

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

Feng Chen
Department of Computer Science
Virginia Tech, USA
chenf@vt.edu

Chang-Tien Lu
Department of Computer Science
Virginia Tech, USA
ctlu@vt.edu

Arnold P. Boedihardjo
Department of Computer Science
Virginia Tech, USA
Arnold.P.Boedihardjo@vt.edu

ABSTRACT

Local based approach is a major category of methods for spatial outlier detection (*SOD*). Currently, there is a lack of systematic analysis on the statistical properties of this framework. For example, most methods assume identical and independent normal distributions (i.i.d. normal) for the calculated local differences, but no justifications for this critical assumption have been presented. The methods' detection performance on geostatistic data with linear or nonlinear trend is also not well studied. In addition, there is a lack of theoretical connections and empirical comparisons between local and global based *SOD* approaches. This paper discusses all these fundamental issues under the proposed *Generalized Local Statistical (GLS)* framework. Furthermore, robust estimation and outlier detection methods are designed for the new *GLS* model. Extensive simulations demonstrated that the *SOD* method based on the *GLS* model significantly outperformed all existing approaches when the spatial data exhibits a linear or nonlinear trend.

Categories and Subject Descriptors

D.2.8 [Database Management]: Database Applications – data mining. I.5.3 [Pattern Recognition]: Outlier Detection.

General Terms

Algorithms, Theory, and Experimentation

Keywords

Spatial Outlier Detection, Spatial Gaussian Random Field.

1. INTRODUCTION

The ever-increasing volume of spatial data has greatly challenged our ability to extract useful but implicit knowledge from them. As an important branch of spatial data mining, spatial outlier detection aims to discover the objects whose non-spatial attribute values are significantly different from the values of their spatial neighbors [1]. In contrast to traditional outlier detection, spatial outlier detection must differentiate spatial and non-spatial attributes, and consider the spatial continuity and autocorrelation between nearby samples. By the first law of geography, "Everything is related to everything else, but nearby things are more related than distant things [3]."

There are two main classes of spatial outlier detection (*SOD*) methods: local and global based approaches. Local based approaches [4] first calculate the local difference (statistic) for

each object which is the difference between the non-spatial attribute of the object and the aggregated value (e.g., average) of its spatial neighbors. By assuming i.i.d. normal distributions for these local differences, the local based approaches discover outlier objects by robust estimation of model parameters such as the aggregated values, mean, and standard deviation. Various methods have been presented by using different spatial neighborhood definitions and robust estimation techniques [5-9]. The second class, global based methods, is to identify outliers using the robust estimator of a global kriging model which is the best linear unbiased estimator for geostatistical data. Particularly, Christensen et al. [10] proposed diagnostics to detect spatial outliers on the estimation of covariance function. Cerioli and Riani [11] developed a forward search procedure to identify spatial outliers for an ordinary kriging model. Militino et al. [12] further generalized the forward search method in [11] to a universal kriging model. This paper focuses on local based methods, because local based methods can achieve higher computational efficiency with minimal loss of accuracy. This feature of the local based approaches is demonstrated through extensive simulations described in Section 5.

This work is primarily motivated by the current situation where there is no systematic study on the statistical properties of local based *SOD* methods. For example, existing works assume i.i.d. on local differences, but justifications for the assumption have never been proposed. Also their performance on spatial data with linear or nonlinear trends has not been well studied. There is also a lack of research on the theoretical connections and empirical comparisons between local and global based *SOD* methods. To that end, this paper provides a generalized framework for local based *SOD* methods and theoretically and empirically compares it against global based *SOD* methods. The proposed framework is cast within the statistical abstraction of a spatial Gaussian random field which is the most popular model for geostatistical data [1,2]. A major reason for its popularity is that the optimal solution based on the Gaussian random field is equivalent to a best linear unbiased estimator (BLUE) for non-Gaussian data. It has been shown to provide accurate results in a variety of practical situations [1,2]. Sections 5.8, 6.3.3, and 7.4 in [2] give an in-depth discussion on the applicability of Gaussian random field.

A spatial Gaussian random field refers to a collection of dependent random variables that are associated with a set of spatial indexes, $\{Z(\mathbf{s}) \mid \mathbf{s} \in D \subset \mathbb{R}^2\}$, where D refers to a continuous fixed region. This family of random variables can be characterized by a joint Gaussian probability density or distribution. In real applications, only partial observations of one realization (or a partial sample of size one) are available. In order to make this model operational, the requirements for stationarity and isotropy, such as second-order or intrinsic stationarity, are further imposed. Imposing such assumptions helps reduce the number of

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

KDD'10, July 25–28, 2010, Washington, DC, USA.

Copyright 2010 ACM 978-1-4503-0055-1/10/07...\$10.00.

model parameters required to be estimated. When the data is second-order stationary and isotropic, the spatial correlation structure is described by a semivariogram or covariance function, in which the correlations between variables are dependent on their spatial distance. Statistical inferences are then performed by assuming explicit forms of the covariance and mean functions. Our major contributions are as follows:

- **Design of a generalized local statistical (GLS) framework:** The general local statistical model is a generalized statistical framework for existing local based *SOD* methods. It can effectively handle complex situations where the spatial data exhibit a global trend or non-negligible dependences between local differences.
- **Robust estimation and outlier detection methods based on the proposed GLS framework:** We analyze the contamination issues that lead to masking and swamping effects. Based on the analysis, two robust algorithms, *GLS*-backward search and *GLS*-forward search, are proposed to estimate the parameters for the *GLS* model.
- **In-depth study on the connection between different SOD methods:** We present theoretical foundations for existing local based *SOD* methods and discuss the crucial connections between local and global based *SOD* methods.
- **Comprehensive simulations to validate the effectiveness and efficiency of GLS.** This is the first work that provides extensive comparisons between existing popular methods through a systematic simulation study. The results showed the proposed *GLS-SOD* approach significantly outperformed all existing methods when the spatial data exhibits a linear or nonlinear trend.

This paper is organized as follows. Section 2 reviews spatial local statistics and related works. Section 3 introduces the generalized local statistical model and presents a rigorous theoretical treatment of its fundamental statistical properties. In Section 4, we introduce several robust estimation and outlier detection methods for the *GLS* model, and analyze the connection between different *SOD* methods. Section 5 provides the simulations and discussions. Section 6 gives the conclusion.

2. SPATIAL LOCAL STATISTICS AND RELATED WORKS

Given a set of observations $\{Z(\mathbf{s}_1), Z(\mathbf{s}_2), \dots, Z(\mathbf{s}_n)\}$, a local spatial statistic [4] is defined as

$$S(\mathbf{s}) = [Z(\mathbf{s}) - E_{\mathbf{s}_i \in N(\mathbf{s})}(Z(\mathbf{s}_i))], \quad (1)$$

where $\mathbf{G} = \{\mathbf{s}_1, \dots, \mathbf{s}_n\} \subset \mathbb{R}^2$ is a set of spatial locations, $\mathbf{s} \in \mathbf{G}$, $Z(\mathbf{s}) \in \mathbb{R}$ represents the value of Z attribute at location \mathbf{s} , $N(\mathbf{s})$ is the set of spatial neighbors of \mathbf{s} , and $E_{\mathbf{s}_i \in N(\mathbf{s})}(Z(\mathbf{s}_i))$ represents the average Z attribute value of the neighbors of \mathbf{s} . It is assumed that the set of local spatial statistics $\{S(\mathbf{s}_1), \dots, S(\mathbf{s}_n)\}$ are independently and identically normally distributed (i.i.d. normal). Then the popular Z -test [4] for detecting spatial outliers can be described as follows: Spatial statistic $Z_{S(\mathbf{s})} = \left| \frac{S(\mathbf{s}) - \mu_s}{\sigma_s} \right| > \Phi^{-1}\left(\frac{\alpha}{2}\right)$, where Φ is the cumulative distribution function (CDF) of a standard normal distribution, α refers to a significance level and is usually set to 0.05, and μ_s and σ_s refer to the sample mean and sample standard deviation, respectively.

Lu et al. [5] pointed out that the Z -test is susceptible to the well-known masking and swamping effects. When multiple

outliers exist in the data, the quantities $E_{\mathbf{s}_i \in N(\mathbf{s})}(Z(\mathbf{s}_i))$, μ_s , and σ_s are biased estimators of the population means and standard deviation. As a result, some true outliers are "masked" as normal objects and some normal objects are "swamped" and misclassified as outliers. The authors proposed an iterative approach that detects outliers by multi-iterations. Each iteration identifies only one outlier and then modifies its attribute value so that it will not impact the results of subsequent iterations. Later, Chen et al. [6] proposed a median based approach that uses median estimator for the quantities $E_{\mathbf{s}_i \in N(\mathbf{s})}(Z(\mathbf{s}_i))$ and μ_s , and median absolute deviation (MAD) estimator for σ_s . Hu and Sung [7] proposed an approach similar to [6], but using trimmed mean to estimate $E_{\mathbf{s}_i \in N(\mathbf{s})}(Z(\mathbf{s}_i))$, instead of the median. Sun and Chawla [8] presented a spatial local outlier measure to capture the local behavior of data in their neighborhood. Shekhar et al. [9] employed a graph-based method to define spatial neighborhoods ($N(\mathbf{s})$) and their method is applied to a special case of transportation network.

Most existing local based methods assume that the set of local statistics $\{S(\mathbf{s}_1), \dots, S(\mathbf{s}_n)\}$ are i.i.d. normal, but no justifications for this assumption have ever been proposed. As we will discuss in next sections, this i.i.d. assumption is only approximately true in certain scenarios, and the dependencies between different local differences (statistics) must be considered when the spatial data exhibit linear or nonlinear trend or the selected neighborhood size for each object is small. As shown in our simulations in Section 5, the violation of i.i.d. assumption can significantly impact the accuracies of the outlier detection methods.

3. GENERALIZED LOCAL SPATIAL STATISTICS

This section first introduces some preliminary background on spatial Gaussian random field, then presents the generalized local statistical (*GLS*) model, and finally discusses the statistical properties of the *GLS* model. Table 1 summarizes the key notations used in this paper.

