

Real-time PM_{2.5} mapping and anomaly detection from AirBoxes in Taiwan

G. Huang¹ L.-J. Chen² W.-H. Hwang³ S. Tzeng⁴ H.-C. Huang⁵

Introduction

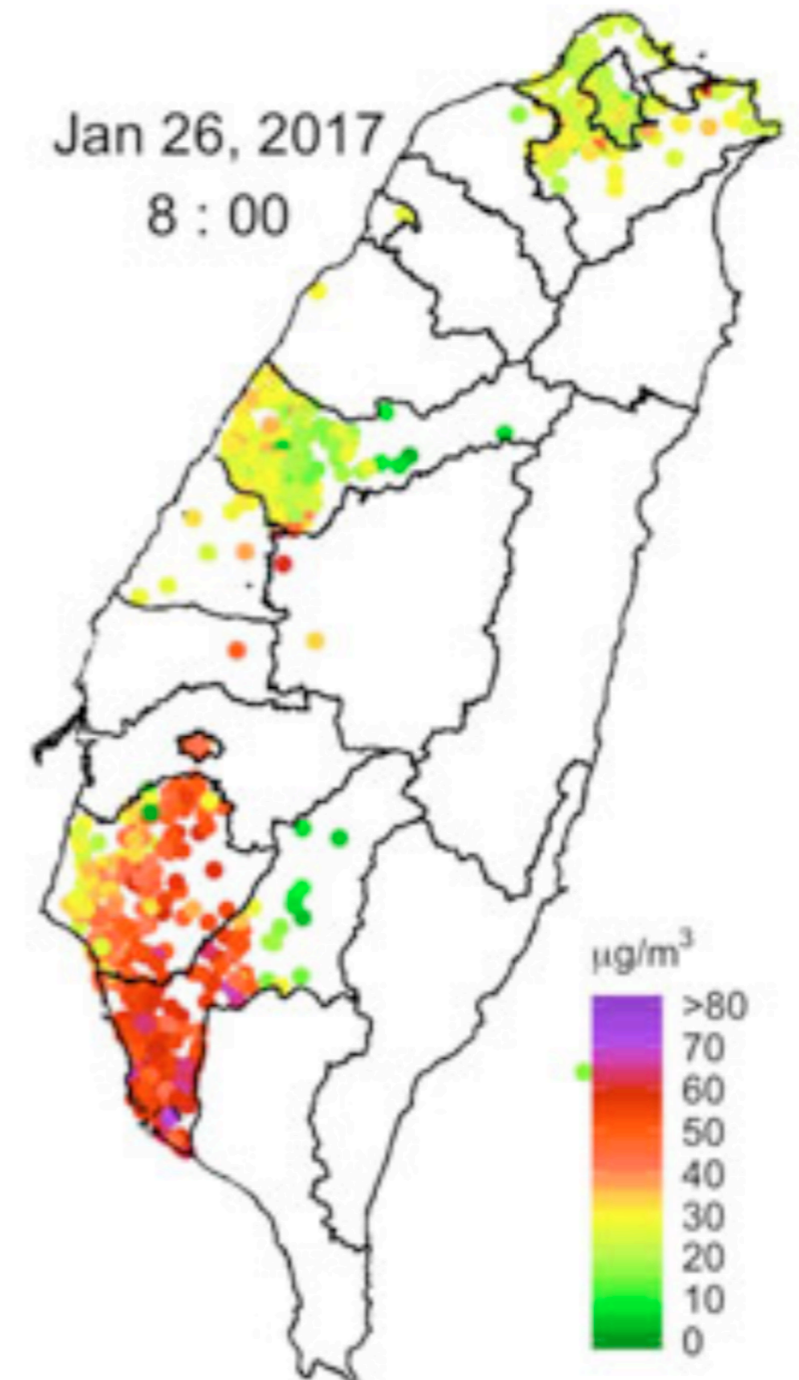
- What is PM2.5
- What is AirBox
- Anomaly detection
 - Clustering-based methods
 - Machine-learning methods
 - Statistical methods

Contribution

- Real time PM2.5 concentration at any location with its estimation error bar
 - Kriging approach
- A spatio-temporal control chart that can automatically monitor anomalous measurements by utilizing neighboring AirBox information.

