i1}, \dots, x_{id})$  corresponding to  $\mathbf{s}_i$
- $\lambda_k$  and  $\mathbf{v}_k$  are the  $k$ -th eigenvalue and eigenvector of  $\mathbf{Q} \Psi \mathbf{Q}$  with  $\lambda_1 \geq \dots \geq \lambda_n$
- $\mathbf{Q} = \mathbf{I} - \mathbf{X} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}'$  and  $\Psi = (\psi_j(\mathbf{s}_i))_{n \times n}$
- $\psi(\mathbf{s}) = (\psi_1(\mathbf{s}), \dots, \psi_n(\mathbf{s}))'$
- $\psi_j(\mathbf{s}) = \frac{1}{8\pi} \|\mathbf{s} - \mathbf{s}_j\|^2 \log (\|\mathbf{s} - \mathbf{s}_j\|), \quad \text{if } d = 2$

## Improvement 3: Selection of $K$

- Akaike's information criterion (Akaike, 1973)

$$\begin{aligned} \text{AIC}(K) &= -2 \underbrace{\ell(\hat{\mathbf{M}}, \hat{\sigma}^2)}_{\text{log-likelihood}} + 2 \underbrace{\text{df}(K, T)}_{\text{\# free parameters}} \\ &= \frac{T \text{tr}(\mathbf{S})}{\hat{\sigma}^2} + T \sum_{k=1}^K \left\{ \log (\hat{d}_k + \hat{\sigma}^2) - \frac{d_k \hat{d}_k}{\hat{\sigma}^2 (\hat{d}_k + \hat{\sigma}^2)} \right\} \\ &\quad + 2 \text{df}(K, T) \end{aligned}$$

$$\bullet \text{df}(K, T) = \begin{cases} K(K+1)/2 + 1; & \text{if } K \leq T \\ KT - T(T-1)/2 + 1; & \text{if } K > T \end{cases}$$

- $\hat{K} = \arg \min_{d+1 \leq K \leq K^*} \text{AIC}(K)$

# Resolution Adaptive FRK

- Automatic FRK (Tzeng and Huang, 2017+)

$$\hat{y}_t(\mathbf{s}) = \varphi(\mathbf{s})' \hat{\mathbf{M}} \Phi' (\Phi \hat{\mathbf{M}} \Phi' + \hat{\sigma}^2 \mathbf{I})^{-1} \mathbf{Z}_t,$$

$$(\Phi \hat{\mathbf{M}} \Phi' + \hat{\sigma}^2 \mathbf{I})^{-1} = \frac{1}{\hat{\sigma}^2} \mathbf{I}_n - \frac{1}{\hat{\sigma}^4} \Phi \left\{ \hat{\mathbf{M}}^{-1} + \frac{1}{\hat{\sigma}^2} \Phi' \Phi \right\}^{-1} \Phi'$$

- $\Phi$ : simple closed-form expression (multiresolution spline basis functions)
- $\hat{\mathbf{M}}$  and  $\hat{\sigma}^2$ : simple closed-form expressions  $\Rightarrow$  allow large  $K$
- AIC( $K$ ): simple closed-form expression
- Numerically stable with no undesirable artifacts in predicted surfaces
- Directly applicable to sparse or irregularly spaced data

# Resolution Adaptive FRK

- Automatic FRK (Tzeng and Huang, 2017+)

$$\hat{y}_t(\mathbf{s}) = \varphi(\mathbf{s})' \hat{\mathbf{M}} \Phi' (\Phi \hat{\mathbf{M}} \Phi' + \hat{\sigma}^2 \mathbf{I})^{-1} \mathbf{Z}_t,$$

$$(\Phi \hat{\mathbf{M}} \Phi' + \hat{\sigma}^2 \mathbf{I})^{-1} = \frac{1}{\hat{\sigma}^2} \mathbf{I}_n - \frac{1}{\hat{\sigma}^4} \Phi \left\{ \hat{\mathbf{M}}^{-1} + \frac{1}{\hat{\sigma}^2} \Phi' \Phi \right\}^{-1} \Phi'$$

