



GLS-SOD: A Generalized Local Statistical Approach for Spatial Outlier Detection

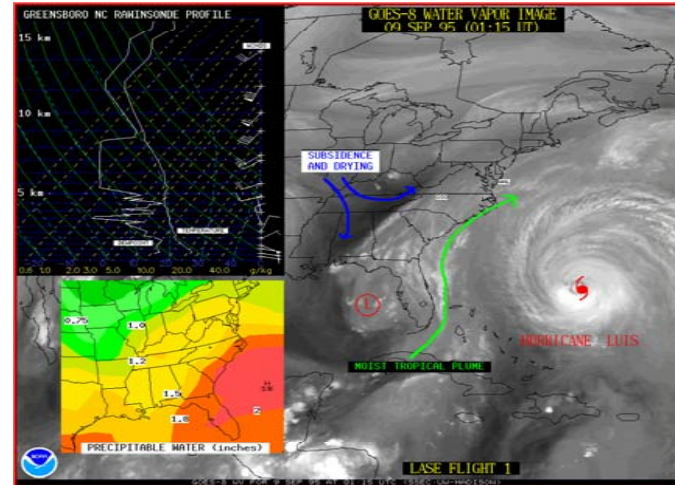
Feng Chen, Chang-Tien Lu, Arnold P. Boedihardjo
Virginia Tech, Computer Science Department
July, 28, 2010



Outline

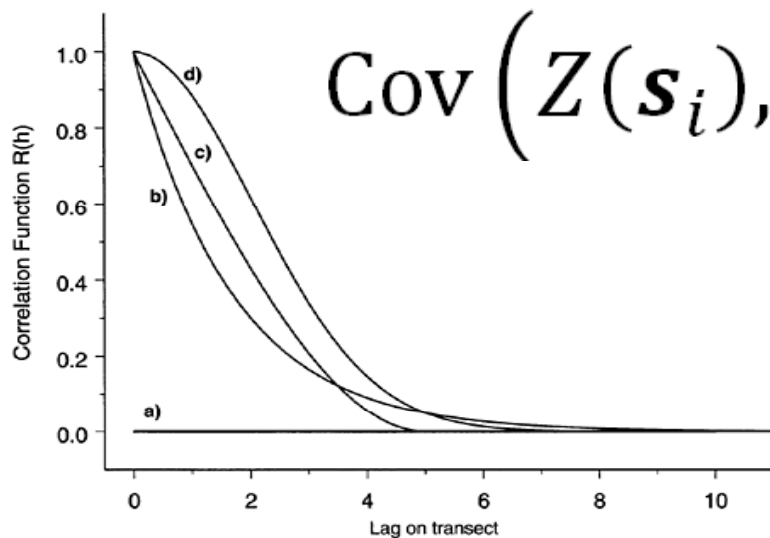
- Motivations
- Generalized Local Statistical (GLS) Model
- GLS Robust Estimation & Inferences
- Simulations
- Summary

Applications



What's Special for Spatial Data?

- By the first law of Geography, "*Everything is related to everything else, but nearby things are more related than distant things*" [Tobler 79]
- Spatial autocorrelation
 - Correlation of a variable with itself through space



$$\text{Cov} \left(Z(\mathbf{s}_i), Z(\mathbf{s}_j) \right) = f \left(\|\mathbf{s}_i - \mathbf{s}_j\|; \boldsymbol{\theta} \right)$$

Global based Spatial Outlier Detection

Given a partial sample $\{Z(\mathbf{s}_i)\}_{i=1}^n$ of a Gaussian random field, let

$$\mathbf{Z} = [Z(\mathbf{s}_1), \dots, Z(\mathbf{s}_n)]^T, \mathbf{X} = [\mathbf{x}(\mathbf{s}_1), \dots, \mathbf{x}(\mathbf{s}_n)]^T$$

$$\mathbf{Z} \sim N(\mathbf{X}\boldsymbol{\beta}, \boldsymbol{\Sigma}(\boldsymbol{\theta}))$$

$$\min_{\boldsymbol{\beta}, \boldsymbol{\theta}} \{(\mathbf{Z} - \mathbf{X}\boldsymbol{\beta})^T [\boldsymbol{\Sigma}(\boldsymbol{\theta})]^{-1} (\mathbf{Z} - \mathbf{X}\boldsymbol{\beta})\}$$

subject to some constraints on $\boldsymbol{\theta}$.