AirBox Data

- 1283 AirBoxes across Taiwan from Jan. 1 to Feb. 28
- Precision is 111m * 102m
- Aggregate the data into hourly data at each location using the simple average
- Few unusual measurement that are either much higher or lower than nearby



Methodologies

$$y(\mathbf{s}) = \beta' \phi(\mathbf{s}) + \eta(\mathbf{s}) = \sum_{k=1}^K \beta_k \phi_k(\mathbf{s}) + \eta(\mathbf{s}); \quad \mathbf{s} \in D,$$

↙ regression coefficients
↑

↖ first K multiresolution spline basis functions
↑ zero-mean spatial process with an exponential covariance function

$$C(\mathbf{s}, \mathbf{u}) = \text{cov}(y(\mathbf{s}), y(\mathbf{u})) = v_y^2 \exp(-\|\mathbf{s} - \mathbf{u}\|/\lambda)$$

- Suppose data $\mathbf{z} \equiv (z(\mathbf{s}_1), \dots, z(\mathbf{s}_n))'$
- With additive white noise $\mathbf{z} = \mathbf{y} + \boldsymbol{\varepsilon}$
 where $\mathbf{y} = (y(\mathbf{s}_1), \dots, y(\mathbf{s}_n))'$, $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, v_\varepsilon^2 \mathbf{I})$

Robust Method

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^K} \sum_{i=1}^n \rho \left(\frac{e_i}{\hat{\sigma}} \right)$$

$$e_i = z(\mathbf{s}_i) - \beta' \phi(\mathbf{s}_i)$$

- Where

$$\rho(x) = \begin{cases} x^2/2, & \text{if } |x| \leq c \\ c|x| - c^2/2, & \text{if } |x| > c \end{cases}$$

Choose $c = 1.345$ which gives an efficiency of 95% if the error are normal distributed

$$\hat{y}(\mathbf{s}_0) = \hat{\beta}' \phi(\mathbf{s}_0) + \mathbf{c}(\mathbf{s}_0; \hat{\theta})' \left(\Sigma(\hat{\theta}) + \hat{v}_\epsilon^2 \mathbf{I} \right)^{-1} \left(\mathbf{z} - \left(\hat{\beta}' \phi(\mathbf{s}_1), \dots, \hat{\beta}' \phi(\mathbf{s}_n) \right)' \right)$$

Anomaly Detection

- Baseline: $\{\hat{y}_t(\mathbf{s}) : \mathbf{s} \in D\}$
- Standardized residuals $r_t(\mathbf{s}_i) = \frac{z_t(\mathbf{s}_i) - \hat{y}_t(\mathbf{s}_i)}{\hat{\sigma}_t(\mathbf{s}_i)}$; $i = 1, \dots, n, t = 1, \dots, T$

$r_t(\mathbf{s}_i)$ has normal distribution if the parameter are known

High positive r \longrightarrow It is higher than neighbor observation

Low negative r \longrightarrow It is lower than neighbor observation

do not need to specify a specific neighborhood range

Anomaly Detection

- Control chart of each AirBox: $r_t(\mathbf{s}_i); i = 1, \dots, T$
- Control limits: 3 standard deviation: $|r_t(\mathbf{s}_i)| > 3$
- Ranking: $\text{RMSE}_i = \left\{ \frac{1}{|T_i|} \sum_{t \in T_i} (r_t(\mathbf{s}_i))^2 \right\}^{1/2}$
- AirBoxes with high RMSE indicate that they tend to produce outlying observations