Table 1: Description of Major Symbols

Symbol	Descriptions
$\{Z(\mathbf{s}_i)\}_{i=1}^n$	A given set of observations, where $\mathbf{s}_i \in \mathbb{R}^2$ is the spatial location and $Z(\cdot)$ is the Z attribute value.
$\{\mathbf{x}(\mathbf{s}_i)\}_{i=1}^n$	$\mathbf{x}(\mathbf{s}_i)$ is a vector of covariates of \mathbf{s}_i , such as the bases of spatial coordinates of \mathbf{s}_i .
\mathbf{Z}	$\mathbf{Z} = [Z(\mathbf{s}_1), \dots, Z(\mathbf{s}_n)]^T$
\mathbf{X}	$\mathbf{X} = [\mathbf{x}(\mathbf{s}_1), \dots, \mathbf{x}(\mathbf{s}_n)]^T$
\mathbf{F}	Neighborhood weight matrix; See Equation 4
$N(\mathbf{s})$	A general definition of spatial neighbors of \mathbf{s} .
$N_K(\mathbf{s})$	K -nearest neighbors of \mathbf{s} . This paper considers $N_K(\mathbf{s})$ as the specification of $N(\mathbf{s})$.
K	Neighborhood size. It is the major parameter to define spatial neighbors ($N_K(\mathbf{s})$).
<i>SOD</i>	S patial O utlier D etection
<i>GLS</i>	G eneralized L ocal S tatistics M odel
β, σ, σ_0	The unknown parameters in the <i>GLS</i> model

3.1 Generalized Local Statistic Model (GLS)

Given a spatial Gaussian random field $\{Z(\mathbf{s}), \mathbf{s} \in D \subset \mathbb{R}^2\}$, consider the following decomposition of the process [1]

$$Z(\mathbf{s}) = f(\mathbf{x}(\mathbf{s}), \boldsymbol{\beta}) + \omega(\mathbf{s}) + \epsilon(\mathbf{s}), \quad (2)$$

where D is a fixed region, $f(\mathbf{x}(\mathbf{s}), \boldsymbol{\beta})$ is the large scale trend (mean) of the process, $\omega(\mathbf{s})$ is the smooth-scale variation that is a Gaussian stationary process, and $\epsilon(\mathbf{s})$ is the white noise with the variance σ_0^2 .

The large scale trend $f(\mathbf{x}(\mathbf{s}), \boldsymbol{\beta}) = \mathbf{x}(\mathbf{s})^T \boldsymbol{\beta}$, where $\boldsymbol{\beta}$ is a vector of trend parameters, and $\mathbf{x}(\mathbf{s})$ is a vector of covariates that are the basis functions of spatial coordinates of \mathbf{s} (See Section 5.1 for illustrations). The nonlinear degree of the trend is dependent on the polynomials of covariates in $\mathbf{x}(\mathbf{s})$. For the smooth-scale variation $\omega(\mathbf{s})$, we assume that it is an isotropic second-order stationary process, in which the covariance $\text{Cov}(Z(\mathbf{s}_i), Z(\mathbf{s}_j))$ is a function of the spatial distance between \mathbf{s}_i and \mathbf{s}_j : $C(\|\mathbf{s}_i - \mathbf{s}_j\|; \boldsymbol{\theta})$, where $\boldsymbol{\theta}$ are function parameters. A variety of distance metrics may be selected, such as L_2 (Euclidean distance), L_1 (Manhattan distance), and graph distance [10]. There are two popular models for the covariance function C , including spherical model and exponential model (see Equations 8 and 11).

Given a sample set $\{Z(\mathbf{s}_1), \dots, Z(\mathbf{s}_n)\}$ that are partial observations of a particular realization of the spatial Gaussian random field, let $\mathbf{Z} = [Z(\mathbf{s}_1), \dots, Z(\mathbf{s}_n)]^T$, $\mathbf{e} = [e(\mathbf{s}_1), \dots, e(\mathbf{s}_n)]^T$, $\boldsymbol{\omega} = [\omega(\mathbf{s}_1), \dots, \omega(\mathbf{s}_n)]^T$, and $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^T$. By Equation 2, the random vector \mathbf{Z} has the decomposition form as

$$\mathbf{Z} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\omega} + \mathbf{e} \sim N(\mathbf{X}\boldsymbol{\beta}, \boldsymbol{\Sigma} + \sigma_0^2 \mathbf{I}), \quad (3)$$

where $\boldsymbol{\omega} \sim N(\mathbf{0}_{n \times 1}, \boldsymbol{\Sigma}_{n \times n})$ and $\mathbf{e} \sim N(\mathbf{0}_{n \times 1}, \sigma_0^2 \mathbf{I}_{n \times n})$.

The vector of local spatial statistics calculated by Equation 1 can be reformulated as the matrix form

$$\text{diff}(\mathbf{Z}) = \mathbf{F}\mathbf{Z}, \quad (4)$$

where $\mathbf{F} \in \mathbb{R}^{n \times n}$ is a neighborhood weight matrix with $F_{ij} = 1$ if $i = j$; $F_{ij} = -1/K$ if $\mathbf{s}_j \in N_K(\mathbf{s}_i)$; and $F_{ij} = 0$, otherwise.

By Equations 3 and 4, we can readily derive the generalized local statistical (GLS) model as

$$\text{diff}(\mathbf{Z}) \sim N(\mathbf{F}\mathbf{X}\boldsymbol{\beta}, \mathbf{F}\boldsymbol{\Sigma}\mathbf{F}^T + \sigma_0^2 \mathbf{F}\mathbf{F}^T). \quad (5)$$

Recall that $\boldsymbol{\Sigma}_{ij} = C(\|\mathbf{s}_i - \mathbf{s}_j\|; \boldsymbol{\theta})$. The GLS model has the unknown components $\boldsymbol{\beta}$, σ_0 , and $\boldsymbol{\theta}$, including $(|\boldsymbol{\beta}| + 1 + |\boldsymbol{\theta}|)$ parameters. Because the covariance function $C(\cdot)$ (e.g., spherical or exponential model) is nonlinear and nonconvex, it requires iterative reweighted generalized least squares algorithm to estimate all these parameters which is computationally expensive and can only guarantee a local optimum [14].

As shown in Section 3.2, an important property of the GLS model is that the component $\mathbf{F}\boldsymbol{\Sigma}\mathbf{F}^T$ can be approximated by $\sigma_0^2 \mathbf{I}$. Hence the GLS form (5) becomes asymptotically equivalent to

$$\text{diff}(\mathbf{Z}) \sim N(\mathbf{F}\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I} + \sigma_0^2 \mathbf{F}\mathbf{F}^T). \quad (6)$$

As discussed in Section 4.1, the model fitting for the GLS form (6) is a convex problem and can be efficiently solved.

By Section 3.2 Theorem 1, when the neighborhood size is relatively large with $K \geq 10$, the component $\sigma_0^2 \mathbf{F}\mathbf{F}^T$ can be further approximated by $\sigma_0^2 \mathbf{I}$. This leads to a simpler form of GLS

$$\text{diff}(\mathbf{Z}) \sim N(\mathbf{F}\mathbf{X}\boldsymbol{\beta}, (\sigma^2 + \sigma_0^2) \mathbf{I}). \quad (7)$$

Discussion: Local statistics is a popular technique used to reduce the dependence between sample points. However, by employing the decomposition form as indicated in Equations 2-4, we observe that local statistics help reduce the correlations between sample points caused by smooth-scale random variations, but at the same

time it also induces "new" correlations due to the averaging of white noise variations. As discussed in [2], correlated data can be expressed as linear combination of uncorrelated data. The approximate GLS form (6) explicitly models the "new" correlations caused by the averaging of white noises variations. The approximate GLS form (7) essentially ignores these "new" correlations. The form (7) may be considered when users expect high efficiency and allow some loss of accuracy. This tradeoff is studied in Section 5 by simulations.

3.2 Theoretical Properties of GLS

This sub-section studies the properties of two major covariance components $\sigma_0^2 \mathbf{F}\mathbf{F}^T$ and $\mathbf{F}\boldsymbol{\Sigma}\mathbf{F}^T$, and discusses the situations where they can be approximated by $\sigma_0^2 \mathbf{I}$ and $\sigma^2 \mathbf{I}$, respectively. As shown in Equation 3, $\sigma_0^2 \mathbf{F}\mathbf{F}^T$ and $\mathbf{F}\boldsymbol{\Sigma}\mathbf{F}^T$ are the covariance matrices of the random vectors $\mathbf{e}^* = \mathbf{F}\mathbf{e}$ and $\boldsymbol{\omega}^* = \mathbf{F}\boldsymbol{\omega}$, respectively. We focus on the study of their correlation structures. Because \mathbf{e}^* and $\boldsymbol{\omega}^*$ are both multivariate normally distributed, the correlation structure gives important information about the related dependence structure (e.g., in-correlation implies independence). Three related theorems are stated as follows:

Theorem 1: *The random vector \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}^* .

(Readers are referred to [15] for the proof.)

Theorem 1 indicates that when the neighborhood size is relative large, the correlations between the components in \mathbf{e}^* are very low (e.g., smaller than 0.2 when $K = 10$) and the variance of each component is very close to σ_0^2 . In this case, $\sigma_0^2 \mathbf{F}\mathbf{F}^T \approx \sigma_0^2 \mathbf{I}$. However, for a small neighborhood size, as shown in simulations (Section 5), the dependence between the components in \mathbf{e}^* must be considered.

The next two theorems are related to the random vector $\boldsymbol{\omega}^*$. It is difficult to analytically evaluate $\boldsymbol{\omega}^*$, because it is generated by an isotropic second order stationary process, and even when the explicit form of the covariance function is known, the statistical properties of $\boldsymbol{\omega}^*$ are still not straightforward. For this reason, several additional assumptions need to be considered. The following are three assumptions required for Theorem 2:

A1. *If $N_K(\mathbf{s}_i) \cap N_K(\mathbf{s}_d) \neq \Phi$, then, $\forall \mathbf{s}_i, \mathbf{s}_j, \mathbf{s}_t \in N_K(\mathbf{s}_i) \cup$*

$N_K(\mathbf{s}_d)$, their between spatial distances are approximately equivalent: $\|\mathbf{s}_j - \mathbf{s}_i\| \approx \|\mathbf{s}_t - \mathbf{s}_i\| \approx \|\mathbf{s}_j - \mathbf{s}_t\|$.