Standardized Z Test for an observation $Z(\mathbf{s})$.

First calculate $E[Z(\mathbf{s})] = \mathbf{x}(\mathbf{s})^T \boldsymbol{\beta}$, $\text{Var}(Z(\mathbf{s})) = \sigma_0^2$, $\text{Cov}[\mathbf{Z}, Z(\mathbf{s})] = \boldsymbol{\sigma}$,

$\text{Var}[\mathbf{Z}] = \boldsymbol{\Sigma}$, and $E[\mathbf{Z}(\mathbf{s})] = \boldsymbol{\mu}$. Then

$$\frac{Z(\mathbf{s}) - \mathbf{x}(\mathbf{s})^T \boldsymbol{\beta} - \boldsymbol{\sigma}^T \boldsymbol{\Sigma}^{-1} (\mathbf{Z} - \mathbf{X}\boldsymbol{\beta})}{\sqrt{\sigma_0^2 - \boldsymbol{\sigma}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\sigma}}} \geq \Phi^{-1} \left(\frac{\alpha}{2} \right), \text{ e. x. } \alpha = 0.05$$

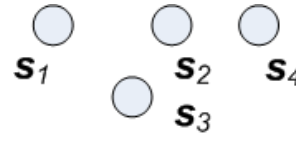
Local based Spatial Outlier Detection

$$\mathbf{Z} \sim N(\mathbf{X}\boldsymbol{\beta}, \boldsymbol{\Sigma}(\boldsymbol{\theta}))$$

↓ Local (laplacian) smoothing

$$\mathbf{FZ} \sim N(\mathbf{FX}\boldsymbol{\beta}, \mathbf{F}\boldsymbol{\Sigma}(\boldsymbol{\theta})\mathbf{F}^T)$$

$$\mathbf{FZ} \sim N(\mu \cdot \mathbf{I}_{n \times 1}, \sigma^2 \mathbf{I}_{n \times n})$$

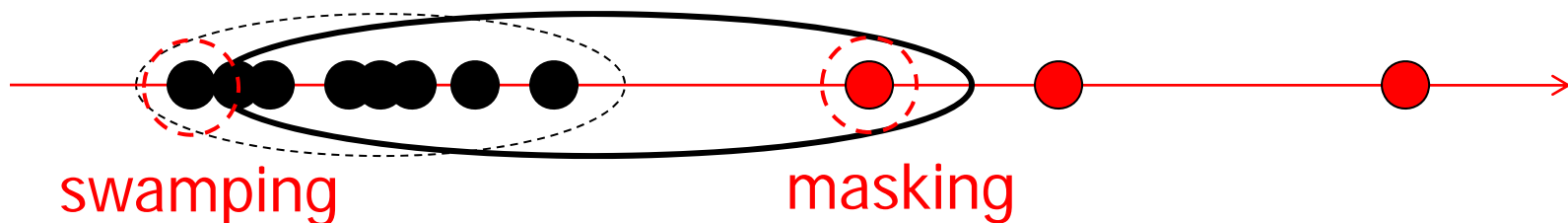


$$\mathbf{F} = \begin{bmatrix} 1 & -1/2 & -1/2 & 0 \\ 0 & 1 & -1/2 & -1/2 \\ -1/2 & -1/2 & 1 & 0 \\ 0 & -1/2 & -1/2 & 1 \end{bmatrix}$$

$$\mathbf{FX}\boldsymbol{\beta} \approx \mu \cdot \mathbf{I}_{n \times 1}$$

$$\mathbf{F}\boldsymbol{\Sigma}(\boldsymbol{\theta})\mathbf{F}^T \approx \sigma^2 \mathbf{I}_{n \times n}$$

$$\text{Z Test: } \frac{(\mathbf{Z}(\mathbf{s}) - \mu)}{\sigma} \geq \Phi^{-1}\left(\frac{\alpha}{2}\right), \text{ e.x. } \alpha = 0.05$$





Why GLS-SOD?

