

Article

Detection of Influential Observations in Spatial Regression Model Based on Outliers and Bad Leverage Classification

Ali Mohammed Baba ^{1,2}, Habshah Midi ^{1,3,*}, Mohd Bakri Adam ^{1,3} and Nur Haizum Abd Rahman ^{1,3}

¹ Institute for Mathematical Research, Universiti Putra Malaysia, Serdang 43400, Selangor, Malaysia; mbali@atbu.edu.ng (A.M.B.); bakri@upm.edu.my (M.B.A.); nurhaizum@upm.edu.my (N.H.A.R.)

² Department of Mathematical Sciences, Abubakar Tafawa Balewa University, Bauchi 0248, Nigeria

³ Department of Mathematics and Statistics, Faculty of Science, Universiti Putra Malaysia, Serdang 43400, Selangor, Malaysia

* Correspondence: habshah@upm.edu.my

Abstract: Influential observations (IOs), which are outliers in the x direction, y direction or both, remain a problem in the classical regression model fitting. Spatial regression models have a peculiar kind of outliers because they are local in nature. Spatial regression models are also not free from the effect of influential observations. Researchers have adapted some classical regression techniques to spatial models and obtained satisfactory results. However, masking or/and swamping remains a stumbling block for such methods. In this article, we obtain a measure of spatial Studentized prediction residuals that incorporate spatial information on the dependent variable and the residuals. We propose a robust spatial diagnostic plot to classify observations into regular observations, vertical outliers, good and bad leverage points using a classification based on spatial Studentized prediction residuals and spatial diagnostic potentials, which we refer to as $ISRs - P_{osi}$ and $ESRs - P_{osi}$. Observations that fall into the vertical outliers and bad leverage points categories are referred to as IOs. Representations of some classical regression measures of diagnostic in general spatial models are presented. The commonly used diagnostic measure in spatial diagnostics, the Cook's distance, is compared to some robust methods, H_i^2 (using robust and non-robust measures), and our proposed $ISRs - P_{osi}$ and $ESRs - P_{osi}$ plots. Results of our simulation study and applications to real data showed that the Cook's distance, non-robust H_{si1}^2 and robust H_{si2}^2 were not very successful in detecting IOs. The H_{si1}^2 suffered from the masking effect, and the robust H_{si2}^2 suffered from swamping in general spatial models. Interestingly, the results showed that the proposed $ESRs - P_{osi}$ plot, followed by the $ISRs - P_{osi}$ plot, was very successful in classifying observations into the correct groups, hence correctly detecting the real IOs.

Citation: Baba, A.M.; Midi, H.; Adam, M.B.; Rahman, N.H. A. Detection of Influential Observations in Spatial Regression Model Based on Outliers and Bad Leverage Classification. *Symmetry* **2021**, *13*, 2030. <https://doi.org/10.3390/sym13112030>

Academic Editor: Jinyu Li

Received: 6 August 2021

Accepted: 22 October 2021

Published: 27 October 2021

Keywords: spatial regression model; influential observation; outlier; leverage; prediction residual; masking and swamping; diagnostic

Publisher's Note: MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Copyright: © 2021 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

1. Introduction

Belsley et al. [1] defined an influential observation (IO) as one which, either individually or together with several other observations, has a demonstrably large impact on the calculated values of various estimates. An influential observation could be an outlier in the X -space (leverage points) or outlier in the Y -space (vertical outlier). Leverage points can be classified into good (GLPs) and bad leverage points (BLPs). Unlike BLPs, GLPs follow the pattern of the majority of the data; hence, they are not considered as IOs as they have little or no influence on the calculated values of numerous estimates [2,3]. In this connection, Rashid et al. [2] stated that IOs could be vertical outliers (VO) or BLPs. Thus, it is very crucial to identify IOs as they are responsible for misleading conclusions about the fitted regression models and various other estimates. Once the IOs are identified, there is a need to study their impact on the model and subsequent analyses. There is a handful

of studies on the diagnostic of IOs in linear regression; some examples are [1,3–12]. Other articles in the literature deal with regressions with correlated residuals, e.g., [13–17]. However, only a few articles deal with the detection of IOs in spatial regression models; some examples include [18–22]. Some robust estimation methods in spatial regression are [23–25]. Christensen et al. [18] and Haining [19] adapted one of the diagnostic measures in [3] to detect influential observations in spatial error autoregression model. They achieved this by defining correlated errors through the spatial weight matrix and coefficient of spatial autocorrelation in the error term. They also presented the spatial Studentized prediction residuals and the spatial leverage terms that contain error terms in spatial information.

The presence of high or low attribute value in the neighbourhood of a spatial location may result in the inability to detect the true spatial outlier, or the false identification of a good observation as an outlier [26]. Hadi [27] has also noted that spatial outlier detection methods inherit the problem of masking and swamping. Masking occurs when outlying observations are incorrectly declared as inliers. Swamping on the other hand, occurs when clean observations are incorrectly classified as outliers [28]. Aggarwal [29] observed that spatial outlier breaks the spatial autocorrelation and continuity of spatial locations. Spatial autocorrelation is a systematic pattern in attribute values that are recorded in several locations on a map. Attribute values in one location that are associated with values at neighbouring locations indicate the presence of autocorrelation. Positive autocorrelation indicates similar values that are clustered together. Negative autocorrelation indicates low attribute values in the neighbourhood of high attribute values and vice-versa [30].

Robust estimation methods mostly focus on estimations that are not influenced much by the effects of outliers. Anselin [23] has extended the bootstrap estimation to mixed-regressive spatial autoregressive models, where pseudo error terms are generated by sampling from the vector of error terms. The spatial structure of the data is maintained through the generation of error terms. Politis et al. [31] and Heagerty and Lumley [32] also adopted the bootstrap method on blocks of contiguous locations to generate replicates of the estimates of the asymptotic standard error of statistics. Cerioli and Riana [24] argued that a robust estimator of the spatial autocorrelation parameters did not exist based on all datasets. They proposed a forward search algorithm based on blocks of contiguous spatial locations (BFS). The BFS algorithm are drawn in such a way that the blocks retain the spatial dependence structure of the original data. Yildirim [25] proposed a robust estimation method of the log-likelihood with influence function in the spatial error model. This is achieved iteratively using scoring algorithm to estimate the parameters. Though they succeeded in obtaining robust estimates, identifying spatial outliers, which is vital in spatial statistics [26], was not achieved. Popular graphical techniques to detect spatial outliers are the scatterplot [33], the Moran's scatterplot [30] and the pockets of nonstationarity [34]. Besides being prone to the problem of masking and swamping [26], they focused mainly on spatial outliers in the Y-space only.

Diagnostic works on models that have both spatial autocorrelations in dependent variable and residual terms are missing in the literature. The problem of masking and swamping is prevalent in spatial regression model diagnostics, which may be due to the presence of vertical outliers as well as leverage points, as in the case of linear regression ([27]). This motivates us to represent the spatial Studentized prediction residuals and spatial leverage values in the general spatial model, and to adapt and extend some robust diagnostic measures of detection of outliers and IOs in linear regression, such as Hadi's potential (p_{oi}), Cook's distance (CD_i) [3], the overall potential influence (H_i^2) [10], and the external (ESRs) and internal (ISRs) Studentized residuals [1,9,10], to spatial regression models in order to minimize the problem of masking and swamping in spatial models.

In this article, we propose a robust spatial diagnostic plot and adapt some diagnostic measures in the linear regression model. Representations of the diagnostic measures in the spatial regression model are obtained, with a special emphasis on the general spatial regression model (GSM) that performs autoregression on both the dependent variable and error terms.