A2. *If $\mathbf{s}_j \in N_K(\mathbf{s}_i), \mathbf{s}_t \notin N_K(\mathbf{s}_i)$, and $N_K(\mathbf{s}_t) \cap N_K(\mathbf{s}_i) = \Phi$, then $\|\mathbf{s}_t - \mathbf{s}_i\| \approx \|\mathbf{s}_t - \mathbf{s}_j\|$.*

A3. *The distance between any points that are k -nearest neighbors is approximately constant everywhere.*

The intuition on assumptions A1 and A2 is that, because neighbors are close to each other, they share similar between-distances and similar distances to points that are not their neighbors. The assumption A3 is valid when the spatial locations follow a uniform distribution or a grid structure. The assumption A3 holds in many practical situations [13]. The situations where assumptions A1 and A2 are potentially violated will be discussed in Theorem 3.

Theorem 2: *If the above assumptions A1 and A2 hold, then the random vector $\boldsymbol{\omega}^*$ has two major properties*

- 1) The variance $\text{Var}(\omega_i^*) \approx \frac{1+K}{K}(\sigma^2 - C_{x_i}), i = 1 \dots n$
- 2) The correlation $\rho(\omega_i^*, \omega_j^*) \approx -\frac{1}{K}$, if $s_j \in N_K(s_i)$ or $s_i \in N_K(s_j)$; otherwise, $\rho(\omega_i^*, \omega_j^*) \approx 0$,

where C_{s_i} refers to the average covariance value between s_i and its K -nearest neighbors, and $\sigma^2 = C(0)$ refers to the constant variance for each component of ω . Further, if the assumption A3 also holds, then the variance $\text{Var}(\omega_i^*)$ becomes constant everywhere.

(Readers are referred to [15] for the proof.)

Theorem 2 indicates that the correlations between the components in ω^* are mostly zero, except for neighboring points. Particularly, the correlations between neighboring points are all negative, and their major impact factor is the neighborhood size K . The greater the value of K , the less the neighbor points are correlated. However, K cannot be arbitrary large; otherwise, the assumptions made above will be violated. For example, suppose $n = 200$ and $K = 10$, then only about 5% of pairs are correlated. For these correlated components, the correlations are only close to -0.1 . As shown in Figure 1, 0.1 indicates a negligible correlation.

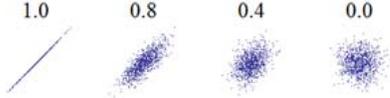


Figure 1: An example of correlation: it reflects the noise and direction of a linear relationship [13].

Theorem 2 states two approximate properties of ω^* . However, it is not directly known how these properties are impacted if assumptions A1 and A2 are violated. The next Theorem 3 will delve deeper into this issue and provide more specific analysis on ω_i^* . For Theorem 3, the following less restrictive assumptions are employed:

B1. The spatial locations $\{s_1, \dots, s_n\}$ follow a grid structure and $n \leq 2500$;

B2. The spatial distance is defined by L_2 (Euclidean) distance;

B3. The covariance function $\text{Cov}(Z(s_i), Z(s_j)) = C(h)$, where $h = \|s_i - s_j\|_2$, follows a popular spherical model;

B4. Consider 4 or 12-nearest neighbors as spatial neighbors for each object.

Assumptions B1 and B2 are generic properties that can be readily applied to spatial data in general [1, 2]. In many applications, the total number of spatial locations is smaller than 300 [1]. Here, we consider a much enlarged range with $n \leq 2500$, for the purpose of generality. For assumption B3, a spherical model is defined as

$$C(h; \theta) = \begin{cases} b & \text{if } h = 0 \\ b \left(1 - \frac{3h}{2c} + \frac{1}{2} \left(\frac{h}{c}\right)^3\right) & \text{if } 0 < h \leq c \\ 0 & \text{if } h > c, \end{cases} \quad (8)$$

where $\theta = (b, c)^T, b \geq 0, c \geq 0$. Note that $b = C(0; \theta)$ refers to the constant variance for each object s , and $C(h; \theta)$ is a decreasing function on the distance h .

The reason for using a spherical model as opposed to exponential or Gaussian models is that the spherical model leads to closed-form analytical results. The closed-form results will provide important insights into its statistical properties. As for assumption B4, K is set to 4 or 12 due to the use of the grid

structure (assumption B1). In a grid data, each object has four nearest objects with the same distance r , eight next-nearest objects with the same distance $2r$, and so on, where r is the grid cell size. Hence, we can select $K = 4, 12, 24, \dots$. We select the first two values with $K = 4$ and 12, which are equivalent to defining neighborhoods with radiuses of r and $2r$, respectively.

Theorem 3: Under the above four assumptions, the random vector ω^* has following properties on the correlation structure

- 1) If $K = 4$, then
 - a) $\rho(\omega_i^*, \omega_j^*) = 0$, if $d(s_j, s_i) > c + 2r$,
 - b) $|\rho(\omega_i^*, \omega_j^*)| \leq 0.4$, if $c \leq 2r$ and $d(s_j, s_i) \leq 2r$,
 - c) $|\rho(\omega_i^*, \omega_j^*)| \leq 0.22$, if $c > 2r$ and $d(s_j, s_i) \leq 2r$,
 - d) $|\rho(\omega_i^*, \omega_j^*)| \leq 0.05$, if $d(s_j, s_i) > 2r$.
- 2) If $K = 12, d(s_j, s_i) > c + 4r$, then $\rho(\omega_i^*, \omega_j^*) = 0$
- 3) If $K = 12, c < 4r$, then
 - a) $|\rho(\omega_i^*, \omega_j^*)| \leq 0.220$, if $d(s_j, s_i) \leq 2r$
 - b) $|\rho(\omega_i^*, \omega_j^*)| \leq 0.110$, if $2r < d(s_j, s_i) \leq 3r$
 - c) $|\rho(\omega_i^*, \omega_j^*)| \leq 0.050$, if $d(s_j, s_i) > 3r$
- 4) If $K = 12, c \geq 4r$ and $\text{row}(s_j) = \text{row}(s_i)$ (or $\text{col}(s_j) = \text{col}(s_i)$), then
 - a) $|\rho(\omega_i^*, \omega_j^*)| \leq 0.4741 - \frac{0.1179 \cdot c^2/r^2}{1+c^2/(2.707 \cdot r^2)}$, if $d(s_j, s_i) = r$
 - b) $|\rho(\omega_i^*, \omega_j^*)| \leq 0.1203$, if $d(s_j, s_i) = 2r$
 - c) $|\rho(\omega_i^*, \omega_j^*)| \leq 0.1719 - \frac{0.0158 \cdot h_{ij}^2/r^2}{1+c^2/(10.5174 \cdot r^2)}$, otherwise.
- 5) If $K = 12, c \geq 4r, \text{row}(s_j) \neq \text{row}(s_i), \text{col}(s_j) \neq \text{col}(s_i)$, then
$$|\rho(\omega_i^*, \omega_j^*)| \leq 0.1085 - \frac{0.0028 \cdot h_{ij}^2/r^2}{1+h_{ij}^2/(37.6723 \cdot r^2)}$$

where r refers to the grid cell size; $\text{row}(s_i)$ and $\text{col}(s_i)$ refer to the row and column locations of the object s_i in the grid structure; $h_{ij} = d(s_j, s_i)$ is the Euclidean distance between s_i and s_j .

(Readers are referred to [15] for the proof.)

Theorem 3 implies similar patterns as drawn by Theorem 2 although Theorem 2 provides only approximate properties. Theorem 3 is a further justification of these patterns. In the following discussions, we consider the situation with $c \geq 5$. The situation with $c < 5$ will be discussed separately. By Theorem 3, if $c \geq 5$, then $|\rho(\omega_i^*, \omega_j^*)| \leq 0.22$ when $K = 4$; and $|\rho(\omega_i^*, \omega_j^*)| \leq 0.18$ when $K = 12$. It indicates small absolute correlation values for different K values. The correlation values slightly decreases when K increases. It can also be shown that most correlations are negative and are close or equal to zero. Readers are referred to [15] for detailed information about $\rho(\omega_i^*, \omega_j^*)$. All these observations are consistent with the Theorem 2.

We have a comparison between $\sigma_0^2 \mathbf{F}\mathbf{F}^T$ and $\mathbf{F}\mathbf{Z}\mathbf{F}^T$. Consider two typical situations: $K = 4$ to represent a small neighborhood; and $K = 12$ to represent a relatively large neighborhood. If $K =$

4, then $|\rho(e_i^*, e_j^*)| \leq 0.4$ and $|\rho(\omega_i^*, \omega_j^*)| \leq 0.22$. If $K = 12$, then $|\rho(e_i^*, e_j^*)| \leq 0.2$ and $|\rho(\omega_i^*, \omega_j^*)| \leq 0.18$. The impacts of these correlation values (degrees) are shown in Figure 1. Although both $|\rho(e_i^*, e_j^*)|$ and $|\rho(\omega_i^*, \omega_j^*)|$ increase when the neighborhood size K decreases, the absolute correlation $|\rho(e_i^*, e_j^*)|$ increases drastically. Based on these results, we will approximate $\mathbf{F}\mathbf{\Sigma}\mathbf{F}^T$ by $\sigma^2\mathbf{I}$ for different settings of K , but will only approximate $\sigma_0^2\mathbf{F}\mathbf{F}^T$ by $\sigma_0^2\mathbf{I}$, when K is relatively large, e.g., $K \geq 10$.