Decompose RMSE

$$\text{RMSE}_i = (b_i^2 + V_i)^{1/2},$$

$$b_i = \frac{1}{|T_i|} \sum_{t \in T_i} r_t(\mathbf{s}_i),$$

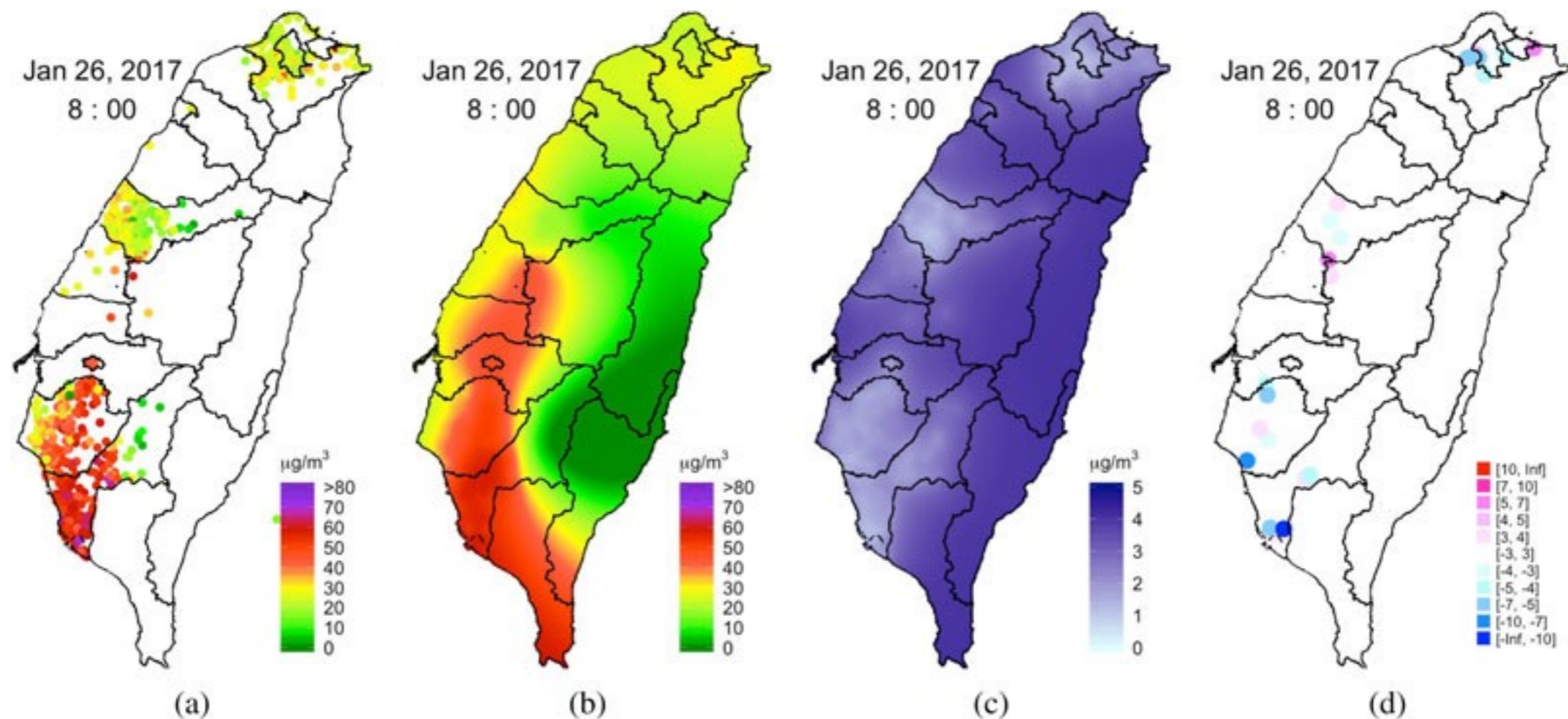
$$V_i = \frac{1}{|T_i|} \sum_{t \in T_i} \left\{ r_t(\mathbf{s}_i) - \frac{1}{|T_i|} \sum_{t \in T_i} r_t(\mathbf{s}_i) \right\}^2$$