- Global based

- Pros: high accuracy with statistical justifications.
- Cons: very slow; complicate estimation process; non-convex optimization

- Local based

- Pros: very fast; simplicity
- Cons: heuristic-driven, lack of statistical justifications

- Questions

- Local (laplacian) smoothing vs. spatial dependence?
- Statistical connections between local and global methods?
- When will existing local based methods perform poorly and how to handle these situations?

Gaussian Random Field

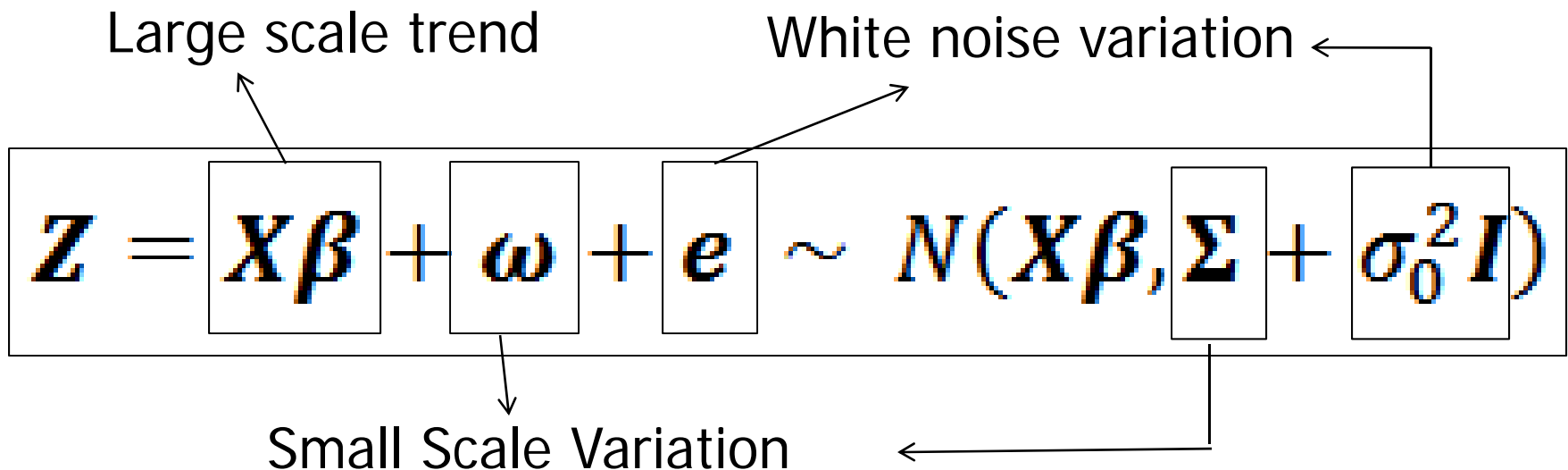
Consider a Gaussian Random Field $\{Z(\mathbf{s}), \mathbf{s} \in D \subset \mathbb{R}^2\}$

$$Z(\mathbf{s}) = \mathbf{x}(\mathbf{s})^T \boldsymbol{\beta} + \omega(\mathbf{s}) + \epsilon(\mathbf{s})$$

Given a partial sample $\{Z(\mathbf{s}_i)\}_{i=1}^n$ of the random field, let

$$\mathbf{Z} = [Z(\mathbf{s}_1), \dots, Z(\mathbf{s}_n)]^T, \mathbf{X} = [\mathbf{x}(\mathbf{s}_1), \dots, \mathbf{x}(\mathbf{s}_n)]^T$$

$$\mathbf{e} = [e(\mathbf{s}_1), \dots, e(\mathbf{s}_n)]^T, \boldsymbol{\omega} = [\omega(\mathbf{s}_1), \dots, \omega(\mathbf{s}_n)]^T$$



Generalized Local Statistical Model

$$\mathbf{Z} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\omega} + \mathbf{e} \sim N(\mathbf{X}\boldsymbol{\beta}, \boldsymbol{\Sigma} + \sigma_0^2 \mathbf{I})$$