The main objective of this study is to propose a robust spatial diagnostic plot. Other objectives are: (1) to represent the leverage values of the hat matrix of the linear regression in the GSM model; (2) to extend the ISR of the linear regression to the GSM model; (3) to extend the *ESR* of the linear regression to the GSM model; (4) to extend the Cook's distance and the overall potential influence of the linear regression to the GSM model (5) to develop a method of identification of the influential observations of the GSM model by proposing a procedure of classification of the observations into regular observations, vertical outliers, good and bad leverage points, and hence IOs; (6) to evaluate the performances of the proposed methods by using simulation studies; (7) to apply the proposed methods on gasoline price data for retail sites in Sheffield, UK, COVID-19 data in Georgia, USA, and the life expectancy data from USA counties. The significance of this study is that it can contribute to the development of a method of identification of influential observations in spatial regression models.

2. Identification of Influential Observations in a Linear Regression Model

Consider a k -variable regression model:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad (1)$$

where \mathbf{y} is an $n \times 1$ vector of observations of dependent variables, \mathbf{X} is an $n \times k$ matrix of independent variables, $\boldsymbol{\beta}$ is a $k \times 1$ vector of unknown regression parameters, $\boldsymbol{\varepsilon}$ is an $n \times 1$ vector of random errors with identical normal distributions, that is, $\boldsymbol{\varepsilon} \sim NID(0, \sigma^2)$.

The ordinary least squares (OLS) estimates in Equation (1) are given by:

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (2)$$

The vector of predicted values can be written as:

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{P}\mathbf{y},$$

where $\mathbf{P} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ is the hat/leverage matrix. The diagonal elements of the leverage matrix are called the hat values, denoted as p_{ii} , and given by:

$$p_{ii} = \mathbf{x}_i^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i, \quad i = 1, 2, \dots, n.$$

The hat matrix is often used as diagnostics to identify leverage points. Leverage is the amount of influence exerted by the observed response y_i on the predicted variable \hat{y}_i . As a result, a large leverage value indicates that the observed response has a large effect on the predicted response.

Hoaglin and Welsh [3] suggested that an observation which exceeds $\frac{2k}{n}$, where $\frac{2k}{n}$ is the average value of p_{ii} , is considered as a leverage point, while Vellman and Welsh suggested $\frac{3k}{n}$ as a cut-off point for leverage points. Huber [7] suggested that the ranges $p_{ii} \leq 0.2$, $0.2 < p_{ii} \leq 0.5$ and $p_{ii} > 0.5$ are safe, risky and to be avoided, respectively, for leverage values.

Unfortunately, the hat matrix suffers from the masking effect. As a result, p_{ii} often fails to detect high leverage points. Hadi [10] suggested a single-case-deleted measure called potentials or Hadi's potentials. The diagonal element of a potential denoted as p_{0ii} is given by:

$$p_{0ii} = \mathbf{x}_i^T (\mathbf{X}_{(i)}^T \mathbf{X}_{(i)})^{-1} \mathbf{x}_i, \quad i = 1, 2, \dots, n \quad (3)$$

where $\mathbf{X}_{(i)}$ is the matrix \mathbf{X} with the i^{th} row deleted. We can rewrite p_{0ii} as a function of p_{ii} as:

$$p_{0ii} = \frac{p_{ii}}{1 - p_{ii}}, \quad i = 1, 2, \dots, n.$$

The vector of the residuals, \mathbf{r} , can be written as:

$$\mathbf{r} = \mathbf{y} - \hat{\mathbf{y}} = (\mathbf{I} - \mathbf{P})\mathbf{y} = \mathbf{Q}\mathbf{y},$$

The Studentized residuals (internally Studentized residuals) denoted as ISRs and R-Student residuals (externally Studentized residuals) denoted as ESRs are widely used measures for the identification of outliers (see [7]). The *ISR*, denoted as t_i , is defined as:

$$t_i = \frac{r_i}{\hat{\sigma}\sqrt{1 - p_{ii}}}$$

where $\hat{\sigma}$ is the standard deviation of the residuals, r_i and p_{ii} are the i^{th} residual and diagonal element of the matrix \mathbf{P} , respectively (see [9] for details). Meanwhile, Chatterjee and Hadi [9] defined *ESR* denoted as t_i^* and given by:

$$t_i^* = \frac{r_i}{\hat{\sigma}_{(i)}\sqrt{1 - p_{ii}}}$$

where $\hat{\sigma}_{(i)}$ is the residuals mean square excluding the i^{th} case. The *ESR* follows a Student's t -distribution with $(n - k - 1)$ degrees of freedom [9].

One of the most employed measures of influence in linear regression is the Cook's distance [3]. It measures the influence on the regression coefficient estimate or the predicted values. The Cook's distance is given by

$$\widehat{CD}_i(\mathbf{X}^T\mathbf{X}, k\hat{\sigma}^2) = \frac{(\hat{\boldsymbol{\beta}}^{(-i)} - \hat{\boldsymbol{\beta}})^T (\mathbf{X}^T\mathbf{X}) (\hat{\boldsymbol{\beta}}^{(-i)} - \hat{\boldsymbol{\beta}})}{k\hat{\sigma}^2}, \quad (4)$$

where $\hat{\boldsymbol{\beta}}$ is the vector of estimates of $\boldsymbol{\beta}$ using the full data, $\hat{\boldsymbol{\beta}}^{(-i)}$ is the vector of estimates of $\boldsymbol{\beta}$ with the i^{th} observation of y_i and x_i omitted, k is the number of parameters and $\hat{\sigma}^2$ is the estimate of variance. Any i^{th} observation is declared influential observation (IO) if $\widehat{CD}_i > F[0.5; k, (n - k)]$. Meloun [12] noted that any observation in which $CD_i > 1$ is considered as an influential observation. The Cook's distance can also be written as [8,9]:

$$\widehat{CD}_i(\mathbf{X}^T\mathbf{X}, k\hat{\sigma}^2) = \frac{(\hat{\mathbf{y}} - \hat{\mathbf{y}}_i)^T (\hat{\mathbf{y}} - \hat{\mathbf{y}}_i)}{k\hat{\sigma}^2} \quad (5)$$

Computing the $\widehat{CD}_i(\mathbf{X}^T\mathbf{X}, k\hat{\sigma}^2)$ does not require fitting a regression equation for each of the i^{th} observations and the full model; instead, Equation (3) can further be simplified as ([3,8,9]):

$$\widehat{CD}_i(\mathbf{X}^T\mathbf{X}, k\hat{\sigma}^2) = \frac{1}{k} t_i^2 \frac{p_{ii}}{q_{ii}} \quad (6)$$

where $t_i = \frac{e_i}{\hat{\sigma}\sqrt{q_{ii}}}$ is the *ISR* and $\frac{p_{ii}}{q_{ii}}$ ($q_{ii} = 1 - p_{ii}$) is referred to as the potential [7–9]. Interestingly, the Cook's distance is a measure of influence based on the potential ($\frac{p_{ii}}{q_{ii}}$) and Studentized residual (t_i).

Hadi [10] demonstrated the drawback of methods that are multiplicative of functions, such as the Cook's distance [3], Andrews–Pregibon statistic [5], Cook and Weisberg statistic [8], etc. (see [10] for details), and proposed a method that is additive of the functions. Though both the multiplicative and additive methods are functions of residuals and leverage values, the former diminishes towards zero for smaller value of any of the two functions or both, while in the latter case, the measure is large if one of the two functions or both are large. He proposed a measure of overall potential influence, denoted as H_i^2 , and defined as follows:

$$H_i^2 = \frac{k}{m} \frac{\mathbf{e}_I^T (\mathbf{I}_m - \mathbf{P}_I)^{-1} \mathbf{e}_I}{\mathbf{e}^T \mathbf{e} - \mathbf{e}_I^T \mathbf{e}_I} + \frac{1}{m} \text{tr}(\mathbf{P}_I (\mathbf{I}_m - \mathbf{P}_I)^{-1}), \quad (7)$$

with k , the number of the parameters in the model, $\mathbf{I} = \{i_1, i_2, \dots, i_m\}$ the set of indices of observations of length m , and \mathbf{P}_I the leverage indexed by \mathbf{I} .