Theorem 3 also indicates that when c is small (e.g., $c < 5r$), some correlations are relatively high (e.g., $|\rho(\omega_i^*, \omega_j^*)| = 0.4$ if $K = 4, c = 1r$, and $d(\mathbf{s}_j, \mathbf{s}_i) = r$). In this case, an important observation is that the correlation matrix of ω^* exhibits similar structure as that of e^* . Particularly, if $c < r$, these two correlation matrices become identical. In this situation, it is still reasonable to approximate the correlation matrix of ω^* as identity or unit matrix. The lost structure information by this approximation will be recovered while estimating the parameter σ_0 for the vector e^* , because of the similar structure between the covariance matrices $\text{Var}(\omega^*)$ and $\text{Var}(e^*)$. For example, suppose $c < r$ and the constant variance for each component of e is σ_e^2 , then we have $\text{Var}(e) = \mathbf{\Sigma} = \sigma_e^2\mathbf{I}$, and $\text{Var}(e^*) = \text{Var}(Fe) = \mathbf{F}\mathbf{\Sigma}\mathbf{F}^T = \sigma_e^2\mathbf{F}\mathbf{F}^T$. By Equation 5, the true distribution model is: $\text{diff}(\mathbf{Z}) \sim \mathbf{N}(\mathbf{F}\mathbf{X}\boldsymbol{\beta}, \mathbf{F}\mathbf{\Sigma}\mathbf{F}^T + \sigma_0^2\mathbf{F}\mathbf{F}^T) = \mathbf{N}(\mathbf{F}\mathbf{X}\boldsymbol{\beta}, (\sigma_0^2 + \sigma_e^2)\mathbf{F}\mathbf{F}^T)$. If we approximate $\mathbf{F}\mathbf{\Sigma}\mathbf{F}^T$ as $\sigma^2\mathbf{I}$ instead, then by Equation 6 the approximate model becomes $\text{diff}(\mathbf{Z}) \sim \mathbf{N}(\mathbf{F}\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I} + \sigma_0^2\mathbf{F}\mathbf{F}^T)$. Using robust parameter estimation, the approximate model can completely recover the true distribution, e.g., by setting the estimated parameters $\hat{\sigma} = 0$ and $\hat{\sigma}_0 = \sqrt{\sigma_0^2 + \sigma_e^2}$.

4. ESTIMATION AND INFERENCES

Spatial outlier detection (SOD) is usually coupled with a robust estimation process for the related statistical model. This section presents robust estimation and outlier detection methods to reduce the masking and swamping effects, and then discusses the connection between the proposed GLS-SOD methods and existing representative methods, such as kriging-based and Z-test SOD methods. We are focused on the estimation techniques for the GLS form (6) that is regarded as the default model. The GLS form (7) will be explicitly stated when discussing its techniques.

4.1 Generalized Least Squares Regression

Given a set of observations $\{Z(\mathbf{s}_1), \dots, Z(\mathbf{s}_n)\}$, the objective is to estimate the parameters $\boldsymbol{\beta}, \sigma$, and σ_0 for the proposed GLS model. We consider mean squared error (MSE) as the score function which is the most popular error function in spatial statistics [11]. This leads to a generalized least square problem for the GLS form (6) and can be formulated as:

$$\arg \min_{\boldsymbol{\beta}, \sigma_0, \sigma} [(\mathbf{F}\mathbf{Z} - \mathbf{F}\mathbf{X}\boldsymbol{\beta})^T (\sigma^2\mathbf{I} + \sigma_0^2\mathbf{F}\mathbf{F}^T)^{-1} (\mathbf{F}\mathbf{Z} - \mathbf{F}\mathbf{X}\boldsymbol{\beta})],$$

subject to $\sigma_0^2 + \sigma^2 = 1$ and $\sigma_0, \sigma \geq 0$. (9)

Note that we scale σ_0 and σ by a factor c with $\sigma_0^* = \sigma_0/c$ and $\sigma^* = \sigma/c$, such that $\sigma_0^{*2} + \sigma^{*2} = 1$. Without this constraint, the objective function in (9) will always be minimized by setting $\sigma_0 = \sigma = \infty$, and $\boldsymbol{\beta}$ to any value. For simplicity, we directly use the original symbols σ_0 and σ , rather than σ_0^* and σ^* . As shown in Theorem 4, the problem (9) is a convex optimization problem which can be efficiently solved by numerical optimization methods such as interior point method [14]. Note

that when the neighborhood size K is large, the following approximation holds: $\sigma_0^2\mathbf{F}\mathbf{F}^T \approx \sigma_0^2\mathbf{I}$ (see Section 3.2 Equation 7). Then the problem (9) reduces to a regular least squares regression problem and an explicit solution is available with $\boldsymbol{\beta} = (\mathbf{X}^T\mathbf{F}^T\mathbf{F}\mathbf{X})^{-1}\mathbf{X}^T\mathbf{F}^T\mathbf{F}\mathbf{Z}$, and $(\sigma^2 + \sigma_0^2) = \|\mathbf{F}\mathbf{X}\boldsymbol{\beta} - \mathbf{F}\mathbf{Z}\|_2^2 / (n - p - 1)$, where p is the size of the vector $\boldsymbol{\beta}$. For the purpose of outlier detection, it is unnecessary to further derive the explicit forms of σ and σ_0 .

Theorem 4: *The problem (9) is a convex optimization problem.*

Proof Sketch: Suppose λ_i and \mathbf{q}_i are the eigenvalues and corresponding (orthonormal) eigenvectors of the matrix $\mathbf{F}\mathbf{F}^T$. It can be readily shown that the problem (9) is equivalent to

$$\arg \min_{\boldsymbol{\beta}, \sigma_0, \sigma} \left[\sum_{i=1}^n \frac{(\mathbf{F}\mathbf{Z} - \mathbf{F}\mathbf{X}\boldsymbol{\beta})^T \mathbf{q}_i)^2}{\sigma^2 + \sigma_0^2 \lambda_i} \right], \text{ s.t. } \sigma_0, \sigma \geq 0 \quad (10)$$

Let $f_i = \frac{(\mathbf{F}\mathbf{Z} - \mathbf{F}\mathbf{X}\boldsymbol{\beta})^T \mathbf{q}_i)^2}{\sigma^2 + \sigma_0^2 \lambda_i}$, It suffices to prove that f_i is a convex function, or equivalently $\frac{\partial^2 f_i}{\partial \boldsymbol{\theta}^2} \succcurlyeq 0$, $\boldsymbol{\theta} = [\boldsymbol{\beta}^T, \sigma^2, \sigma_0^2]^T$.

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

When the parameters $\boldsymbol{\beta}, \sigma$, and σ_0 are estimated by generalized least squares, we can calculate standard residuals and use standard statistic test procedure to identify outliers. This method works well for sample data with small data contamination, but is susceptible to the well-known masking and swamping effects when multiple outliers exist. For the GLS model, the masking and swamping effects originate from two phases of the estimation process:

1) **Phase I contamination** occurs in the process of calculating local differences $\mathbf{F}\mathbf{Z}$. For example, suppose we define neighbors by the K-nearest-neighbor rule. Consider an outlier object $Z^*(\mathbf{s}_1) = Z(\mathbf{s}_1) + \zeta_1$, where $Z(\mathbf{s}_1)$ is the normal value but it is contaminated by a large error ζ_1 , and suppose only one of its neighbors is an outlier with $Z^*(\mathbf{s}) = Z(\mathbf{s}) + \zeta$, where ζ is the contamination error. The local difference $\text{diff}(Z^*(\mathbf{s}_1)) = \left[Z(\mathbf{s}_1) - \frac{1}{K} \sum_{\mathbf{s}_i \in N(\mathbf{s}_1)} Z(\mathbf{s}_i) \right] + \zeta_1 - \frac{\zeta}{K}$. If $\zeta = K \cdot \zeta_1$, then the error is marginalized and we obtain a normal local difference for an outlier object $Z^*(\mathbf{s}_1)$ which will be identified as a normal object. If $Z^*(\mathbf{s}_1)$ is a normal object with $\zeta = 0$, then the related local difference is contaminated by the error $-\frac{\zeta}{K}$. This leads to the swamping effect where the normal object $Z^*(\mathbf{s}_1)$ may be misclassified as an outlier. For a relatively large K (e.g., 8), it can be readily shown that Phase I contamination is more significant for a spatial sample with clusters of outliers than a spatial sample with isolated outliers. Another important observation is that the masking and swamping effects will not completely distort the ordering of true outliers. The top ranking outliers are still usually a subset of the true outliers. This observation motivates the backward algorithm presented in Section 4.3. 2) **Phase II contamination** occurs in the generalized regression process, where we regard $\mathbf{Z}^* = \mathbf{F}\mathbf{Z}$ as the pseudo "observed" values. The masking and swamping effects in this phase are the same effects occurred in a general least squares regression process. This is consequence of the biased estimates of the regression parameters (e.g., $\boldsymbol{\beta}, \sigma$, and σ_0) due to abnormal observations in \mathbf{Z}^* .

Drawbacks of existing robust estimation techniques: Most existing robust regression techniques are designed to reduce the effect of Phase II contamination. There are two major categories of estimators [13]. The first category (also called M-estimators) is to replace the MSE function by more robust score function such as L1 norm and Huber penalty function. The second category is to estimate parameters based on a robustly selected subset of data, such as least median of square (LMS), least trimmed square (LTS), and forward search (FS) method. Unfortunately, all these robust techniques cannot be directly applied to address both Phase I and Phase II contaminations concurrently. As with the M-estimators, the application of robust penalty function (e.g., L1) will lead to a non-convex optimization problem where local optimal solution may be found. With the second type of estimators based on subset selection, the estimation results are highly sensitive to the selected objects which can detrimentally impact neighborhood quality. The next sub-section will adapt existing robust methods to resolve the problem of concurrently handling Phase I and Phase II contaminations.