- $\Phi$ : simple closed-form expression (multiresolution spline basis functions)
- $\hat{\mathbf{M}}$  and  $\hat{\sigma}^2$ : simple closed-form expressions  $\Rightarrow$  allow large  $K$
- AIC( $K$ ): simple closed-form expression
- Numerically stable with no undesirable artifacts in predicted surfaces
- Directly applicable to sparse or irregularly spaced data

# Extensions

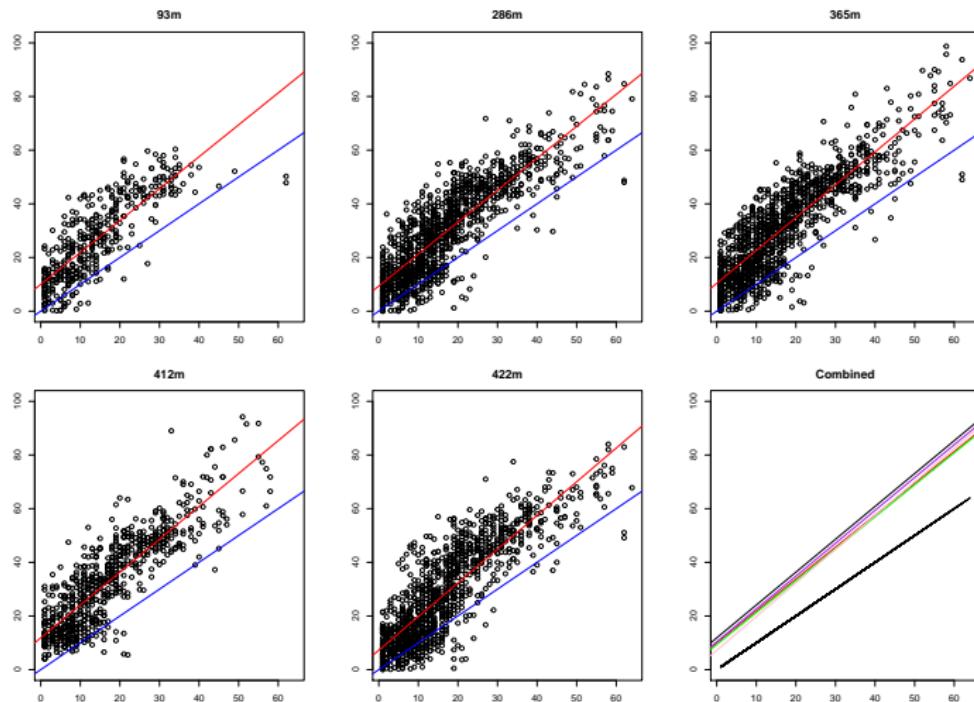
# Spatial Random Effects Model with Nonzero Means

- Spatial PM<sub>2.5</sub> process at time  $t = 1, \dots, T$

$$y_t(\mathbf{s}) = \mathbf{w}'_t \varphi(\mathbf{s}) = \sum_{k=1}^K w_{t,k} \varphi_k(\mathbf{s}); \quad \mathbf{s} \in D$$

- $\mathbf{w}_t = (w_{t,1}, \dots, w_{t,K})' \sim N(\mu, \mathbf{M})$
- $\varphi_1(\mathbf{s}), \dots, \varphi_K(\mathbf{s})$  is an order set of multi-resolution basis functions
- Parameters:  $\mu$ ,  $\mathbf{M}$  and  $\sigma^2$ 
  - Given  $\mu$ , the ML estimators of  $\mathbf{M}$  and  $\sigma^2$  have simple closed-form expressions
  - Given  $\mathbf{M}$  and  $\sigma^2$ , the ML estimator of  $\mu$  has a simple closed-form expression

# Spatial Prediction Based on Two Types of Data



- x-axis: PM<sub>2.5</sub> measurements at an EPA monitoring station
- y-axis: PM<sub>2.5</sub> measurements at nearby AirBoxes

# Spatial Prediction Based on Two Types of Data

$$y_t(\mathbf{s}) = \mathbf{w}'_t \varphi(\mathbf{s}), \quad \mathbf{w}_t \sim N(\mu, \mathbf{M})$$