For $m = 1$ and $\mathbf{I} = i$, Equation (7) simplifies to:

$$H_i^2 = \frac{k}{(1 - p_{ii})} \frac{e_i^2}{(\mathbf{e}^T \mathbf{e} - e_i^2)} + \frac{p_{ii}}{1 - p_{ii}} = \frac{k}{(1 - p_{ii})} \frac{d_i^2}{(1 - d_i^2)} + \frac{p_{ii}}{1 - p_{ii}}, \quad (8)$$

where $\sum p_{ii} = k$, $\sum d_i^2 = 1$, $d_i^2 = \frac{e_i^2}{\mathbf{e}^T \mathbf{e}}$ is the square of the i^{th} normalized residual.

Hadi [10] suggested a cut-off point for Hadi's potential (p_{oi}) and H_i^2 denoted as (l_1) which is given as follows:

$$\begin{aligned} l_1 &= \text{mean}(\mathbf{p}_{oi}) + c\sqrt{\text{Var}(\mathbf{p}_{oi})} \\ &= \frac{k}{n} + c \sqrt{\frac{ns - k^2}{n(n-1)}}, \end{aligned}$$

where $c = 2, 3$, $s = \sum p_{ii}$ and \mathbf{p}_{oi} is the vector of Hadi's potential. Since both the mean and the standard deviation are easily affected by outliers, he suggested to employ such a confidence-bound type of cut-off points by replacing the mean and the standard deviation by robust estimators, namely the median and normalized median absolute deviation, respectively. The resulting cut-off point is denoted as l_2 ;

$$l_2 = \text{Med}(\mathbf{p}_{oi}) + c\text{MAD}(\mathbf{p}_{oi}),$$

3. Influential Observations in Spatial Regression Models

The general spatial autoregressive model (GSM) ([21,35,36]) includes the spatial lag term and spatially correlated error structure. The data generating process (DGP) of the general spatial model is given by:

$$\mathbf{y} = \rho \mathbf{W}_1 \mathbf{y} + \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\xi}, \quad \boldsymbol{\xi} = \lambda \mathbf{W}_2 \boldsymbol{\xi} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim N(0, \sigma^2 \mathbf{I}_n), \quad (9)$$

where \mathbf{y} is an $n \times 1$ vector of dependent variables. \mathbf{X} is an $n \times k$ matrix of explanatory variables. \mathbf{W}_1 and \mathbf{W}_2 are $n \times n$ spatial weight matrices. \mathbf{I}_n is an $n \times n$ identity matrix. $\boldsymbol{\xi}$ is the spatially correlated error term, and $\boldsymbol{\varepsilon}$ is the random residual term. The parameter ρ is the coefficient of the spatially lagged dependent variables $\mathbf{W}_1 \mathbf{y}$, and λ is the coefficient of the spatially correlated errors.

The general spatial autoregressive model in Equation (9) can be rewritten as:

$$\mathbf{A} \mathbf{y} = \mathbf{X} \boldsymbol{\beta} + \mathbf{B}^{-1} \boldsymbol{\varepsilon}, \quad (10)$$

where $\mathbf{A} = \mathbf{I}_n - \rho \mathbf{W}_1$, $\boldsymbol{\xi} = \mathbf{B}^{-1} \boldsymbol{\varepsilon}$, $\mathbf{B} = \mathbf{I}_n - \lambda \mathbf{W}_2$, $\boldsymbol{\xi} \sim N(0, \sigma^2 \mathbf{V})$, and $\mathbf{V} = (\mathbf{B}^T \mathbf{B})^{-1}$. Estimation of the parameters is achieved using the maximum likelihood estimation method.

The log-likelihood function (L) is given by:

$$L = -\frac{n}{2} \ln(\sigma^2) + \ln|\mathbf{A}| + \ln|\mathbf{B}| - \frac{1}{2\sigma^2} (\mathbf{A} \mathbf{y} - \mathbf{X} \boldsymbol{\beta})^T \mathbf{B}^T \mathbf{B} (\mathbf{A} \mathbf{y} - \mathbf{X} \boldsymbol{\beta}) \quad (11)$$

Let $\hat{\rho}, \hat{\lambda}, \hat{\sigma}^2, \hat{\boldsymbol{\beta}}$ be the maximum likelihood estimates (MLEs) of $\rho, \lambda, \sigma^2, \boldsymbol{\beta}$, respectively. The MLEs are obtained iteratively using numerical methods in the maximum likelihood estimation. Anselin [35] and LeSage [36] discussed the maximum likelihood estimation procedure of the parameters.

3.1. Leverage in Spatial Regression Model

Denote the vector of parameters in Equation (11) as $\boldsymbol{\beta}_{ay}$. The estimate of $\boldsymbol{\beta}_{ay}$, $\hat{\boldsymbol{\beta}}_{ay}$, is given by:

$$\hat{\boldsymbol{\beta}}_{ay} = (\mathbf{X}^T \hat{\mathbf{V}}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \hat{\mathbf{V}}^{-1} \hat{\mathbf{A}} \mathbf{y}.$$

The model (11) is viewed as fitting a general linear model, $\mathbf{A} \mathbf{y}$ on \mathbf{X} , that has correlated residual terms. Set $\mathbf{z} = \mathbf{A} \mathbf{y}$, where $\text{var}(\mathbf{A} \mathbf{y}) = \sigma^2 \mathbf{V}$. Therefore,

$$\hat{\mathbf{z}} = \mathbf{X} \hat{\boldsymbol{\beta}}_{ay} = \mathbf{X} (\mathbf{X}^T \hat{\mathbf{V}}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \hat{\mathbf{V}}^{-1} \mathbf{z} = \mathbf{P}_{ay} \mathbf{z}$$

The hat matrix, in this case, is given by P_{ay} ,

$$P_{ay} = X(X^T \hat{V}^{-1} X)^{-1} X^T \hat{V}^{-1}.$$

Let $Q_{ay} = I_n - P_{ay}$. Though P_{ay} and Q_{ay} have satisfied the idempotence property and their sum of diagonal elements equals k and $n - k$, respectively, they are not symmetric. As a result, they are not positive semi-definite, and as such, the diagonal elements of P_{ay} will have negative values. The hat matrices P_{ay} and Q_{ay} are not symmetric, and their diagonal values do not lie between 0 and 1 (inclusive).

Martins [15] proposed a measure of leverage that is orthogonal, in the models with correlated residuals, whose diagonal values lie in the interval $[0, 1]$, which we denote by P_{ay}^* , such that:

$$P_{ay}^* = \hat{V}^{-1} P_{ay} = \hat{V}^{-1} X(X^T \hat{V}^{-1} X)^{-1} X^T \hat{V}^{-1}$$

Let $Q_{ay}^* = I_n - P_{ay}^*$. P_{ay}^* and Q_{ay}^* are idempotent, symmetric and orthogonal with respect to V , i.e.,

1. $P_{ay}^* \hat{V} P_{ay}^* = P_{ay}^*$
2. $Q_{ay}^* \hat{V} Q_{ay}^* = Q_{ay}^*$
3. $P_{ay}^* \hat{V} Q_{ay}^* = 0$

Note that the sum of the diagonal elements of P_{ay}^* and Q_{ay}^* , the leverage, does not sum to k and $n - k$.

Again, consider a new set of dependent variables obtained by pre-multiplying Equation (11) by the matrix B (B as defined in Equation (10)) so that $z^* = BAy$. Schall and Dunne [14] defined the matrix V^{-1} as a singular value decomposition such that $V^{-1} = B\Delta B^T$; where B is of the same order as V^{-1} and Δ is a diagonal matrix. The transformation z^* is the principal component score. Puterman [13] and Haining [19] defined it as canonical variates such that $BX(X^T V^{-1} X)^{-1} X^T B^T$ is positive semi-definite. By setting $z^* = BAy$, Equation (9) is rewritten in a generalized least squares (GLS) form as:

$$z^* = X^* \beta_s + \varepsilon, \varepsilon \sim N(0, \sigma^2 I_n) \quad (12)$$

where $X^* = BX$.