Generalized local statistical model (GLS)

$$\text{diff}(\mathbf{Z}) = \mathbf{FZ} \sim N(\mathbf{FX}\boldsymbol{\beta}, \mathbf{F}\boldsymbol{\Sigma}\mathbf{F}^T + \sigma_0^2 \mathbf{F}\mathbf{F}^T)$$

↓ \approx (See theorem 3)

$$\text{diff}(\mathbf{Z}) = \mathbf{FZ} \sim N(\mathbf{FX}\boldsymbol{\beta}, \sigma^2 \mathbf{I} + \sigma_0^2 \mathbf{F}\mathbf{F}^T)$$

Convolution effects

When neighbor size ≥ 10 , $\sigma_0^2 \mathbf{F}\mathbf{F}^T \approx \sigma_0^2 \mathbf{I}$ (See Theorem 1)

Generalized least Squares

$$\mathbf{diff}(\mathbf{Z}) \sim \mathbf{N}(\mathbf{FX}\boldsymbol{\beta}, \sigma^2 \mathbf{I} + \sigma_0^2 \mathbf{FF}^T)$$



$$\arg \min_{\boldsymbol{\beta}, \sigma_0, \sigma} \left[(\mathbf{FZ} - \mathbf{FX}\boldsymbol{\beta})^T (\sigma^2 \mathbf{I} + \sigma_0^2 \mathbf{FF}^T)^{-1} (\mathbf{FZ} - \mathbf{FX}\boldsymbol{\beta}) \right],$$

subject to $\sigma_0^2 + \sigma^2 = 1$ and $\sigma_0, \sigma \geq 0$.



$$\arg \min_{\boldsymbol{\beta}, \sigma_0, \sigma} \left[\sum_{i=1}^n \frac{\{(\mathbf{FZ} - \mathbf{FX}\boldsymbol{\beta})^T \mathbf{q}_i\}^2}{\sigma^2 + \sigma_0^2 \lambda_i} \right], \text{ s. t. } \sigma_0^2 + \sigma^2 = 1; \sigma_0, \sigma \geq 0.$$

$$\frac{\partial^2 f_i}{\partial \theta^2} = \begin{bmatrix} \mathbf{X}^T \mathbf{F} \mathbf{q}_i (\sigma^2 + \sigma_0^2 \lambda_i) \\ (\mathbf{q}_i^T \mathbf{Z} - \mathbf{q}_i^T \mathbf{FX}\boldsymbol{\beta})^T \\ \lambda_i (\mathbf{q}_i^T \mathbf{Z} - \mathbf{q}_i^T \mathbf{FX}\boldsymbol{\beta})^T \end{bmatrix} \begin{bmatrix} \mathbf{X}^T \mathbf{F} \mathbf{q}_i (\sigma^2 + \sigma_0^2 \lambda_i) \\ (\mathbf{q}_i^T \mathbf{Z} - \mathbf{q}_i^T \mathbf{FX}\boldsymbol{\beta})^T \\ \lambda_i (\mathbf{q}_i^T \mathbf{Z} - \mathbf{q}_i^T \mathbf{FX}\boldsymbol{\beta})^T \end{bmatrix}^T \succcurlyeq 0.$$

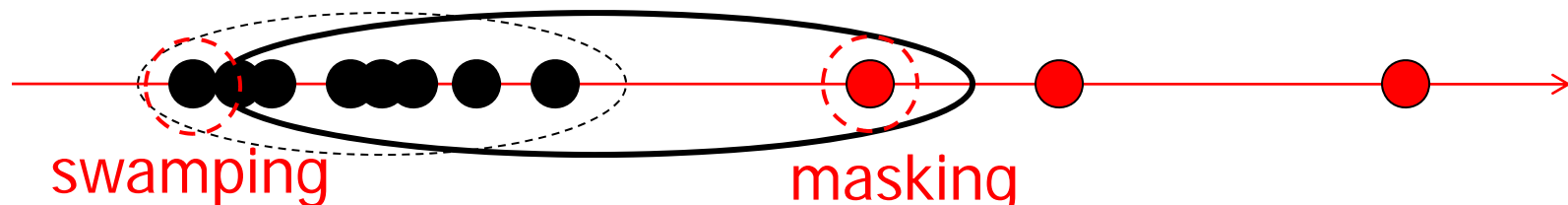
Forward / Backward Search

■ GLS Backward Search

- Model estimation by generalized least squares
- Remove the most probable outlier and update all local differences
- Repeat until the p values of all existing objects are greater than a threshold (e.g., 0.025)