4.2 GLS-Backward Search Algorithm

As discussed above, existing methods only address Phase II contamination. The motivation for our proposed backward search algorithm is to address both Phase I and Phase II contaminations concurrently. The algorithm is described as follows:

Algorithm 1 (Backward search algorithm) Given a spatial data set $\{Z(\mathbf{s}_1), \dots, Z(\mathbf{s}_n)\}$, the covariate vectors $\{\mathbf{x}(\mathbf{s}_1), \dots, \mathbf{x}(\mathbf{s}_n)\}$, the value of K for defining K -nearest neighbors, and the confidence interval $\alpha \in (0,1)$,

1. Set $\mathbf{S}_Z = \{Z(\mathbf{s}_1), \dots, Z(\mathbf{s}_n)\}$, $\mathbf{S}_x = \{\mathbf{x}(\mathbf{s}_1), \dots, \mathbf{x}(\mathbf{s}_n)\}$, and \mathbf{S}_{output} be an empty set.
2. Estimate the parameters $\boldsymbol{\beta}, \sigma, \sigma_0$ of the GLS model by solving the generalized least squares regression problem (9).
3. Calculate the absolute values of standard estimated residuals
$$\mathbf{e} = [e_1, \dots, e_{|S_Z|}]^T = \left| (\sigma^2 \mathbf{I} + \sigma_0^2 \mathbf{F}\mathbf{F}^T)^{-\frac{1}{2}} (\mathbf{F}\mathbf{Z} - \mathbf{F}\mathbf{X}\boldsymbol{\beta}) \right|$$
4. Set $e_m = \max\{e_i\}_{i=1}^{|S_Z|}$.

If $e_m \geq \Phi^{-1}(\alpha/2)$, where Φ is the CDF of the standard normal distribution, then update $\mathbf{S}_Z = \mathbf{S}_Z - \{Z(\mathbf{s}_m)\}$, $\mathbf{S}_x = \mathbf{S}_x - \{\mathbf{x}(\mathbf{s}_m)\}$, and $\mathbf{S}_{output} = \mathbf{S}_{output} + \{Z(\mathbf{s}_m)\}$, and go to Step 2.

Otherwise, stop the algorithm and return \mathbf{S}_{output} as the ordered set of candidate outliers.

In the above algorithm, the confidence interval α can be set to 0.001, 0.01, and 0.05. In step 2, we apply interior point [14] method to solve the optimization problem (9). When the neighborhood size is large, we may approximate $\sigma_0^2 \mathbf{F}\mathbf{F}^T$ as $\sigma_0^2 \mathbf{I}$. The parameters $\boldsymbol{\beta}, \sigma, \sigma_0$ can be efficiently estimated by least squares regression: $\boldsymbol{\beta} = (\mathbf{X}^T \mathbf{F}^T \mathbf{F} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{F}^T \mathbf{F} \mathbf{Z}$, and $(\sigma^2 + \sigma_0^2) = \|\mathbf{F}\mathbf{X}\boldsymbol{\beta} - \mathbf{F}\mathbf{Z}\|_2^2 / (n - p - 1)$, where p is the size of the vector $\boldsymbol{\beta}$.

This backward search algorithm's design is based on the observation that top ranked outliers identified by the least squares techniques are still true outliers (in most cases) under both Phase I and II contaminations. Suppose a true outlier \mathbf{s} is removed after the first iteration, then both Phase I and Phase II contaminations in the next iteration will be reduced. To illustrate this process, we use the same example in Section 4. Recall that an outlier object $Z^*(\mathbf{s})$ is decomposed into two additive components $Z^*(\mathbf{s}) =$

$Z(\mathbf{s}) + \zeta$, where $Z(\mathbf{s})$ represents the normal value and ζ represents the contamination error. Suppose \mathbf{s} is the only outlier neighbor of an object \mathbf{s}_1 that happens to be an outlier as well. Then the local difference $\text{diff}(Z^*(\mathbf{s}_1)) = \left[Z(\mathbf{s}_1) - \frac{1}{K} \sum_{\mathbf{s}_i \in N(\mathbf{s})} (Z(\mathbf{s}_i)) \right] + \zeta_1 - \frac{\zeta}{K}$ will be marked as normal if $\zeta = K \cdot \zeta_1$. Suppose now that the true outlier $Z(\mathbf{s})$ is removed and the newly replaced neighbor for \mathbf{s}_1 is normal, then $\text{diff}(Z^*(\mathbf{s}_1)) = \left[Z(\mathbf{s}_1) - \frac{1}{K} \sum_{\mathbf{s}_i \in N(\mathbf{s})} (Z(\mathbf{s}_i)) \right] + \zeta_1$. This local difference becomes an abnormal value and the masking effect is removed. Similarly, suppose $Z^*(\mathbf{s}_1)$ is a normal object, then its local difference is contaminated (swamped) by the error $-\frac{\zeta}{K}$, because of its outlier neighbor $Z(\mathbf{s})$. The removal of \mathbf{s} will make $-\frac{\zeta}{K} = 0$, thus reducing the swamping effect. For Phase II contamination, the removal of $Z(\mathbf{s})$ leads to the removal of an abnormal difference $\text{diff}(Z^*(\mathbf{s}))$. The set of remaining local differences will therefore have less contamination. The center of the distribution is less attracted by outliers, and the distributional shape becomes less distorted. As a result, outliers tend to be more separated and normal objects tend to be closer together. The masking and swamping effects are therefore reduced.

4.3 GLS-Forward Search Algorithm

This section adapts the popular Forward Search (FR) algorithm [13] to the GLS parameters estimation problem. There are several restrictions to apply FR here. As discussed in Section 4.1, FR starts from a robustly selected subset of sample, but GLS is a statistical model based on neighborhood aggregations. Considering only a subset of the observations $\{Z(\mathbf{s}_1), \dots, Z(\mathbf{s}_n)\}$ will significantly impact the quality of the calculated local differences. To apply FR algorithm, we make the assumption that Phase I contamination is negligible compared to Phase II contamination. As discussed in Section 4.1, this is reasonable for the case of isolated outliers. Based on this assumption, we consider the local differences $\{\text{diff}(Z(\mathbf{s}_1)), \dots, \text{diff}(Z(\mathbf{s}_n))\}$ as pseudo "observations", and then apply FR algorithm to estimate the model parameters. By simulations, we also noticed that in this case there is no significant difference on accuracy between the GLS forms (6) and (7). For the sake of efficiency, we only consider the GLS form (7) and apply regular least squares regression to estimate the parameters $\boldsymbol{\beta}, \sigma$, and σ_0 . The forward search algorithm is described as follows:

Algorithm 2 (Forward Search algorithm) Given a spatial data set $\{Z(\mathbf{s}_1), \dots, Z(\mathbf{s}_n)\}$, the covariate vectors $\{\mathbf{x}(\mathbf{s}_1), \dots, \mathbf{x}(\mathbf{s}_n)\}$, and the value of K for defining K -nearest neighbors,

1. Calculate the local differences: $\text{diff}(\mathbf{Z}) = \mathbf{F}\mathbf{Z}$, and set \mathbf{S}_{output} be an empty set.
2. Set $\mathbf{S} = \{\mathbf{s}_1, \dots, \mathbf{s}_n\}$, $\mathbf{Z}^*(\mathbf{S}) = [Z^*(\mathbf{s}_1), \dots, Z^*(\mathbf{s}_n)] = \text{diff}(\mathbf{Z})$, and $\mathbf{X}^*(\mathbf{S}) = [\mathbf{x}^*(\mathbf{s}_1), \dots, \mathbf{x}^*(\mathbf{s}_n)] = \mathbf{F}\mathbf{X}$ as the vector of pseudo "observations" and pseudo "covariates".
3. Apply least trimmed squares (LTS) [13] to identify a robust subset of \mathbf{S} , defined as \mathbf{S}^* , and set $\mathbf{S}_{test}^* = \mathbf{S} - \mathbf{S}^*$. The size of the subset \mathbf{S}^* is $\lfloor (n + p + 1)/2 \rfloor$ by default.
4. Estimate the parameter $\boldsymbol{\beta}$ based on $\mathbf{Z}^*(\mathbf{S}^*)$ and $\mathbf{X}^*(\mathbf{S}^*)$. Then calculate the absolute standard residuals of \mathbf{S}_{test}^* as $\mathbf{e} = \sqrt{(n - p - 1)} | \mathbf{Z}^*(\mathbf{S}_{test}^*) - \mathbf{X}^*(\mathbf{S}_{test}^*) \boldsymbol{\beta} | / \| \mathbf{Z}^*(\mathbf{S}) - \mathbf{X}^*(\mathbf{S}) \boldsymbol{\beta} \|_2$.
5. Find the minimal residual of the test set \mathbf{S}_{test}^* :

$$e_m = \min\{e_i\}_{e_i \in S_{test}^*}$$

6. Update $\mathbf{S}_{output} = \mathbf{S}_{output} + \{\mathbf{s}_m\}$, $\mathbf{S}^* = \mathbf{S}^* + \{\mathbf{s}_m\}$, $\mathbf{S}_{test}^* = \mathbf{S}_{test}^* - \{\mathbf{s}_m\}$. If \mathbf{S}_{test}^* is not empty, go to step 4; otherwise, output the ordered set \mathbf{S}_{output} and terminate the algorithm.

The proposed *FR* algorithm provides an ordering of objects based on their agreements with the *GLS* model. To identify outliers, it plots and monitors the change of the minimal residual with the increasing size of the normal set \mathbf{S}^* . A drastic drop implies that an outlier was added to \mathbf{S}^* . This plot could also help identify masked or swamped objects. Readers are referred to [13] for details. A direct method for calculating the local differences can be achieved via robust mean functions such as median and trimmed mean. However, as indicated by our simulation study, this direct approach will deteriorate the performance of *GLS*. Recall that the statistical model of *GLS* has the form $\mathbf{diff}(\mathbf{Z}) \sim \mathbf{N}(\mathbf{F}\mathbf{X}\boldsymbol{\beta}, \mathbf{F}\boldsymbol{\Sigma}\mathbf{F}^T + \sigma_0^2\mathbf{F}\mathbf{F}^T)$. If we replace the left hand side $\mathbf{diff}(\mathbf{Z}) = \mathbf{F}\mathbf{Z}$ by medians or trimmed means, the right side will remain unchanged and thus still employs the average matrix \mathbf{F} . The increased bias caused by this inconsistency is much larger than the reduction of contamination effects achieved through robust means.