The estimate $\hat{\beta}_s$ of β_s is now given by:

$$\hat{\beta}_s = (X^{*T} X^*)^{-1} X^{*T} z^*$$

Thus,

$$\hat{z}^* = X^* (X^{*T} X^*)^{-1} X^{*T} z^* \quad (13)$$

where, $\hat{A} = I_n - \hat{\rho}W_1$ and $\hat{B} = I_n - \hat{\lambda}W_2$. Note that \hat{y} is deduced from Equation (13) as follows:

$$\begin{aligned} \hat{B}\hat{A}\hat{y} &= \hat{B}X(X^T \hat{B}^T \hat{B}X)^{-1} X^T \hat{B}^T \hat{B}\hat{A}y \\ \xrightarrow{\text{yields}} \hat{y} &= \hat{A}^{-1} X(X^T \hat{V}^{-1} X)^{-1} X^T \hat{V}^{-1} \hat{A}y \end{aligned}$$

Denote the projection matrix in the transformed spatial regression model as P_s , then:

$$\begin{aligned} P_s &= X^* (X^{*T} X^*)^{-1} X^{*T} \\ &= \hat{B}X(X^T \hat{V}^{-1} X)^{-1} X^T \hat{B}^T, \quad \hat{V} = (\hat{B}^T \hat{B})^{-1} \end{aligned}$$

The properties of the leverage in the transformed spatial model in Equation (13) are:

Property I: idempotent and symmetric.

Property Ia: idempotence

$$\begin{aligned}
\mathbf{P}_s^2 &= \widehat{\mathbf{B}}\mathbf{X}(\mathbf{X}^T\widehat{\mathbf{V}}^{-1}\mathbf{X})^{-1}\mathbf{X}^T\widehat{\mathbf{B}}^T\widehat{\mathbf{B}}\mathbf{X}(\mathbf{X}^T\widehat{\mathbf{V}}^{-1}\mathbf{X})^{-1}\mathbf{X}^T\widehat{\mathbf{B}}^T \\
&= \widehat{\mathbf{B}}\mathbf{X}(\mathbf{X}^T\widehat{\mathbf{V}}\mathbf{X})^{-1}\mathbf{X}^T\widehat{\mathbf{B}}^T \\
&= \mathbf{P}_s
\end{aligned}$$

Hence, \mathbf{P}_s is idempotent.

Property Ib: symmetric

$$\begin{aligned}
\mathbf{P}_s^T &= \left(\widehat{\mathbf{B}}\mathbf{X}(\mathbf{X}^T\widehat{\mathbf{V}}^{-1}\mathbf{X})^{-1}\mathbf{X}^T\widehat{\mathbf{B}}^T\right)^T \\
&= \widehat{\mathbf{B}}\mathbf{X}(\mathbf{X}^T\widehat{\mathbf{V}}^{-1}\mathbf{X})^{-1}\mathbf{X}^T\widehat{\mathbf{B}}^T \\
&= \mathbf{P}_s
\end{aligned}$$

The matrix \mathbf{P}_s is symmetric. Therefore, \mathbf{P}_s in the transformation $\mathbf{z}^* = \widehat{\mathbf{B}}\widehat{\mathbf{A}}\mathbf{y}$ is both idempotent and symmetric.

Property II: the sum of the diagonal terms of the projection matrix is k , the number of parameters including the constant term.

$$\begin{aligned}
\text{trace}(\mathbf{P}_s) &= \text{trace}\left(\widehat{\mathbf{B}}\mathbf{X}(\mathbf{X}^T\widehat{\mathbf{V}}^{-1}\mathbf{X})^{-1}\mathbf{X}^T\widehat{\mathbf{B}}^T\right) \\
&= \text{trace}\left(\widehat{\mathbf{B}}^T\widehat{\mathbf{B}}\mathbf{X}(\mathbf{X}^T\widehat{\mathbf{V}}^{-1}\mathbf{X})^{-1}\mathbf{X}^T\right) \text{ (cyclic permutation of trace of matrix)} \\
&= \text{trace}\left(\mathbf{X}^T\widehat{\mathbf{V}}^{-1}\mathbf{X}(\mathbf{X}^T\widehat{\mathbf{V}}^{-1}\mathbf{X})^{-1}\right) \text{ (cyclic permutation of trace of matrix)} \\
&= \text{trace}(\mathbf{I}_k) \\
&= k,
\end{aligned}$$

where \mathbf{I}_k is an $k \times k$ identity matrix.

Therefore, $\sum_{i=1}^k ps_{ii} = k$. ps_{ii} is the i^{th} diagonal element of the leverage \mathbf{P}_s .

Property III: bounds on the spatial leverage.

The bound on the leverage of the classical regression is $0 \leq p_{ii} \leq 1$ due to the fact that the hat matrix \mathbf{P} satisfies all the orthogonal properties, including symmetry. As such, it is positive semi-definite. However, the spatial leverage \mathbf{P}_{ay} is not symmetric because positive semi-definite matrix is symmetric [37–39]. The transformation in Equation (11) yields the projection \mathbf{P}_s that satisfies the symmetry condition.

From the idempotent property of \mathbf{P}_s ,

$$\mathbf{P}_s = \mathbf{P}_s^2.$$

Equating diagonal terms of LHS and RHS, we have:

$$ps_{ii} = ps_{ii}^2 + \sum_{j \neq i} ps_{ij} ps_{ji}, \quad \sum_{j \neq i} ps_{ij} ps_{ji} \geq 0, \quad (14)$$

where ps_{ij} are the off-diagonal terms. Equation (14) implies that $ps_{ii} \geq 0$. Therefore,

$$\begin{aligned}
ps_{ii} &\geq ps_{ii}^2 \\
&\xrightarrow{\text{yields}} ps_{ii} \leq 1.
\end{aligned}$$

Note that \mathbf{P}_s and \mathbf{Q}_s are orthogonal:

1. $\mathbf{P}_s\mathbf{P}_s = \mathbf{P}_s$
2. $\mathbf{Q}_s\mathbf{Q}_s = \mathbf{Q}_s$
3. $\mathbf{P}_s\mathbf{Q}_s = \mathbf{0}$

The model in Equation (9) gives rise to different special spatial regressions in accordance with different restrictions. Such special spatial regression models are the spatial autoregressive-regressive model (SAR) and the spatial error model (SEM). While the former has spatial autoregression in the response variable, the latter has spatial autoregression in the model residual; model (9) (GSM) combines both features.

The spatial autoregressive-regressive model is obtained when the coefficient of the lagged spatial autoregression in the residuals of Equation (9) is zero, i.e., $\lambda = 0$. Thus, the SAR model is given by:

$$\mathbf{y} = \rho \mathbf{W}_1 \mathbf{y} + \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim N(0, \sigma^2 \mathbf{I}_n). \quad (15)$$

The \mathbf{P}_s corresponding to the model in Equation (13) reduces to:

$$\mathbf{P}_s = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T,$$

with the transformation in Equation (11) simplifying to $\mathbf{z}^* = \mathbf{A}\mathbf{y}$, since $\mathbf{V} = (\mathbf{B}^T \mathbf{B})^{-1}$ and $\mathbf{B} = \mathbf{I}_n$, when $\lambda = 0$. Clearly, the hat matrix in the SAR model preserves the features of the hat matrix in the classical regression model.