■ GLS Forward Search

- Estimation by a robust subset \mathbf{S} of local differences
- Add test objects one by one to the training set \mathbf{S}
- Check the change of the smallest p value in \mathbf{S}
- A large drop in the smallest p value indicates an outlier





GLS Z-Test Statistics

Test if an observation $Z(\mathbf{s}_n)$ is abnormal:

$$\left\{ [\sigma^2 \mathbf{I} + \sigma_0^2 \mathbf{F}\mathbf{F}^T]^{-\frac{1}{2}} [\mathbf{F}\mathbf{Z} - \mathbf{F}\mathbf{X}\boldsymbol{\beta}] \right\}_n \geq \Phi^{-1} \left(\frac{\alpha}{2} \right), \quad e.x. \quad \alpha = 0.05$$



Connections with Existing Methods

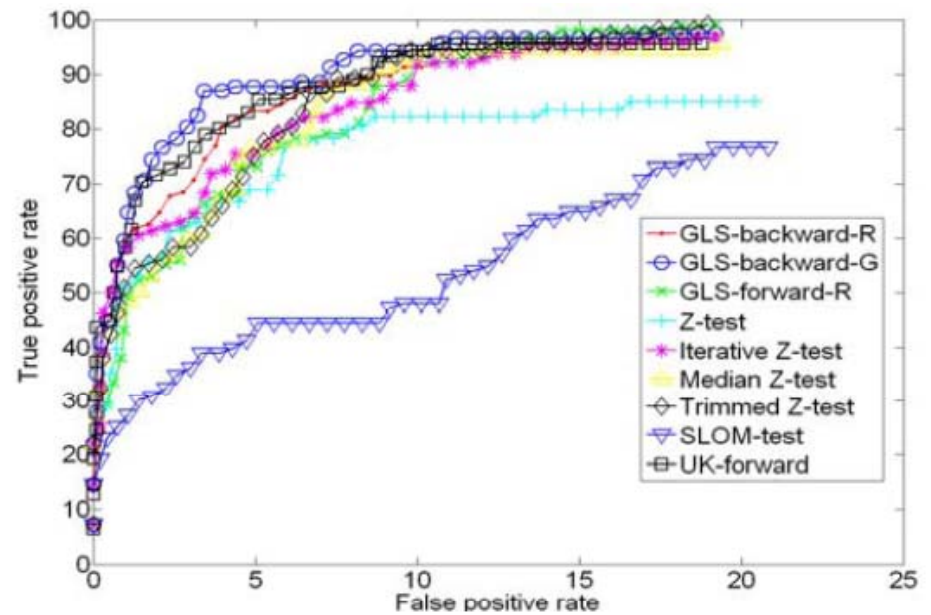
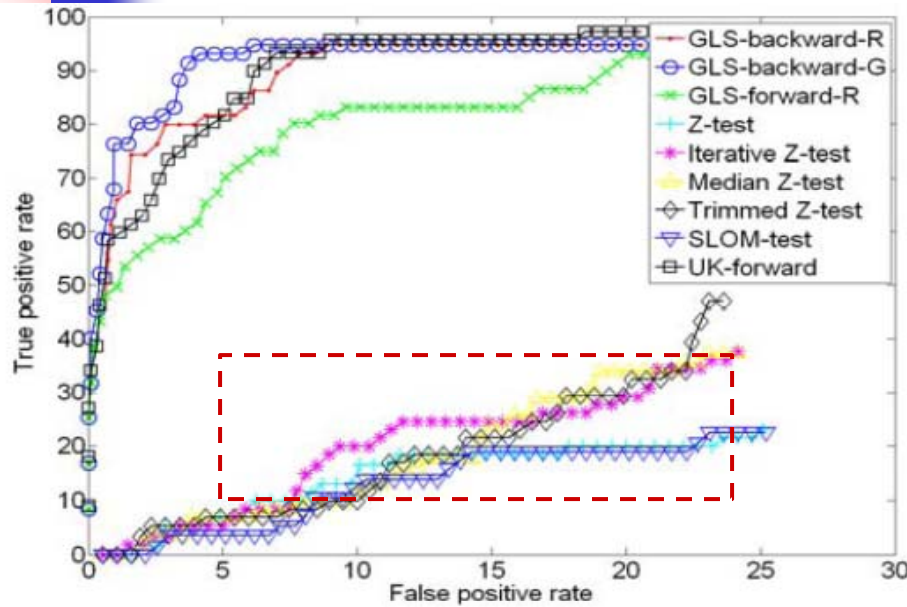
- If $F\Sigma F^T = \sigma^2 I$, then *GLS-SOD* is equivalent to *Universal Kriging SOD*.
 - Local vs. global estimator for a spatial Gaussian random field. The key is which estimator is more robust.
- When $F\Sigma F^T = \sigma^2 I$, $FXB = \mu I$ and $FF = \sigma_0^2 I$, then *Local based SOD* is equivalent to *GLS-SOD* and *Universal Kriging SOD*.
 - *Local based SOD* is a special case of *GLS-SOD*