- Data

- Airboxes:  $\mathbf{Z}_t = (z_t(\mathbf{s}_1), \dots, z_t(\mathbf{s}_m))'$
- Monitoring stations:  $\mathbf{Z}_t^* = (z_t^*(\mathbf{s}_{m+1}), \dots, z_t^*(\mathbf{s}_n))'$

- Measurement equations

$$z_t(\mathbf{s}_i) = y_t(\mathbf{s}_i) + \varepsilon_t(\mathbf{s}_i); \quad i = 1, \dots, m,$$

$$z_t^*(\mathbf{s}_j) = \alpha_0 + \alpha_1 y_t(\mathbf{s}_i) + \varepsilon_t^*(\mathbf{s}_i); \quad i = m+1, \dots, n$$

- $\varepsilon_t(\mathbf{s}_1), \dots, \varepsilon_t(\mathbf{s}_m) \sim N(0, \sigma^2)$
- $\varepsilon_t^*(\mathbf{s}_{m+1}), \dots, \varepsilon_t^*(\mathbf{s}_n) \sim N(0, r\sigma^2)$  with  $r < 1$

- Parameters:  $\alpha_0, \alpha_1, r, \mathbf{M}$  and  $\sigma^2$

- Given  $\alpha_1$  and  $r$ , the ML estimators of  $\alpha_0$ ,  $\mathbf{M}$  and  $\sigma^2$  have closed-form expressions

# Spatial Prediction Based on Two Types of Data

$$y_t(\mathbf{s}) = \mathbf{w}'_t \varphi(\mathbf{s}), \quad \mathbf{w}_t \sim N(\mu, \mathbf{M})$$

- Data

- Airboxes:  $\mathbf{Z}_t = (z_t(\mathbf{s}_1), \dots, z_t(\mathbf{s}_m))'$
- Monitoring stations:  $\mathbf{Z}_t^* = (z_t^*(\mathbf{s}_{m+1}), \dots, z_t^*(\mathbf{s}_n))'$

- Measurement equations

$$z_t(\mathbf{s}_i) = y_t(\mathbf{s}_i) + \varepsilon_t(\mathbf{s}_i); \quad i = 1, \dots, m,$$

$$z_t^*(\mathbf{s}_j) = \alpha_0 + \alpha_1 y_t(\mathbf{s}_i) + \varepsilon_t^*(\mathbf{s}_i); \quad i = m+1, \dots, n$$

- $\varepsilon_t(\mathbf{s}_1), \dots, \varepsilon_t(\mathbf{s}_m) \sim N(0, \sigma^2)$
- $\varepsilon_t^*(\mathbf{s}_{m+1}), \dots, \varepsilon_t^*(\mathbf{s}_n) \sim N(0, r\sigma^2)$  with  $r < 1$

- Parameters:  $\alpha_0, \alpha_1, r, \mathbf{M}$  and  $\sigma^2$

- Given  $\alpha_1$  and  $r$ , the ML estimators of  $\alpha_0$ ,  $\mathbf{M}$  and  $\sigma^2$  have closed-form expressions

# Spatial Prediction Based on Two Types of Data

$$y_t(\mathbf{s}) = \mathbf{w}'_t \varphi(\mathbf{s}), \quad \mathbf{w}_t \sim N(\mu, \mathbf{M})$$

- Data

- Airboxes:  $\mathbf{Z}_t = (z_t(\mathbf{s}_1), \dots, z_t(\mathbf{s}_m))'$
- Monitoring stations:  $\mathbf{Z}_t^* = (z_t^*(\mathbf{s}_{m+1}), \dots, z_t^*(\mathbf{s}_n))'$

- Measurement equations

$$z_t(\mathbf{s}_i) = y_t(\mathbf{s}_i) + \varepsilon_t(\mathbf{s}_i); \quad i = 1, \dots, m,$$

$$z_t^*(\mathbf{s}_j) = \alpha_0 + \alpha_1 y_t(\mathbf{s}_i) + \varepsilon_t^*(\mathbf{s}_i); \quad i = m+1, \dots, n$$