In the spatial error model (SEM), the coefficient of the spatial autoregression on the lagged dependent variable is zero, i.e., $\rho = 0$. This yields the model:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\xi}, \quad \boldsymbol{\xi} = \lambda \mathbf{W}_2 \boldsymbol{\xi} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim N(0, \sigma^2 \mathbf{I}_n) \quad (16)$$

The transformation in Equation (11) simplifies to $\mathbf{z}^* = \mathbf{B}\mathbf{y}$, and the projection matrix remains:

$$\mathbf{P}_s = \mathbf{B}\mathbf{X}(\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{B}^T.$$

It can be observed that the leverage measure in the spatial regression model is dominated by the autocorrelation in the residual term.

Works on spatial regression diagnostics in the literature mainly focus on the autocorrelation in the residuals, mostly using a time series analogy [13–15]. Some remarkable works on the spatial regression model can be found in [18,19,21].

3.2. Influential Observations in Spatial Regression Model

The leverages \mathbf{P}_s and \mathbf{Q}_s in Equation (11) satisfy all the properties of a projection matrix, including that the sum of the diagonal terms of \mathbf{P}_s and \mathbf{Q}_s equal k and $n - k$, respectively. It also incorporates the autocorrelation in the dependent variables, $\mathbf{W}\mathbf{y}$. Hence, it can be used as a diagnostic measure of leverage points in a spatial regression model.

By extending the results of linear regression to spatial regression with slight modification, the Cook's distance in the spatial regression of Equation (13), denoted as CD_{si} , can be formulated as follows:

$$\begin{aligned} \widehat{CD}_{si} &= \frac{(\widehat{\boldsymbol{\beta}}_s^{(-i)} - \widehat{\boldsymbol{\beta}}_s)^T (\mathbf{X}^{*T} \mathbf{X}^*) (\widehat{\boldsymbol{\beta}}_s^{(-i)} - \widehat{\boldsymbol{\beta}}_s)}{k\widehat{\sigma}^2} \\ &= \frac{(\widehat{\boldsymbol{\beta}}_s^{(-i)} - \widehat{\boldsymbol{\beta}}_s)^T ((\mathbf{B}\mathbf{X})^T (\mathbf{B}\mathbf{X})) (\widehat{\boldsymbol{\beta}}_s^{(-i)} - \widehat{\boldsymbol{\beta}}_s)}{k\widehat{\sigma}^2} \\ &= \frac{(\widehat{\boldsymbol{\beta}}_s^{(-i)} - \widehat{\boldsymbol{\beta}}_s)^T (\mathbf{X}^T \mathbf{B}^T \mathbf{B} \mathbf{X}) (\widehat{\boldsymbol{\beta}}_s^{(-i)} - \widehat{\boldsymbol{\beta}}_s)}{k\widehat{\sigma}^2} \\ &= \frac{(\widehat{\boldsymbol{\beta}}_s^{(-i)} - \widehat{\boldsymbol{\beta}}_s)^T (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X}) (\widehat{\boldsymbol{\beta}}_s^{(-i)} - \widehat{\boldsymbol{\beta}}_s)}{k\widehat{\sigma}^2}, \end{aligned}$$

where:

$$\hat{\beta}_s^{(-i)} = X_{(i)}(X_{(i)}^T \hat{V}_{(i,i)}^{-1} X_{(i)})^{-1} X_{(i)}^T \hat{V}_{(i,i)}^{-1} \hat{A}_{(i,i)} Y_{(i)}.$$

$\hat{V}_{(i,i)}$ and $\hat{A}_{(i,i)}$ denote \hat{V} and \hat{A} with the i^{th} row and the i^{th} column deleted, respectively.

The spatial Cook distance, CD_{si} , is declared large if $CD_{si} > 0.70$ [19]. In its simplified form, the Cook's distance in spatial regression is written as

$$\widehat{CD}_{si}(X^T V^{-1} X, k\hat{\sigma}^2) = \frac{1}{k} t_{si}^2 \frac{p_{si}}{q_{si}}, \quad (17)$$

where t_{si} is the spatial Studentized prediction residual (also called spatial internally Studentized residual), p_{si} is the spatial leverage, which is the i^{th} diagonal element of \mathbf{P}_s , and $q_{si} = 1 - p_{si}$. Let $r_{si} = y_i - \hat{y}_i$, then:

$$t_{si} = \frac{b_i^T a_i r_{si}}{\hat{\sigma} \sqrt{q_{si}}}, \quad (18)$$

where b_i and a_i are the i^{th} columns of matrices \mathbf{B} and \mathbf{A} , respectively. The spatial Studentized residual has a cut-off point of 2 to declare a point large [19,40].

Similarly, the spatial externally Studentized residual (ESRs), is defined as:

$$t_{si}^* = \frac{r_{si}}{\hat{\sigma}_{(i)} \sqrt{1 - p_{si}}} \\ = t_{si} \sqrt{\frac{n - k - 1}{n - k - t_{si}^2}}, \quad \hat{\sigma}_{(i)} = \hat{\sigma} \left(\frac{n - k - t_{si}}{n - k - 1} \right).$$

where $\hat{\sigma}_{(i)}$ is the residuals mean square excluding the i^{th} case. The ESRs follow a Student's t-distribution with $(n - k - 1)$ degrees of freedom. Thus, the spatial Studentized prediction residuals contain the neighbourhood information of both the dependent variable and the residual of each r_{si} , and the leverage \mathbf{P}_s contains the residual autocorrelation effect. The spatial potential, which is analogous to the potential in [10], is defined in Equation (19) as:

$$p_{osi} = \frac{p_{si}}{q_{si}} \quad (19)$$

where $q_{si} = 1 - p_{si}$. Let $q_{osi} = 1 - p_{osi}$.

We define the spatial measure of overall potential influence as

$$H_{si}^2 = \frac{k}{q_{osi}} \frac{d_i^2}{(1 - d_i^2)} + \frac{p_{osi}}{q_{osi}} \quad (20)$$

When measuring the influence of an observation in a linear regression model by using the Cook's distance [3], the observation in question is deleted, and the model is then refitted. In a similar way, usually a group of suspected influential observations is deleted in the linear regression and admitted into the model if it is proven clean (BACON [41], [42], DGRP [11]). This is because IOs in linear regression are global in nature; however, in a spatial regression model, IOs are local. Haining [20] noted that spatial outliers are local in nature; their attribute values are outliers if they are extreme relative to the set of values in their neighbourhood on the map. IOs in spatiotemporal statistics usually carry vital information in applications. Kou et al. [26] further pointed out that detecting spatial outliers can help in locating extreme meteorological events such as tornadoes and hurricanes, identify aberrant genes or tumour cells, discover highway traffic congestion points, pinpoint military targets in satellite images, determine possible locations of oil reservoirs and detect water pollution incidents. Thus, measuring the influence of multiple spatial locations requires a contiguous set of points to reveal the unusual features related to that neighbourhood.

Although methods that detect multiple outliers in spatial regression work well (see [21]), we refer to methods that group observations as clean or suspect, irrespective of their positions (with reference to spatial data), and admit them into the model as clean observations according to some conditions.

According to Hadi [10], examining each value of influence measure alone, such as P_{si} , $ISRs$, $ESRs$, CD_{si} and H_{si}^2 , might not be successful to indicate the IOs or the source of influence. Imon [43] and Mohammed [44] noted that one should consider both outliers and leverage points when identifying IOs. The easiest way to capture IOs is by using diagnostic plots. Following [43,45], we adopt their rules for the classification of observations into four categories, namely regular observations, vertical outliers, GLPs and BLPs. Once observations are classified accordingly, those observations that fall in the vertical outliers and BLPs categories are referred to as IOs. However, due to the local nature of spatial IOs, we have to make some modifications to the classification scheme. In this paper, a new diagnostic plot is proposed by plotting the $ISRs$ (or $ESRs$) on the Y-axis against the spatial potential, P_{osi} , on the X-axis. We consider the $ISRs$ and $ESRs$ because both measures contain spatial information. On the other hand, the potentials that are obtained from the transformed model in Equation (13) are considered in order to reflect spatial dependence. Hence, the proposed diagnostic plots are denoted as $ISRs - P_{osi}$ and $ESRs - P_{osi}$ plot, and they are based on the following classification schemes:

- (a) $ISRs - P_{osi}$
- i^{th} observation is declared RO if $|ISRs| < 2.0$ and $p_{osi} < l_2$.
 - i^{th} observation is GLP if $|ISRs| < 2.0$ and $p_{osi} > l_2$.
 - i^{th} observation is BLP if $|ISRs| > 2.0$ and $p_{osi} > l_2$.
 - i^{th} observation is IO if $|ISRs| \geq 2.0$ and $p_{osi} \leq l_2$.