Experiments

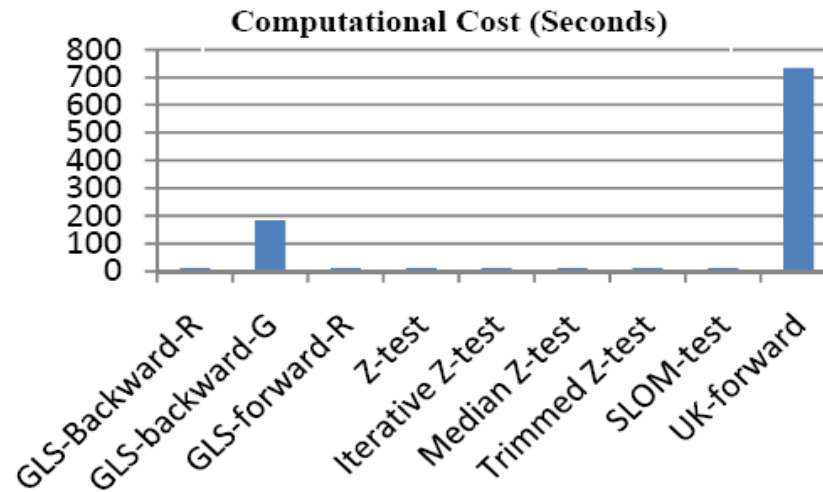
- Standard Simulation Model for SOD
 - 864 different simulation settings. Six repetitions for each setting; consider average error
 - Existing statistical SOD methods only consider 10 to 15 simulation settings in their experiments.

Simulation Results



f) Nonlinear trend, clustered outliers, $\alpha = 0.15, \sigma_0^2 = 10, c = 5, K$

a) Constant trend, isolated outliers, $\alpha = 0.1, \sigma_0^2 = 2, c = 15, K = 4$





Summary

- Design of a generalized local statistical framework
- Robust estimation and outlier detection methods based on the proposed GLS framework
- In-depth study on the connection between different SOD methods
- Comprehensive simulations to validate the effectiveness and efficiency of GLS



Thank you !

chenf@vt.edu



Property of "FΣF^T" (Theorem 2)