4.4 Connections with Existing Methods

This section studies the connection between global (kriging) based [11, 12, 13], local statistics (*LS*) based [4-10], and the proposed *GLS-SOD* methods. First, we review the first two approaches: *Kriging-SOD* and *LS-SOD*. *Kriging-SOD* basically applies robust methods to estimate the parameters of a global kriging model. The method then uses the estimated statistical model to predict the Z attribute value of each sample location \mathbf{s} , denoted as $\hat{Z}(\mathbf{s})$, based on the Z values of other locations. The standardized residual ($|\hat{Z}(\mathbf{s}) - Z(\mathbf{s})|/\sigma_s$) follows a standard normal distribution, where σ_s is the estimated standard deviation. If a residual is outside the range $[-\Phi^{-1}(\alpha/2), \Phi^{-1}(\alpha/2)]$, the corresponding object is reported as an outlier, where Φ is the *CDF* and α is usually set as 0.05. The *LS-SOD* method assumes that $\mathbf{diff}(\mathbf{Z}) \sim \mathbf{N}(\boldsymbol{\mu} \cdot \mathbf{1}, \sigma^2\mathbf{I})$. The components in $\mathbf{diff}(\mathbf{Z})$ can be regarded as i.i.d. sample points of a normal distribution $N(\boldsymbol{\mu}, \sigma^2)$. Robust techniques are then designed to estimate $\boldsymbol{\mu}$ and σ . The remaining steps are similar to *Kriging-SOD*.

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 estimations, then *Kriging-SOD* and *GLS-SOD* are equivalent.

Proof: For *Kriging-SOD*, we consider a universal kriging model [1], since other kriging models (e.g., ordinary kriging) are simply special cases. It suffices to prove that the standardized residuals calculated by *Kriging-SOD* and *GLS-SOD* are identical. Without loss of generality, we test the standardized residual of one particular sample point $Z(\mathbf{s}_n)$. Let $\mathbf{Z}^* = [Z(\mathbf{s}_1), \dots, Z(\mathbf{s}_{n-1})]^T$ and $\mathbf{Z} = [\mathbf{Z}^{*T}, Z(\mathbf{s}_n)]^T$. By Section 3.1 Equation 3, we have that $\mathbf{Z} \sim \mathbf{N}(\mathbf{X}\boldsymbol{\beta}, \mathbf{D})$, where $\mathbf{D} = \boldsymbol{\Sigma} + \sigma_0^2\mathbf{I} = \begin{bmatrix} \boldsymbol{\Sigma}^* & \boldsymbol{\sigma} \\ \boldsymbol{\sigma}^T & \sigma_n^2 \end{bmatrix}$, $\text{Var}(\mathbf{Z}^*) = \boldsymbol{\Sigma}^*$, $\text{Cov}(Z(\mathbf{s}_1), \mathbf{Z}^*) = \boldsymbol{\sigma}$, and $\text{Var}(Z(\mathbf{s}_n)) = \sigma_n^2$.

Then, the standard residual by *Kriging-SOD* is

$$\text{StdRsd}_{\text{Kriging-SOD}}(Z(\mathbf{s}_n)) = \frac{\mathbf{x}_n^T\boldsymbol{\beta} + \boldsymbol{\sigma}^T\boldsymbol{\Sigma}^{*-1}(\mathbf{Z}^* - \mathbf{X}^*\boldsymbol{\beta})}{\sigma_n - \boldsymbol{\sigma}^T\boldsymbol{\Sigma}^{*-1}\boldsymbol{\sigma}}$$

The standard residual by *LS-SOD* is

$$\text{StdRsd}_{\text{GLS-SOD}}(\text{diff}(Z(\mathbf{s}_n))) = \left[(\boldsymbol{\sigma}\mathbf{I} + \sigma_0^2\mathbf{F}\mathbf{F}^T)^{-\frac{1}{2}}(\mathbf{F}\mathbf{Z} - \mathbf{F}\mathbf{X}\boldsymbol{\beta}) \right]_n$$

The following will prove that

$$\text{StdRsd}_{\text{Kriging-SOD}}(Z(\mathbf{s}_n)) = \text{StdRsd}_{\text{GLS-SOD}}(\text{diff}(Z(\mathbf{s}_n)))$$

The condition $\mathbf{F}\boldsymbol{\Sigma}\mathbf{F}^T = \sigma^2\mathbf{I}$ implies $\sigma^2\mathbf{I} + \sigma_0^2\mathbf{F}\mathbf{F}^T = \mathbf{F}\boldsymbol{\Sigma}\mathbf{F}^T + \sigma_0^2\mathbf{F}\mathbf{F}^T = \mathbf{F}\mathbf{D}\mathbf{F}^T$, and $(\boldsymbol{\sigma}\mathbf{I} + \sigma_0\mathbf{F}\mathbf{F}^T)^{-\frac{1}{2}} = (\mathbf{F}\mathbf{D}\mathbf{F}^T)^{-\frac{1}{2}} = (\mathbf{F}\mathbf{D}\mathbf{z})^{-1} = \mathbf{D}^{-\frac{1}{2}}\mathbf{F}^{-1}$. It follows that $(\boldsymbol{\sigma}\mathbf{I} + \sigma_0\mathbf{F}\mathbf{F}^T)^{-\frac{1}{2}}(\mathbf{F}\mathbf{Z} - \mathbf{F}\mathbf{X}\boldsymbol{\beta}) = \mathbf{D}^{-\frac{1}{2}}\mathbf{F}^{-1}(\mathbf{F}\mathbf{Z} - \mathbf{F}\mathbf{X}\boldsymbol{\beta}) = \mathbf{D}^{-\frac{1}{2}}(\mathbf{Z} - \mathbf{X}\boldsymbol{\beta})$.

Further, given that $\mathbf{D} = \begin{bmatrix} \boldsymbol{\Sigma}^* & \boldsymbol{\sigma} \\ \boldsymbol{\sigma}^T & \sigma_n \end{bmatrix}$, it can be readily derived that

$$\mathbf{D}^{-\frac{1}{2}} = \begin{bmatrix} \left[\mathbf{C}_1^{-1} + \mathbf{C}_2^{-\frac{1}{2}}\boldsymbol{\Sigma}^{*-1}\boldsymbol{\sigma}\boldsymbol{\sigma}^T\boldsymbol{\Sigma}^{*-1} \right]^{\frac{1}{2}} & \mathbf{0} \\ -\boldsymbol{\sigma}^T\boldsymbol{\Sigma}^{*-1}\mathbf{C}_2^{-\frac{1}{2}} & \mathbf{C}_2^{-\frac{1}{2}} \end{bmatrix}$$

where $\mathbf{C}_1 = \boldsymbol{\Sigma}^{*-1} - \sigma_n\boldsymbol{\sigma}\boldsymbol{\sigma}^T$ and $\mathbf{C}_2 = \sigma_n - \boldsymbol{\sigma}^T\boldsymbol{\Sigma}^{*-1}\boldsymbol{\sigma}$.

Then, $\left[(\boldsymbol{\sigma}\mathbf{I} + \sigma_0\mathbf{F}\mathbf{F}^T)^{-\frac{1}{2}}(\mathbf{F}\mathbf{Z} - \mathbf{F}\mathbf{X}\boldsymbol{\beta}) \right]_n = \left[\mathbf{D}^{-\frac{1}{2}}(\mathbf{Z} - \mathbf{X}\boldsymbol{\beta}) \right]_n = \left[\mathbf{D}^{-\frac{1}{2}} \begin{bmatrix} \mathbf{X}^*\boldsymbol{\beta} \\ \mathbf{x}_n\boldsymbol{\beta} \end{bmatrix} \right]_n = -\mathbf{C}_2^{-\frac{1}{2}}\boldsymbol{\sigma}^T\boldsymbol{\Sigma}^{*-1}\mathbf{X}^*\boldsymbol{\beta} + \mathbf{C}_2^{-\frac{1}{2}}\mathbf{x}_n\boldsymbol{\beta} = \{\mathbf{x}_n^T\boldsymbol{\beta} + \boldsymbol{\sigma}^T\boldsymbol{\Sigma}^{*-1}(\mathbf{Z}^* - \mathbf{X}^*\boldsymbol{\beta})\}/(\sigma_n - \boldsymbol{\sigma}^T\boldsymbol{\Sigma}^{*-1}\boldsymbol{\sigma})$.

The above indicates that

$$\text{StdRsd}_{\text{Kriging-SOD}}(Z(\mathbf{s}_n)) = \text{StdRsd}_{\text{GLS-SOD}}(\text{diff}(Z(\mathbf{s}_n))),$$

We conclude that *Kriging-SOD* and *GLS-SOD* are equivalent. \square

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 estimations, 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.

Proof: For either condition (1) or (2), it can be readily derived that $\mathbf{F}\mathbf{X}\boldsymbol{\beta} = \mathbf{0}$. By conditions $\mathbf{F}\boldsymbol{\Sigma}\mathbf{F}^T = \sigma^2\mathbf{I}$ and $\sigma_0^2\mathbf{F}\mathbf{F}^T = \sigma_0^2\mathbf{I}$, we have $\mathbf{F}\mathbf{Z} \sim \mathbf{N}(0, (\sigma^2 + \sigma_0^2)\mathbf{I})$ which is consistent with the i.i.d. assumption in *LS-SOD*. If we use the same robust methods to estimate the parameters, such as using median and median absolute deviation (*MAD*) to estimate the mean and standard deviation, then *GLS-SOD* becomes equivalent to *LS-SOD*. \square

Discussion: By Theorem 6, *LS-SOD* is a special form of *GLS-SOD*. *LS-SOD* assumes $\text{Var}(\mathbf{diff}(\mathbf{Z})) = \sigma^2\mathbf{I}$ for some constant σ , but no justifications are presented. From this perspective, *GLS-SOD* actually provides a theoretical foundation for *LS-SOD*. Section 3.1 discusses the situations where $\text{Var}(\mathbf{diff}(\mathbf{Z}))$ can be approximated by $(\sigma^2 + \sigma_0^2)\mathbf{I}$. Furthermore, under the conditions of Theorem 6, *LS-SOD* is equivalent to *GLS-SOD* and since the conditions also include " $\mathbf{F}\boldsymbol{\Sigma}\mathbf{F}^T = \sigma^2\mathbf{I}$ ", then by Theorem 4 we have that *GLS-SOD* is equivalent to *Kriging-SOD*. Therefore, *LS-SOD* becomes equivalent to *Kriging-SOD* in this situation. Hence, it can be seen that the proposed *GLS* framework can be parameterized to become instances of *LS-SOD* or *Kriging-SOD*. Further study on various outlier detection methods can be greatly enhanced under the lens of this unifying *GLS* framework.