- Classification

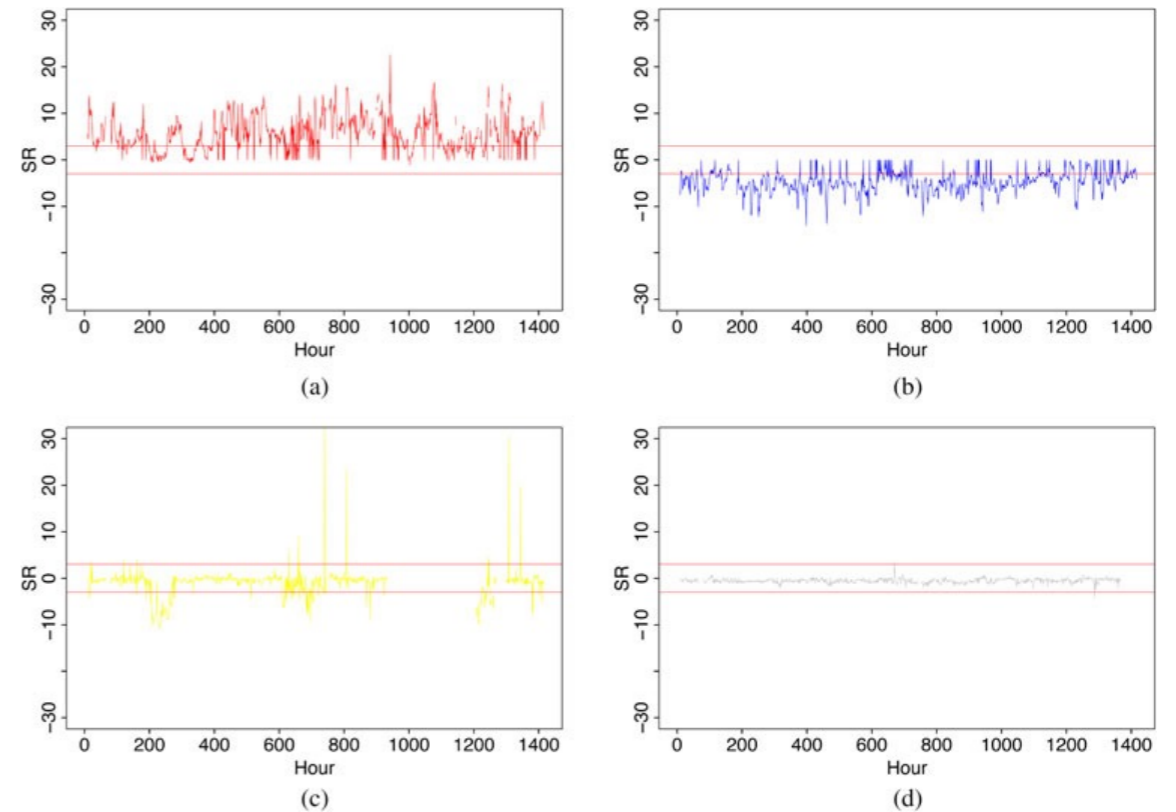
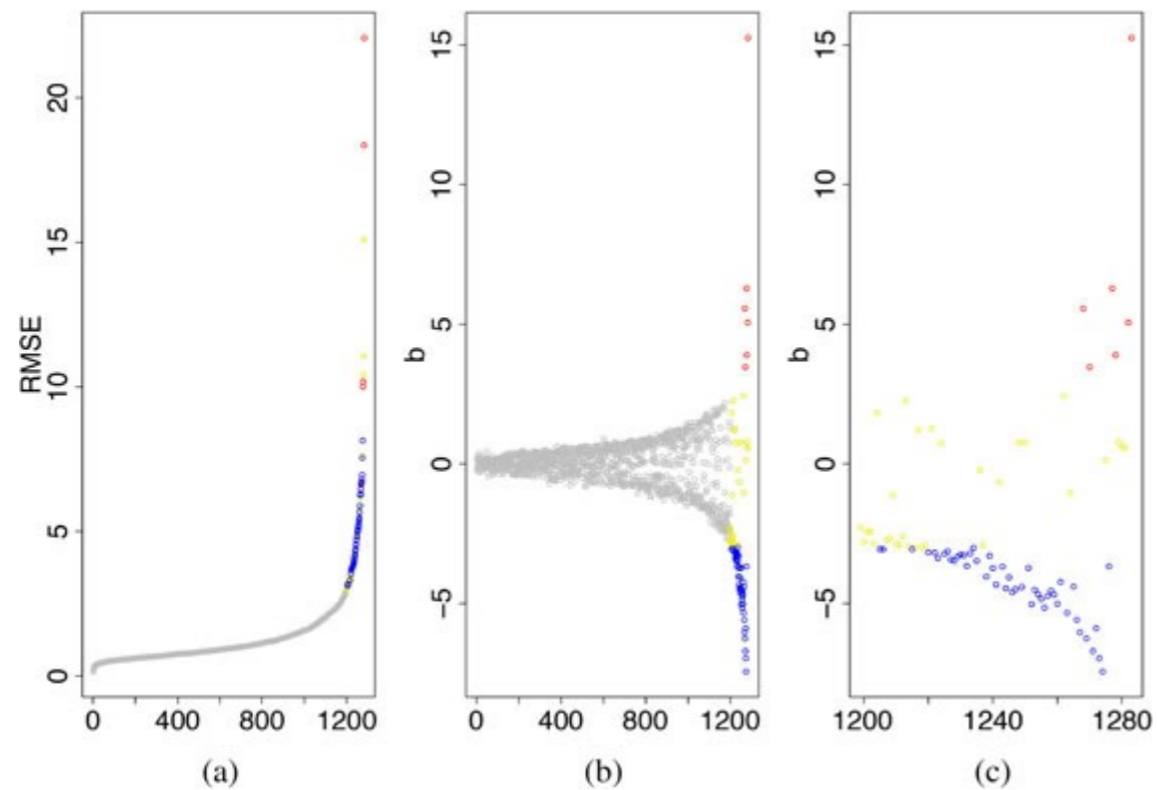
High RMSE	
Bi <= -3	Average standardize residual is above control limit
Bi >= 3	Average standardize residual falls below control limit
-3 < Bi < 3	Tends to have high variation

Analysis Result

- Locations with fewer AirBox has higher error values
- Unusual large or small PM2.5 value is shown in fig. D