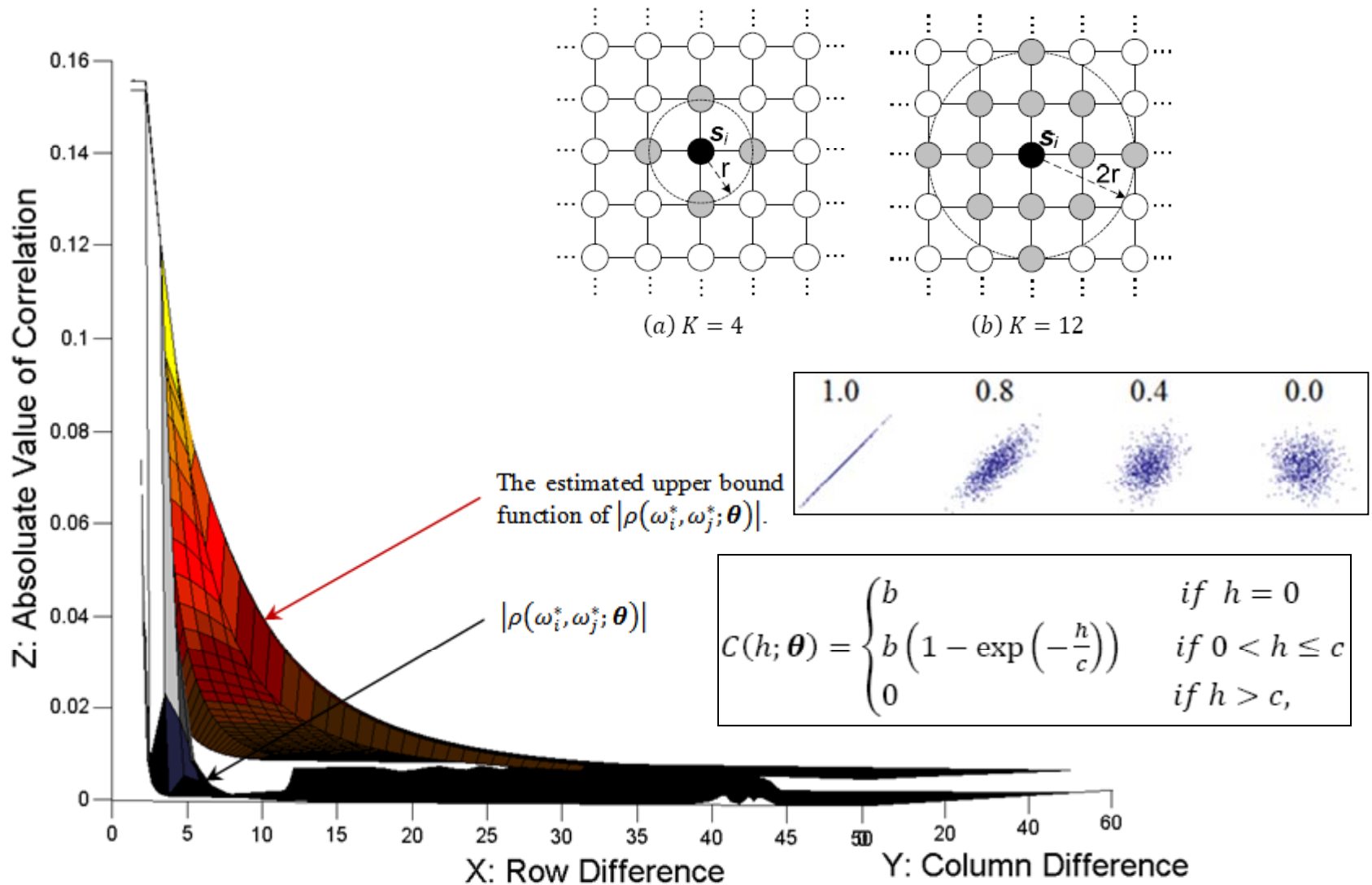


Figure 9: The comparison between the true correlation $|\rho(\omega_i^*, \omega_j^*; \theta)|$ and the estimated bound function. Here, $K = 12, c = 40$.