Figures 1 and 2 show the classification of the observations as RO, GLP and IOs according to $ISRs - P_{osi}$ and $ESRs - P_{osi}$, respectively.

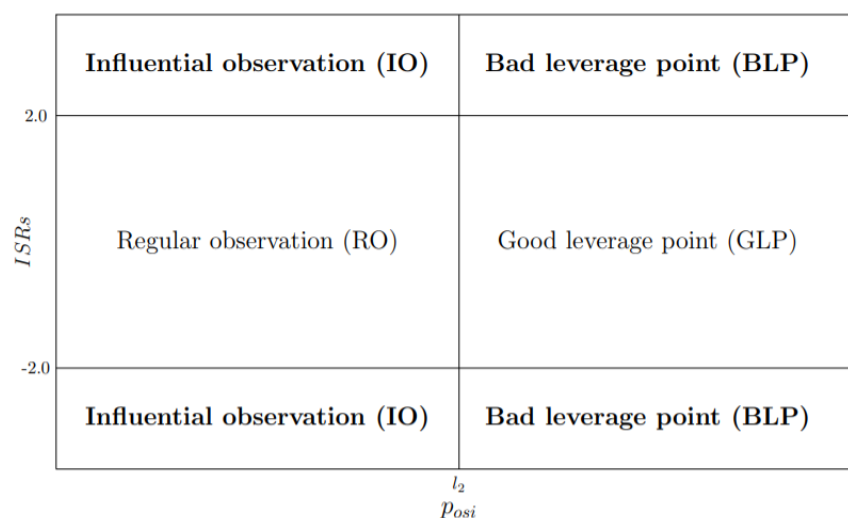


Figure 1. Classification of RO, GLP, and IO according $ISRs - P_{osi}$.

$ESRs$	$t_{(n-k-1)}$	Influential observation (IO)	Bad leverage Point (BLP)
	$ESRs$	Regular Observation (RO)	Good leverage point (GLP)
	$-t_{(n-k-1)}$	Influential observation (IO)	Bad leverage point (BLP)
		l_2	p_{osi}

Figure 2. Classification of RO, GLP, and IO according $ESRs - P_{osi}$.

(b) $ESRs - P_{osi}$

- i. i^{th} observation is declared RO if $|ESRs| < t_{n-k-1}$ and $p_{osi} < l_2$.
- ii. i^{th} observation is GL if $|ESRs| < t_{n-k-1}$ and $p_{osi} > l_2$.
- iii. i^{th} observation is IO if $|ESRs| > t_{n-k-1}$ and $p_{si} > l_2$.
- iv. i^{th} observation is IO if $|ESRs| \geq t_{n-k-1}$ and $p_{si} \leq l_2$.

4. Results and Discussions

In this section, the performance of all the proposed methods, i.e., the Cook's Distance (\widehat{CD}_{si}), H_{si}^2 (H_{si1}^2 (non-robust) and H_{si2}^2 (robust)), $ISRs - P_{osi}$ and $ESRs - P_{osi}$, is evaluated using a simulation study, artificial data and real datasets of gasoline price data in the southwest area of Sheffield, UK, COVID-19 data in the counties of the State of Georgia, USA and the life expectancy data in counties of the USA.

Simulated Data

We simulated the spatial regression model in Equation (9) for a square spatial grid with sample size, $n = 400$, $\rho = 0.4$, $\lambda = 0.5$ and $\mathbf{W}_1 = \mathbf{W}_2$, using row-standardized Queen's contiguity spatial weights. $\mathbf{x}_0 = \mathbf{1}$, $\mathbf{x}_1 \sim N(0,1)$, $\boldsymbol{\beta}_0 = \mathbf{0}$, $\boldsymbol{\beta}_1 = \mathbf{1}$ (bold face 0 and 1 refer to column vectors of values zeros and ones, respectively). The contamination is taken at two percent in each of \mathbf{X} and \mathbf{y} directions. The contamination in the \mathbf{y} direction is taken from the Cauchy distribution because of its fat tails. Contamination in the \mathbf{X} direction is taken from the following multivariate distribution,

$$\mathbf{X} \sim \left(\begin{bmatrix} 0 \\ 2 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right)$$

However, it is important to note that during the contamination, some of the contaminations may have attributes similar to those in their neighbourhood, as noted by Dowd [46], and spatial simulation is conditioned to a real dataset.

Figure 3 shows the graph of average attribute values in the neighbourhood of locations against their attribute values with added contamination. It can be observed that some of the added contamination, in black dots, are in the middle of clean data points while some stand out from the bulk of the data (i.e., away from their average neighbourhood values), which clearly indicates outlyingness.

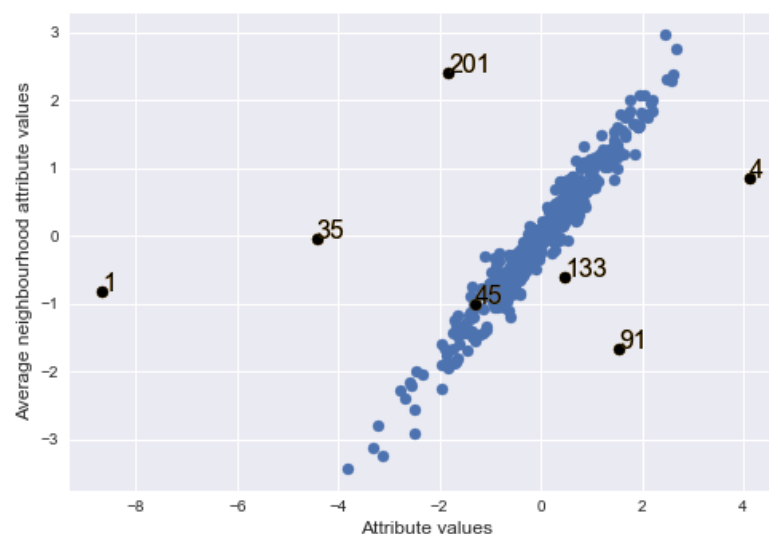


Figure 3. Graph of average attribute values in neighbourhood of locations against the attribute values in the locations with contamination (black points).

Table 1 presents the values of ISRs, ESRs and p_{osi} , where values in parentheses are their corresponding cut-off points. It shows seven locations with large Studentized residuals according to ISRs and ESRs. There are 54 observations with large potentials (>0.0078). Two out of the fifty-four potentials correspond to Studentized residuals greater than the thresholds of ISRs and ESRs (locations 51 and 201).

Table 1. ISRs, ESRs and p_{osi} of locations with large Studentized residuals in the simulated GSM model, with their cut-off points in parentheses.

Location	ISRs (2.00)	ESRs (1.97)	p_{osi} (0.0078)
1	15.0378	22.8179	0.0008
4	4.5847	4.7046	0.0033
35	−7.1434	−7.6397	0.0026
51	−4.4695	−4.5801	0.0430
91	4.7613	4.8965	0.0068
201	−6.9336	−7.3840	0.0280
265	−2.2644	−2.2762	0.0068

In order to confirm the outlyingness of the locations classified as spatial IOs, the threshold of each outlier neighbourhood given by

$$\text{med}_i + 3\text{MAD}_i$$

is computed for the Studentized residuals of the classified location and its immediate neighbourhood, where med_i is the median of the Studentized residuals and MAD_i is the median absolute deviation. The absolute value of the Studentized residuals is compared to the neighbourhood threshold for confirmation as an outlier.