- $\varepsilon_t(\mathbf{s}_1), \dots, \varepsilon_t(\mathbf{s}_m) \sim N(0, \sigma^2)$
- $\varepsilon_t^*(\mathbf{s}_{m+1}), \dots, \varepsilon_t^*(\mathbf{s}_n) \sim N(0, r\sigma^2)$  with  $r < 1$

- Parameters:  $\alpha_0, \alpha_1, r, \mathbf{M}$  and  $\sigma^2$

- Given  $\alpha_1$  and  $r$ , the ML estimators of  $\alpha_0$ ,  $\mathbf{M}$  and  $\sigma^2$  have closed-form expressions

# Spatial Prediction Based on Two Types of Data

$$y_t(\mathbf{s}) = \mathbf{w}'_t \varphi(\mathbf{s}), \quad \mathbf{w}_t \sim N(\mu, \mathbf{M})$$

- Data

- Airboxes:  $\mathbf{Z}_t = (z_t(\mathbf{s}_1), \dots, z_t(\mathbf{s}_m))'$
- Monitoring stations:  $\mathbf{Z}_t^* = (z_t^*(\mathbf{s}_{m+1}), \dots, z_t^*(\mathbf{s}_n))'$

- Measurement equations

$$z_t(\mathbf{s}_i) = y_t(\mathbf{s}_i) + \varepsilon_t(\mathbf{s}_i); \quad i = 1, \dots, m,$$

$$z_t^*(\mathbf{s}_j) = \alpha_0 + \alpha_1 y_t(\mathbf{s}_i) + \varepsilon_t^*(\mathbf{s}_i); \quad i = m + 1, \dots, n$$

- $\varepsilon_t(\mathbf{s}_1), \dots, \varepsilon_t(\mathbf{s}_m) \sim N(0, \sigma^2)$
- $\varepsilon_t^*(\mathbf{s}_{m+1}), \dots, \varepsilon_t^*(\mathbf{s}_n) \sim N(0, r\sigma^2)$  with  $r < 1$

- Parameters:  $\alpha_0, \alpha_1, r, \mathbf{M}$  and  $\sigma^2$

- Given  $\alpha_1$  and  $r$ , the ML estimators of  $\alpha_0$ ,  $\mathbf{M}$  and  $\sigma^2$  have closed-form expressions

# Spatial Prediction Based on Two Types of Data

- Minimum mean squared error predictor of  $y_t(\mathbf{s})$

$$\hat{y}_t(\mathbf{s}) = \varphi(\mathbf{s})' \boldsymbol{\mu} + \varphi(\mathbf{s})' \mathbf{M} \boldsymbol{\Phi}' \mathbf{A}' (\mathbf{A} \boldsymbol{\Phi} \mathbf{M} \boldsymbol{\Phi}' \mathbf{A}' + \mathbf{V})^{-1} \left( \mathbf{Z}_t^* - \alpha_0 \mathbf{1} - \alpha_1 \mathbf{f}_t^* \right)$$

- $\mathbf{f}_t = (\varphi(\mathbf{s}_1), \dots, \varphi(\mathbf{s}_m)) \boldsymbol{\mu}$
- $\mathbf{f}_t^* = (\varphi(\mathbf{s}_{m+1}), \dots, \varphi(\mathbf{s}_n)) \boldsymbol{\mu}$
- $\mathbf{A} = \begin{pmatrix} \mathbf{I}_m & \mathbf{0} \\ \mathbf{0} & \alpha_1 \mathbf{I}_{n-m} \end{pmatrix}$
- $\mathbf{V} = \begin{pmatrix} \sigma^2 \mathbf{I}_m & \mathbf{0} \\ \mathbf{0} & r\sigma^2 \mathbf{I}_{n-m} \end{pmatrix}$