Anomaly Detection Result



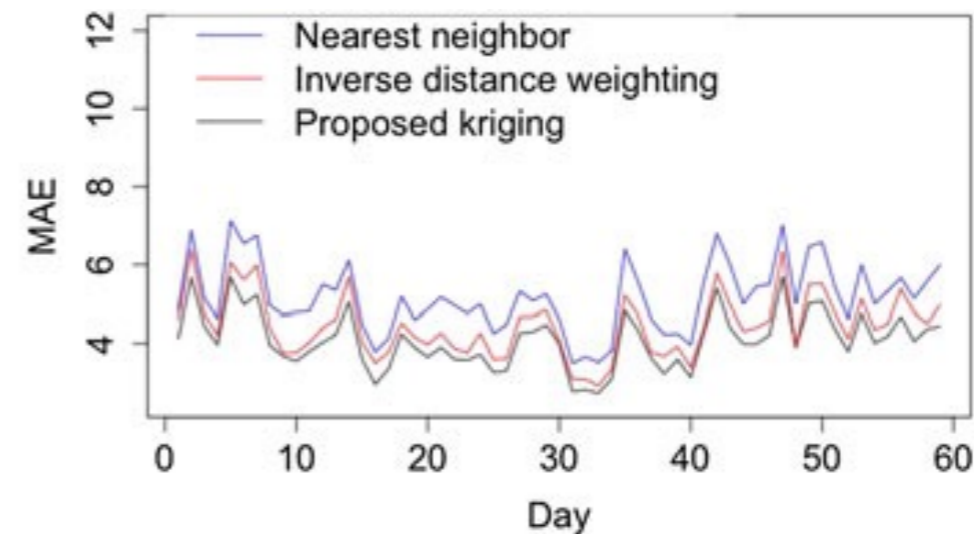
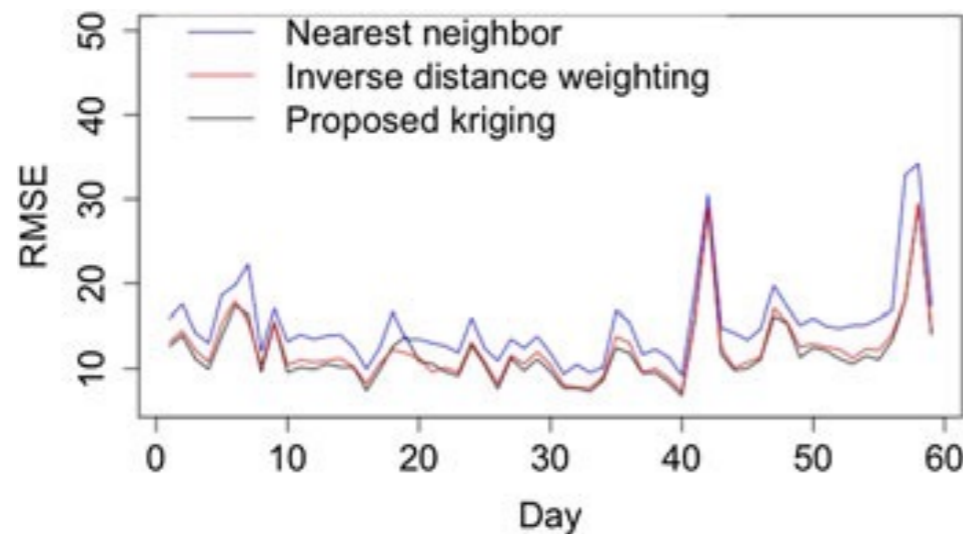
- (a) Ranking RMSE value
- (b) The corresponding biases of (a)
- (c) Classify high RMSE into 3 groups

Compare with Other Method

- Criteria

$$\text{RMSE} = \left\{ \frac{1}{100} \sum_{i=1}^{100} (\tilde{y}(\mathbf{s}_i^*) - z(\mathbf{s}_i^*))^2 \right\}^{1/2},$$

$$\text{MAE} = \text{median} \left\{ \left| \tilde{y}(\mathbf{s}_1^*) - z(\mathbf{s}_1^*) \right|, \dots, \left| \tilde{y}(\mathbf{s}_{100}^*) - z(\mathbf{s}_{100}^*) \right| \right\}$$



Method	Averaged RMSE	Averaged MAE
Nearest neighbor	15.16	5.18
Inverse distance weighting	12.37	4.46
The proposed kriging	11.88	4.09

Conclusion

- Proposed method is able to detect potential emission sources, malfunctioned AirBoxes, and AirBoxes that are put indoors
- AirBoxes provides very high spatial and temporal coverage but they have much more higher error than those large monitoring stations