The CD_{si} detected location 201, which has large ISRs, ESRs and p_{si} . ISRs $-P_{osi}$ and ESRs $-P_{osi}$ classified locations 1, 4, 35, 51, 91, 201 and 265 as IOs. As noted on Figure 4, ISRs $-P_{osi}$ and ESRs $-P_{osi}$ classified locations 1, 4, 35, 91 and 265 as outliers in the y direction, and locations 51 and 201 in both X and y directions. The cut-off limits of ESRs $-P_{osi}$ are narrower than 2 for the 5% cut-off point of the Student's t-distribution, which is around 1.96 for large sample sizes.

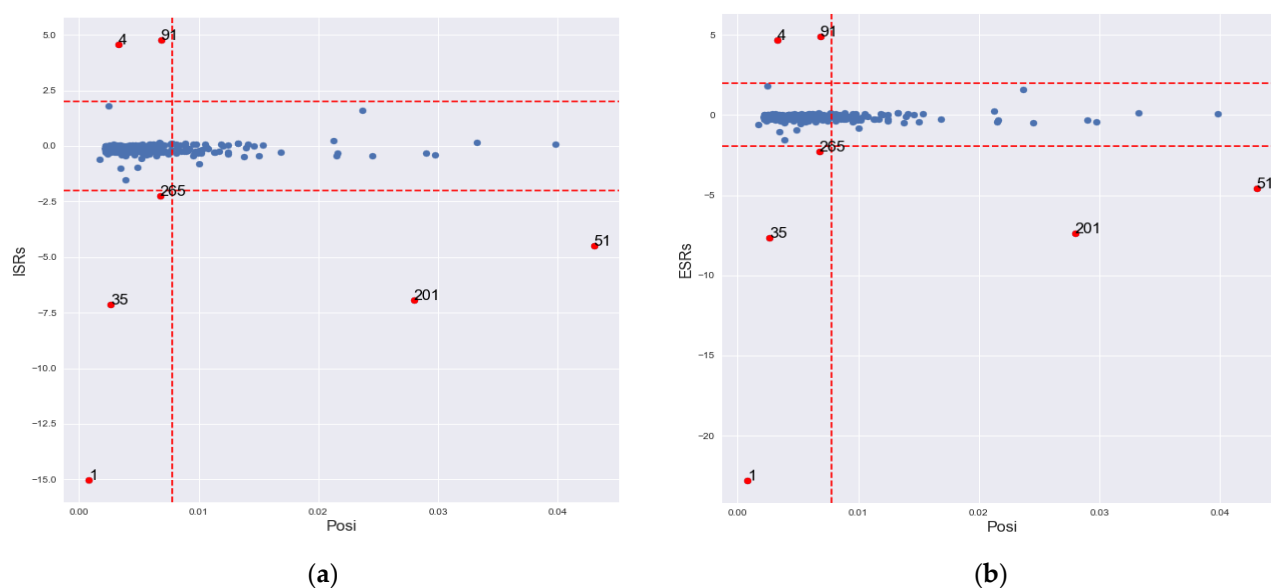


Figure 4. Graph of IO classification according to GLP, BLP and vertical outlier in $ISRs - P_{osi}$ and $ESRs - P_{osi}$ for simulated data. (a) $ISRs - P_{osi}$; (b) $ESRs - P_{osi}$.

H_{si1}^2 classified location 1 only as IO. Location 1 has large ISRs and ESRs with small p_{osi} . It is an outlier in the y direction. H_{si2}^2 identified 60 locations as IOs, including all the locations classified by the other methods. However, a diagnostic examination of the 53 other locations classified by H_{si2}^2 alone reveals that all locations that have small ISRs and ESRs with large potential values are classified as IOs. Moreover, the locations with small Studentized residuals, which show no difference with their neighbourhood, are classified as IOs. This is a clear case of swamping, perhaps due the local nature of the spatial IOs.

In a 1000-run of the simulation described above at different error variances of 0.01, 0.1, 0.2 and 0.3 as shown in Table 2, the CD_{si} consistently maintained low classification of influential observations with consistent swamping rates of 0%. The $ISRs - P_{osi}$ demonstrated a high detection to the tune of 98% while $ESRs - P_{osi}$ had 100% accurate classification of the IOs, both with swamping rates of 0%. H_{si1}^2 had less than 40% accurate classification with zero swamping rate, while the H_{si2}^2 had up to 99% accurate IO classification, but usually with very high swamping rates.

Table 2. Influential observations classification rate based on large prediction Studentized residuals and large potentials.

σ^2	Method	Accurate Classification (%)	Swamping (%)
0.01	CD_{si}	22.25	0.0
	$ISRs - P_{osi}$	98.54	0.0
	$ESRs - P_{osi}$	100.00	0.0
	H_{si1}^2	39.45	0.0
	H_{si2}^2	99.71	81.41
0.1	CD_{si}	20.64	0.0
	$ISRs - P_{osi}$	98.36	0.0
	$ESRs - P_{osi}$	100.00	0.0
	H_{si1}^2	38.09	0.0
	H_{si2}^2	99.14	76.48
0.2	CD_{si}	17.86	0.00
	$ISRs - P_{osi}$	97.51	0.00
	$ESRs - P_{osi}$	100.00	0.00
	H_{si1}^2	37.23	0.00

	H_{si2}^2	97.34	69.25
	CD_{si}	16.36	0.00
	$ISRs - P_{osi}$	96.57	0.00
0.3	$ESRs - P_{osi}$	100.00	0.00
	H_{si1}^2	36.23	0.00
	H_{si2}^2	96.00	64.42

5. Illustrative Examples

5.1. Example 1

The gasoline price data for 61 retail sites in the southwest area of Sheffield from [19] were used in Example 1. The analysis indicated the presence of spatial interaction in the error term with a Moran's I of 0.239.

The fitted SEM model is given by Equation (21):

$$\hat{\mathbf{y}}_M = 35.78 + 0.71\mathbf{X}_F + \hat{\lambda}\mathbf{W}\boldsymbol{\xi} \quad (21)$$

where, \mathbf{y}_M and \mathbf{X}_M are the March and February sales from the southwest Sheffield gasoline sale data, respectively, $\hat{\lambda} = 0.15$ is the estimate of coefficient of correlation in the residuals, \mathbf{W} is the standardized weight matrix and $\boldsymbol{\xi}$ is the vector of correlated residuals.

Table 3 shows the results of the detected IOs in the SEM model for the gasoline data with all the sites detected by the methods. A “yes” under a method column indicates that the site has been detected by the method as IO and a “no” means otherwise. The values in bold in columns ISRs, ESRs and p_{osi} indicate large Studentized residuals and potentials greater than 0.0335, respectively. Figure 5 shows the classification of observations by $ISRs - P_{osi}$ and $ESRs - P_{osi}$.

Table 3. Sites with IOs in the analysis of the southwest Sheffield gasoline data.