Property of “ $\sigma_0^2 FF^T$ ”

Theorem 1: *The random vector $\mathbf{e}^* = \mathbf{F}\mathbf{e}$ has two major properties*

- 1) *The variance $\text{Var}(e_i^*) = \frac{K+1}{K} \sigma_0^2, i = 1 \dots n,$*
- 2) *The correlation $|\rho(e_i^*, e_j^*)| \leq \frac{2}{K+1}, \forall i, j$ with $i \neq j,$*

where e_i^ refers to the i -th element in the vector \mathbf{e}^* .*

Conclusion: When selecting a relatively large neighborhood size to do local smoothing, we can approximate “ FF^T ” as an identity matrix.

$$\mathbf{diff}(\mathbf{Z}) \sim \mathbf{N}(\mathbf{FX}\boldsymbol{\beta}, \mathbf{F}\boldsymbol{\Sigma}\mathbf{F}^T + \sigma_0^2 \mathbf{F}\mathbf{F}^T).$$



Connections with Existing Methods

Theorem 5: *Suppose that $\mathbf{F}\boldsymbol{\Sigma}\mathbf{F}^T = \sigma^2\mathbf{I}$ and the parameters of Kriging-SOD and GLS-SOD are correctly calculated by robust estimation, then Kriging-SOD and GLS-SOD are equivalent.*

Theorem 6. *If $\mathbf{F}\boldsymbol{\Sigma}\mathbf{F}^T = \sigma^2\mathbf{I}$, $\sigma_0^2\mathbf{F}\mathbf{F}^T = \sigma_0^2\mathbf{I}$, the parameters of GLS-SOD and LS-SOD are correctly calculated by robust estimation, and one of the following conditions is true, then GLS-SOD becomes equivalent to LS-SOD.*

- (1) $\mathbf{Z}(\mathbf{s})$ has a constant trend (mean): $\mathbf{X}\boldsymbol{\beta} = c\mathbf{1}$, where c is a constant value.
- (2) $\mathbf{Z}(\mathbf{s})$ is a linear trend of spatial coordinates, and each point \mathbf{s} is the geometric center (or centroid) of its neighbors.

Simulations

- Simulation Model and Settings

- 864 different simulation settings. Six repetitions for each setting; consider average error

$$\mathbf{Z} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\omega} + \mathbf{e} \sim N(\mathbf{X}\boldsymbol{\beta}, \boldsymbol{\Sigma} + \sigma_0^2 \mathbf{I})$$

$$\boldsymbol{\omega} \sim N(\mathbf{0}, \boldsymbol{\Sigma})$$

$$C(h; \boldsymbol{\theta}) = \begin{cases} b & \text{if } h = 0 \\ b \left(1 - \exp\left(-\frac{h}{c}\right)\right) & \text{if } 0 < h \leq c \\ 0 & \text{if } h > c, \end{cases}$$

Exponential Model

$$\epsilon(\mathbf{s}) \sim \begin{cases} N(0, \sigma_0^2) & \text{with probability } 1 - \alpha \\ N(0, \sigma_C^2) & \text{with probability } \alpha \end{cases}$$

Contaminations

$$C(h; \boldsymbol{\theta}) = \begin{cases} b & \\ b \left(1 - \frac{3h}{2c} + \frac{1}{2} \left(\frac{h}{c}\right)^3\right) & \\ 0 & \end{cases}$$

Spherical Model

$$\mathbf{x}(\mathbf{s}) = [1, x(\mathbf{s}), y(\mathbf{s}), x(\mathbf{s}) \cdot y(\mathbf{s}), x(\mathbf{s})^2, y(\mathbf{s})^2]$$



Simulation Results

Table 3: Competition statistics for different combinations of parameter settings. Each cell contains three values, representing the win times for the related method based on the accuracies of top 10, 15, and 20 ranked outlier candidates for all methods.

Algorithm	Constant Trend	Linear Trend	Nonlinear Trend
<i>GLS-backward-R</i>	47, 47, 45	79, 72, 82	76, 81, 77
<i>GLS-backward-G</i>	88, 86, 89	114, 102, 120	141, 144, 138
<i>GLS-forward-R</i>	13, 11, 14	22, 25, 27	40, 36, 47
<i>Z-test</i>	47, 35, 40	29, 30, 13	0, 0, 0
<i>Iterative Z-test</i>	35, 46, 63	16, 20, 21	0, 0, 0
<i>Median Z-test</i>	20, 23, 29	1, 7, 8	0, 0, 0
<i>Trimmed Z-test</i>	15, 23, 32	5, 13, 13	0, 0, 0
<i>SLOM-test</i>	0, 0, 0	0, 0, 0	0, 0, 0