As discussed in Section 3.1, $\mathbf{F}\mathbf{S}\mathbf{F}^T$ can be reasonably approximated by $\sigma^2\mathbf{I}$. From Theorem 5, the major difference between *Kriging-SOD* and *GLS-SOD* is for which approach the related model parameters can be estimated more accurately and efficiently. From this perspective, *GLS-SOD* is superior to *Kriging-SOD* based on three major reasons: First, *GLS-SOD* has less uncertainty than *Kriging-SOD*, since *Kriging-SOD* needs to further assume a semivariogram model. If the semivariogram model is not selected properly, the performance may be significantly impacted. Second, *GLS-SOD* is a convex optimization problem and therefore a global optimal solution exists. However, *Kriging-SOD* is a non-convex optimization problem and relies on an iteratively reweighted generalized least square (*IRWGLS*) approach [12] to determine a local solution. Finally, as shown in Section 5 simulations, *GLS-SOD* runtime performance is superior to *Kriging-SOD*.

5. SIMULATIONS

This section conducts extensive simulations to compare the performance between the proposed *GLS* based *SOD* methods and other related *SOD* methods. The experimental study follows the standard statistical approach for evaluating the performance of spatial outlier detection methods presented in [11, 12, 1, 2].

5.1 Simulation Settings

Data set: The simulation data are generated based on the following statistical model:

$$Z(\mathbf{s}) = \mathbf{x}^T(\mathbf{s})\boldsymbol{\beta} + \omega(\mathbf{s}) + \epsilon(\mathbf{s}), \quad (\text{See Section 3.1})$$

where $\omega(\mathbf{s})$ is a Gaussian random field with covariogram model $C(h; \boldsymbol{\theta})$.

We consider two popular covariogram models: spherical model and exponential model. See Equation 8 in Section 3.2 for the definition of spherical model. The exponential model is defined as

$$C(h; \boldsymbol{\theta}) = \begin{cases} b & \text{if } h = 0 \\ b \left(1 - \exp\left(-\frac{h}{c}\right)\right) & \text{if } 0 < h \leq c \\ 0 & \text{if } h > c, \end{cases} \quad (11)$$

These two models have the same parameters b and c . Recall that b is also the constant variance for each $Z(\mathbf{s})$.

For the trend component $\mathbf{x}^T(\mathbf{s})\boldsymbol{\beta}$, we define $\mathbf{x}(\mathbf{s}) = [1, x(\mathbf{s}), y(\mathbf{s}), x(\mathbf{s}) \cdot y(\mathbf{s}), x(\mathbf{s})^2, y(\mathbf{s})^2]$, where $x(\mathbf{s})$ and $y(\mathbf{s})$ be the X and Y coordinates of the location \mathbf{s} . This implies that the trend $\mathbf{x}(\mathbf{s})\boldsymbol{\beta}$ is a polynomial of order two. The nonlinearity of the trend is decided by the regression parameters $\boldsymbol{\beta}$. For example, if $\boldsymbol{\beta} = [1, 0, 0, 0, 0, 0]^T$, then the trend is constant; if $\boldsymbol{\beta} = [1, 1, 1, 0, 0, 0]^T$, then the trend is linear.

For the white noise component, we employ the standard model:

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

There are three related parameters σ_0, σ_c and α . σ_0^2 is the variance of a normal white noise, σ_c^2 is the variance of contaminated error that generates outliers, and α is used to control the number of outliers. Note that it is possible the distribution $N(0, \sigma_c^2)$ will also generate some normal white noises. All true outliers must be only identified based on standard statistical test by calculating the conditional mean and standard deviation for each observation [2]. We also consider the case of clustered outliers. This can be simulated by constraining that the noises of a random cluster of $n \cdot \alpha$ points follow $N(0, \sigma_c^2)$. In the simulations,

we tested several representative settings for each parameter, which are summarized in Table 2.

Table 2: Combination of Parameter settings

Variable	Settings
n	$n \in 100, 200$. Randomly generate n spatial locations $\{\mathbf{s}_i\}_{i=1}^n$ in the range $[0, 25] \times [0, 25]$.
b, c	$b = 5; c = 5, 15, 25$
$\boldsymbol{\beta}$	For constant trend, $\beta_1 \sim U(0, 1)$ and $\beta_i = 0, i = 2, \dots, 5$; For linear trend, $\{\beta_1, \beta_2, \beta_3\} \in U(0, 1)$, $\beta_i = 0, i = 4, 5, 6$; For nonlinear trend, $\{\beta_i\}_{i=1}^6 \in U(0, 1)$.
σ_0, σ_c	$\sigma_0^2 = 2, 10; \sigma_c^2 = 20$
α	$\alpha = 0.05, 0.10, 0.15$.
K	$K = 4, 8$
Covariance model	Exponential, spherical
Outlier type	Isolated, Clustered

Outlier detection methods: We compared our methods with the state of the art local and global based *SOD* methods, including *Z-test* [4], *Median Z-test* [6], *Iterative Z-test* [5], *trimmed Z-test* [7], *SLOM-test* [8], and universal kriging (*UK*) based forward search [11, 12] (noted as *UK-forward*). Our proposed methods are identified as *GLS-backward-G*, *GLS-backward-R*, and *GLS-forward-R*. *GLS-backward-G* refers to the *GLS* backward algorithm using generalized least squares regression. *GLS-backward-R* refers to the *GLS* backward algorithm using regular least square regression (See Sections 4.2 and 4.3). The implementations of all existing methods are based on their published algorithm descriptions.

Performance metric: We tested the performance of all methods for every combination of parameter setting in Table 2. For each specific combination, we ran the experiments six times and then calculated the mean and standard deviation of accuracy for each method. To compare the accuracies of each method, we used the standard ROC curves. We further collected accuracies of top 10, 15, and 20 ranked outlier candidates for each method, and then the counts of winners are shown in Table 3. To calculate these winning counts, we used as an example the *GLS-backward-R* result in the top left cell of table 4: "47, 47, 45". This column refers to the constant trend cases. If within this particular case, we only consider the true accuracy of the top 10 candidate outliers, then the *GLS-backward-R* has "won" 47 times over all combination of parameters against all other methods. A win is given to the method that exhibits the highest accuracy. Consequently, if we consider the true accuracy of the top 20 candidate outliers, then the *GLS-backward-R* has won 45 times.

All the simulations were conducted on a PC with Intel (R) Core (TM) Duo CPU, CPU 2.80 GHz, and 2.00 GB memory. The development tool is MATLAB 2008.

5.2 Detection Accuracy

We compared the outlier detection accuracies of different methods based on different combinations of parameter settings as shown in Table 2. Six representative results are displayed in Figure 3. First we considered the detection performance between

local based methods. For a constant trend, our methods were competitive with existing techniques. For data sets exhibiting linear trends, our *GLS* algorithms achieved an average 10% improvement over existing local based methods. However, for data sets with nonlinear trends, our *GLS* algorithms exhibited more significant improvement (approximately 50% increase) over existing local methods. For the other combination of parameter settings in Table 2, the winning statistics for each method are displayed in Table 3. These results further justify the preceding performance results.

We also compared our *GLS* algorithms against the global based method *UK-forward*. Overall, our methods were comparable to *UK-forward*. Particularly, *GLS-backward-G* attained better accuracy than *UK-forward* on about half of the data sets. For the remaining data sets, the *GLS-backward-G* was still competitive to the *UK-forward*. Additionally, as shown in Section 5.3, the *UK-forward* incurred a significantly much higher computational cost than the *GLS* algorithms. As discussed in section 4.1, when K is small, the effects of $\sigma_0^2 FF^T$ must be considered and a generalized least regression is necessary. The theorems indicate that *GLS-backward-G* should perform better than *GLS-backward-R*, this was justified in Figure 3 c).

Table 3: Competition statistics for different combinations of parameter settings. Each cell contains 3 values, representing the win times for the related method 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

5.3 Computational Cost

The comparison on computational cost is shown in Figure 2. The results indicate that the time cost of *UK-forward* is much higher than other methods. Even the second slowest method *GLS-backward-G*, is still three times faster than *UK-forward*. The other local methods are approximately equal and hence much faster than *UK-forward*.

From the comparisons of both the accuracy and computational cost, it can be seen that our proposed *GLS SOD* algorithms (especially *GLS-backward-G*) is significantly more accurate than existing local based algorithms when the spatial data exhibits either a linear or nonlinear spatial trend. Our *GLS* algorithms are comparable to the global based method *UK-forward* on accuracy, but significantly faster than *UK-forward*.