- Sherman-Morrison-Woodbury formula

$$(\mathbf{A} \boldsymbol{\Phi} \mathbf{M} \boldsymbol{\Phi}' \mathbf{A}' + \mathbf{V})^{-1} = \mathbf{V}^{-1} - \mathbf{V}^{-1} \boldsymbol{\Phi} \{ \mathbf{M}^{-1} + \mathbf{A} \boldsymbol{\Phi}' \mathbf{V}^{-1} \mathbf{A} \boldsymbol{\Phi} \}^{-1} \boldsymbol{\Phi}' \mathbf{A}' \mathbf{V}^{-1}$$

# Other Extensions

- Some data are missing:  $\mathbf{Z}_t = \begin{pmatrix} \mathbf{z}_t^{(obs)} \\ \mathbf{z}_t^{(mis)} \end{pmatrix}$ 
  - Compute the ML estimators using the EM algorithm
  - Treat  $\mathbf{z}_t^{(mis)}$  as missing data directly instead of  $\mathbf{w}_t$
- Models with covariates

$$\mathbf{Z}_t \sim N(\mathbf{X}\boldsymbol{\beta}, \Phi\mathbf{M}\Phi' + \sigma^2 \mathbf{I}_n)$$

- $\mathbf{X}$ :  $n \times (p+1)$  design matrix of rank  $(p+1)$
- $\boldsymbol{\beta}$ : Regression parameter vector
- Restricted ML estimators of  $\mathbf{M}$  and  $\sigma^2$  have closed-form expressions

- Space-time models

$$y_t(\mathbf{s}) = \mathbf{w}_t' \varphi(\mathbf{s}), \quad \mathbf{w}_t \sim N(\boldsymbol{\mu}, \mathbf{M}),$$

$$\mathbf{w}_t = \mathbf{A}\mathbf{w}_{t-1} + \eta_t \quad (\text{multivariate stationary AR}(1)),$$

$$\mathbf{Z}_t = \Phi\mathbf{w}_t + \varepsilon_t, \quad \varepsilon_t \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$$

- Given  $\mathbf{A}$ , the ML estimators of  $\mathbf{M}$  and  $\sigma^2$  have closed-form expressions

# Other Extensions

- Some data are missing:  $\mathbf{Z}_t = \begin{pmatrix} \mathbf{Z}_t^{(obs)} \\ \mathbf{Z}_t^{(mis)} \end{pmatrix}$ 
  - Compute the ML estimators using the EM algorithm
  - Treat  $\mathbf{Z}_t^{(mis)}$  as missing data directly instead of  $\mathbf{w}_t$
- Models with covariates

$$\mathbf{Z}_t \sim N(\mathbf{X}\boldsymbol{\beta}, \Phi\mathbf{M}\Phi' + \sigma^2\mathbf{I}_n)$$

- $\mathbf{X}$ :  $n \times (p+1)$  design matrix of rank  $(p+1)$
- $\boldsymbol{\beta}$ : Regression parameter vector
- Restricted ML estimators of  $\mathbf{M}$  and  $\sigma^2$  have closed-form expressions

- Space-time models

$$y_t(\mathbf{s}) = \mathbf{w}_t'\varphi(\mathbf{s}), \quad \mathbf{w}_t \sim N(\boldsymbol{\mu}, \mathbf{M}),$$

$$\mathbf{w}_t = \mathbf{A}\mathbf{w}_{t-1} + \eta_t \quad (\text{multivariate stationary AR}(1)),$$

$$\mathbf{Z}_t = \Phi\mathbf{w}_t + \varepsilon_t, \quad \varepsilon_t \sim N(\mathbf{0}, \sigma^2\mathbf{I})$$

- Given  $\mathbf{A}$ , the ML estimators of  $\mathbf{M}$  and  $\sigma^2$  have closed-form expressions

# Other Extensions

- Some data are missing:  $\mathbf{Z}_t = \begin{pmatrix} \mathbf{Z}_t^{(obs)} \\ \mathbf{Z}_t^{(mis)} \end{pmatrix}$ 
  - Compute the ML estimators using the EM algorithm
  - Treat  $\mathbf{Z}_t^{(mis)}$  as missing data directly instead of  $\mathbf{w}_t$
- Models with covariates

$$\mathbf{Z}_t \sim N(\mathbf{X}\boldsymbol{\beta}, \Phi\mathbf{M}\Phi' + \sigma^2\mathbf{I}_n)$$