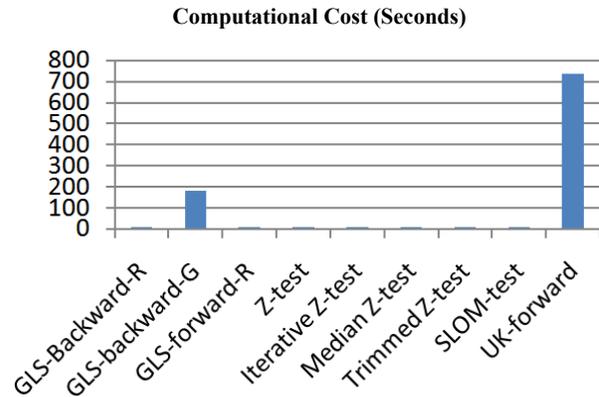


Figure 2: Comparison on computational cost (setting: Linear trend, isolated outliers, $\alpha = 0.1$, $\sigma_0^2 = 2$, $c = 15$, $K = 8$, $n = 200$)

6. CONCLUSION AND FUTURE WORK

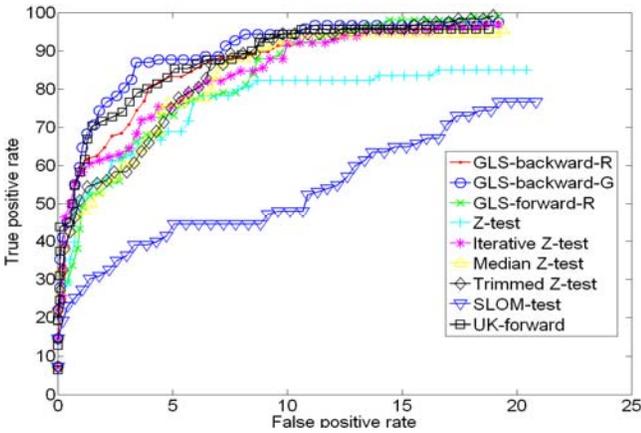
This paper presents a generalized local statistical (*GLS*) framework for existing local based methods. This generalized statistical framework not only provides theoretical foundations for local based methods, but can also significantly enhance spatial outlier detection methods. This is the first paper to present the theoretical connection between local and global based *SOD* methods under the *GLS* framework. As future work we will design other algorithms to further improve the efficiency of the *GLS* backward and forward methods.

7. REFERENCES

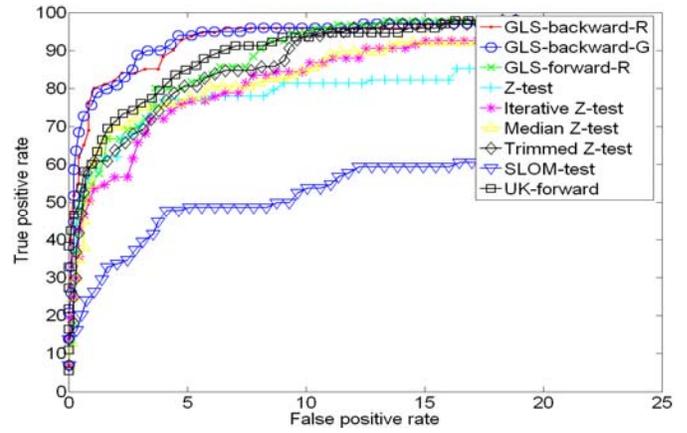
- [1] Cressie, N.A. 1993. Statistics for Spatial Data, Wiley.
- [2] Schabenberger O. and Gotway C. A. 2005 Statistical Methods for Spatial Data Analysis. Boca Raton: Chapman and Hall–CRC, Boca Raton, Florida.
- [3] Tobler, W. R. 1979. Cellular geography, in *Philosophy in Geography*, Reidel, Dordrecht, 379–386.
- [4] Shekhar, S., Lu, C.-T. and Zhang, P. 2003. A Unified Approach to Spatial Outliers Detection, *Journal of GeoInformatica*, vol. 7, No.2, 139-166.
- [5] Lu, C.-T., Chen, D. and Kou, Y. 2003. Algorithms for Spatial Outlier Detection, In *Proceedings of the 3rd IEEE International Conference on Data Mining*, 597-600.
- [6] Lu, C.-T., Chen, D. and Chen, F. 2008. On Detecting Spatial Outliers, *Journal of Geoinformatica*, vol. 12, 455-475.
- [7] Hu, T. and Sung, S.Y. 2004. A trimmed mean approach to finding spatial outliers, *J. Intell Data Anal*, vol. 8, 79-95.
- [8] Sun, P. and Chawla, S. 2004. On Local Spatial Outliers, *Proc. 4th IEEE Int'l Conf. on Data Mining*, 209–216.
- [9] S. Shekhar, Lu, C.-T. and Zhang, P. 2001. Detecting graph-based spatial outliers: algorithms and applications (a summary of results). *SIGKDD*, 365-370.
- [10] Christensen, R., Johnson, W. and Pearson, L.M., 1993. Covariance function diagnostics for spatial linear models. *Math. Geol.* vol. 25, 145–160.
- [11] Cerioli, A. and Riani, M. 1999. The ordering of spatial data and the detection of multiple outliers. *J. Comput. Graphical Statist.* vol. 8, 239–258.
- [12] Militino, A.F., Palacios, M.B. and Ugarte, M.D. 2006. Outlier detection in multivariate spatial linear models, *J. Stat Plann Infer*, vol. 136, 125-146.

[13] Atkinson, A.C. and Riani, M. 2000. Robust Diagnostics Regression Analysis. Springer Series in Statistics. Springer.
 [14] S. Boyd and L. Vanderberghe. 2004. Convex Optimization. Cambridge Univ. Press.

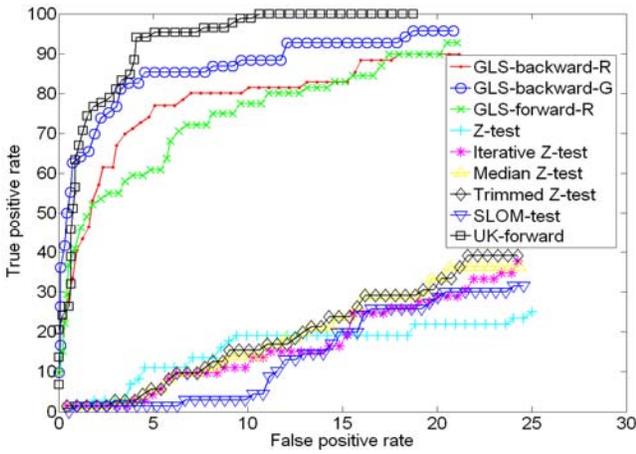
[15] Chen, F, Lu, C-T, and Boedihardjo, Arnold P., 2010. GLS-SOD: A Generalized Local Statistical Approach for Spatial Outlier Detection. Technical Report TR-10-03, Computer Science, Virginia Tech.



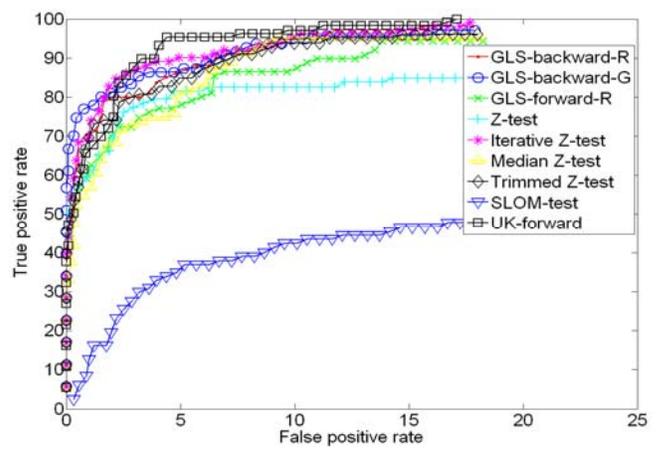
a) Constant trend, isolated outliers, $\alpha = 0.1, \sigma_0^2 = 2, c = 15, K = 4$



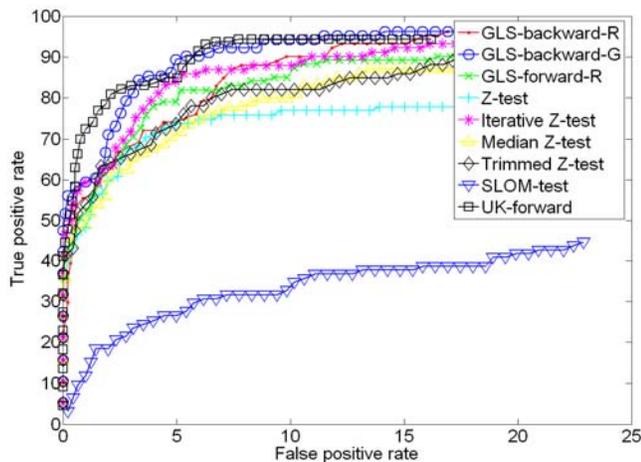
b) Linear trend, isolated outliers, $\alpha = 0.1, \sigma_0^2 = 2, c = 15, K = 8$



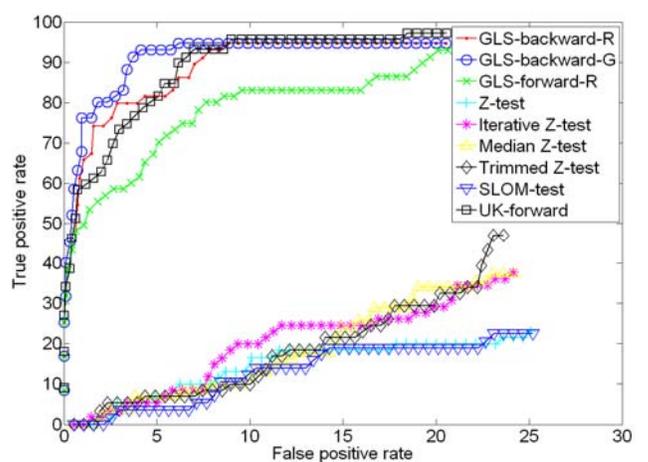
c) Nonlinear trend, isolated outliers, $\alpha = 0.15, \sigma_0^2 = 10, c = 15, K = 4$



d) Constant trend, cluster outliers, $\alpha = 0.1, \sigma_0^2 = 2, c = 25, K = 4$



e) Linear trend, cluster outliers, $\alpha = 0.15, \sigma_0^2 = 2, c = 25, K = 8$



f) Nonlinear trend, cluster outliers, $\alpha = 0.15, \sigma_0^2 = 10, c = 5, K = 8$

Figure 3: Outlier ROC Curve Comparison (the same setting: $n = 200, b = 5, \sigma_c^2 = 20$)