- $\mathbf{X}$ :  $n \times (p+1)$  design matrix of rank  $(p+1)$
  - $\boldsymbol{\beta}$ : Regression parameter vector
  - Restricted ML estimators of  $\mathbf{M}$  and  $\sigma^2$  have closed-form expressions
- Space-time models

$$y_t(\mathbf{s}) = \mathbf{w}'_t \varphi(\mathbf{s}), \quad \mathbf{w}_t \sim N(\boldsymbol{\mu}, \mathbf{M}),$$

$$\mathbf{w}_t = \mathbf{A}\mathbf{w}_{t-1} + \boldsymbol{\eta}_t \quad (\text{multivariate stationary AR(1)}),$$

$$\mathbf{Z}_t = \Phi\mathbf{w}_t + \boldsymbol{\varepsilon}_t, \quad \boldsymbol{\varepsilon}_t \sim N(\mathbf{0}, \sigma^2\mathbf{I})$$

- Given  $\mathbf{A}$ , the ML estimators of  $\mathbf{M}$  and  $\sigma^2$  have closed-form expressions

## PM<sub>2.5</sub> Applications

# Spatial Prediction of PM<sub>2.5</sub>

# Detect Abnormal AirBox Measurements

- AirBoxes closer together in space are likely to produce similar measurements
- $\hat{y}_t(\mathbf{s})$ : Predicted surface obtained from the proposed fixed rank kriging
- Detect outliers by monitoring standardized residuals

$$r_t(\mathbf{s}_i) = \frac{z_t(\mathbf{s}_i) - \hat{y}_t(\mathbf{s}_i)}{v_t}; \quad i = 1, \dots, n, t = 1, \dots, T$$

- Monitor excessively large values of  $r_t^2(\mathbf{s}_i)$ 
  - Small (negative)  $r_t(\mathbf{s}_i)$  corresponds to the AirBoxes that may be put indoors at  $\mathbf{s}_i$
  - Large (positive)  $r_t(\mathbf{s}_i)$  corresponds to potential emission sources of PM<sub>2.5</sub> around location  $\mathbf{s}_i$  and time  $t$

# Detect Abnormal AirBox Measurements

- AirBoxes closer together in space are likely to produce similar measurements
- $\hat{y}_t(\mathbf{s})$ : Predicted surface obtained from the proposed fixed rank kriging
- Detect outliers by monitoring standardized residuals

$$r_t(\mathbf{s}_i) = \frac{z_t(\mathbf{s}_i) - \hat{y}_t(\mathbf{s}_i)}{v_t}; \quad i = 1, \dots, n, t = 1, \dots, T$$

- Monitor excessively large values of  $r_t^2(\mathbf{s}_i)$ 
  - Small (negative)  $r_t(\mathbf{s}_i)$  corresponds to the AirBoxes that may be put indoors at  $\mathbf{s}_i$
  - Large (positive)  $r_t(\mathbf{s}_i)$  corresponds to potential emission sources of PM<sub>2.5</sub> around location  $\mathbf{s}_i$  and time  $t$

# Detect Abnormal AirBox Measurements

- AirBoxes closer together in space are likely to produce similar measurements
- $\hat{y}_t(\mathbf{s})$ : Predicted surface obtained from the proposed fixed rank kriging
- Detect outliers by monitoring standardized residuals

$$r_t(\mathbf{s}_i) = \frac{z_t(\mathbf{s}_i) - \hat{y}_t(\mathbf{s}_i)}{v_t}; \quad i = 1, \dots, n, t = 1, \dots, T$$

- Monitor excessively large values of  $r_t^2(\mathbf{s}_i)$ 
  - Small (negative)  $r_t(\mathbf{s}_i)$  corresponds to the AirBoxes that may be put indoors at  $\mathbf{s}_i$
  - Large (positive)  $r_t(\mathbf{s}_i)$  corresponds to potential emission sources of PM<sub>2.5</sub> around location  $\mathbf{s}_i$  and time  $t$

# R Package

- `autoFRK`: To be submitted to CRAN