S/N	Site	ISRs (2.00)	ESRs (2.00)	p_{osi} (0.0335)	CD_{st}	$ISRs - P_{osi}$	$ESRs - P_{osi}$	H_{st1}^2	H_{st2}^2
1.	3	−1.8879	−1.9301	0.3538	no	No	No	no	Yes
2.	9	1.4810	1.4962	0.0223	no	No	No	no	Yes
3.	22	1.0127	1.0129	0.0779	no	No	No	no	Yes
4.	25	5.4292	7.5481	0.2773	yes	Yes	Yes	yes	Yes
5.	26	1.4438	1.4573	0.1352	no	No	No	no	Yes
6.	30	2.2054	2.2813	0.2489	no	No	No	no	Yes
7.	40	1.5692	1.5890	0.0194	no	No	No	no	Yes
8.	41	1.1974	1.2058	0.0218	no	No	No	no	Yes
9.	42	−1.9150	−1.9598	0.0378	no	No	No	no	Yes
10.	46	0.1003	0.0995	0.1319	no	No	No	no	Yes
11.	55	−1.2042	−1.2089	0.0219	no	No	No	no	Yes
12.	61	−1.8011	−1.8363	0.0319	no	No	No	no	Yes

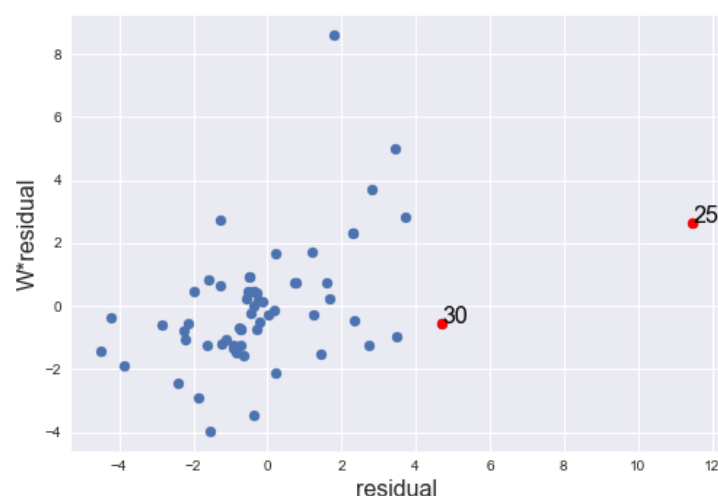


Figure 5. Graph of the lagged residuals against the residuals, of the 61 sites of the southwest Sheffield fitted with SEM, showing the IO points in red dots.

The CD_{si} , $ISRs - P_{osi}$, $ESRs - P_{osi}$, and H_{si1}^2 coincidentally identified site 25 only as IO. H_{si2}^2 detects 11 more sites as IOs in addition to site 25. Haining [19] has made elaborate diagnostic analysis of the data where he emphasized the effect of site 25 as IO in the data. Our methods have classified site 30 in addition to location 25 as IO. Figure 5 shows the graph of the lagged residuals against the residuals. It is noticeable from the graph that site 30 has also been marked as an IO.

Though the H_{si2}^2 has detected all the suspected IOs, it is prone to swamping. The remaining high potentials are classified as GLP by $ISRs - P_{osi}$ and $ESRs - P_{osi}$ since their Studentized values are small.

Figure 6 shows the graph of classification of the $ISRs - P_{osi}$ (a) and $ESRs - P_{osi}$ (b) indicating the outliers in red dots, where both are classified as outliers in both the X and y directions.

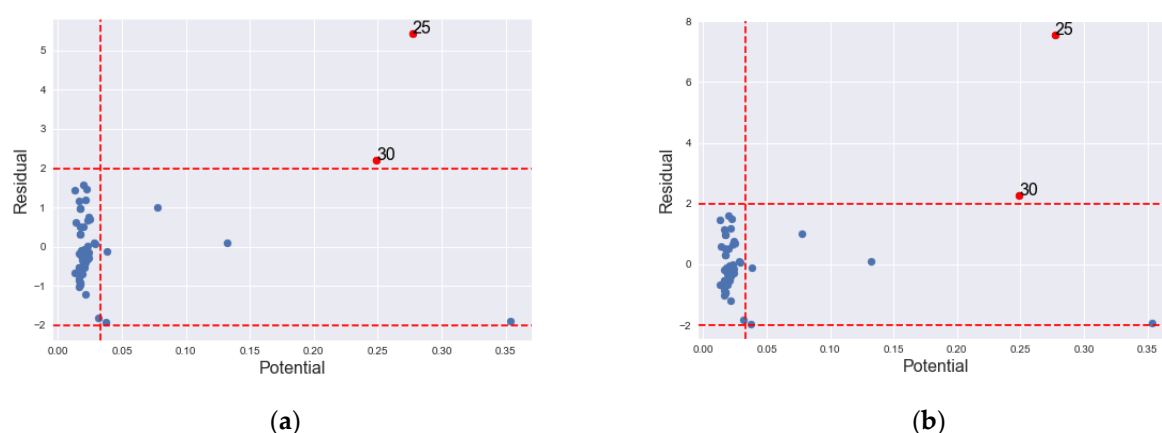


Figure 6. Graph of IO classification according to GLP, BLP and vertical outlier in South West gasoline data. (a) $ISRs - P_{osi}$; (b) $ESRs - P_{osi}$.

5.2. Example 2

The data for example 2 were the COVID-19 data for the 159 counties of the State of Georgia, USA, as of 30 June 2020 (<http://dph.georgia.gov/covid-19-daily-status-report>; accessed on 30 June 2020) and the health ranking (<http://www.countyheathrankings.org>; accessed on 30 June 2020). The case-rate per 100,000 of COVID-19 was the dependent variable. The independent variables were the population of black race in the county (X_1), population of Asians (X_2), population of Hispanic (X_3), population of people that are 65 years and above (X_4), population of female in the county (X_5) and life expectancy (X_6).

The model was fitted with the SAR model (model with the lowest Akaike information criteria (AIC) value of 2192). The SAR model is presented in Equation (22):

$$\hat{y} = \hat{\rho}W\hat{y} + \hat{\beta}_0 + \sum_{i=1}^6 \hat{\beta}_i X_i \quad (22)$$

where $\hat{\rho} = 0.6967$, $\hat{\beta}_0 = 1087.7388$, $\hat{\beta}_1 = 9.7831$, $\hat{\beta}_2 = -6.2210$, $\hat{\beta}_3 = -54.1402$, $\hat{\beta}_4 = -28.5874$, $\hat{\beta}_5 = 4.8288$ and $\hat{\beta}_6 = 40.3323$. X_1 , X_3 and $\hat{\rho}$ are significant at 5%, while X_2 and X_5 are significant at 10%. X_4 and X_6 are not significant.

The Cook's distance only classified county 50 as an IO. The $ISRs - P_{osi}$ and $ESRs - P_{osi}$ coincided in detecting counties 3, 26, 49, 50, 70, 120, 135, 141 and 142 as IOs. The H_{si1}^2 (non-robust) detected 26 and 50 as IOs. The H_{si2}^2 (robust) detected 3, 26, 50, 58, 67, 70, 98, 118, 120, 128, 131, 134, 135, 139, 141, 142, 153 and 155 counties. Table 4 shows the detected locations by the various methods with large ISRs, ESRs and high potentials in bold font.

Table 4. Detected IOs counties by different methods in the Georgia COVID-19 data.

County	ISRs (2.00)	ESRs (1.98)	p_{si} (0.0851)	CD_{si}	$ISRs - P_{osi}$	$ESRs - P_{osi}$	H_{si1}^2	H_{si2}^2
3	2.2245	2.2539	0.0257	no	no	no	no	Yes
26	4.5733	4.9060	0.1956	no	yes	yes	yes	Yes
49	2.7685	2.8313	0.0265	no	yes	Yes	no	Yes
50	5.7504	6.4737	0.2298	yes	yes	Yes	yes	Yes
58	0.7090	0.7079	0.6893	no	no	No	yes	Yes
67	0.1018	0.1015	0.3524	no	no	No	no	Yes
70	3.1334	3.2285	0.0895	no	yes	Yes	no	Yes
98	-1.8549	-1.8699	0.0105	no	no	No	no	Yes
118	-1.5657	-1.5731	0.0827	no	no	No	no	Yes
120	3.0168	3.1006	0.0544	no	yes	Yes	no	Yes
128	-2.0152	-2.0359	0.4557	no	yes	yes	no	Yes
131	-1.6718	-1.6862	0.0565	no	no	No	no	Yes
134	-1.6168	-1.6253	0.0818	no	no	No	no	Yes
135	2.1674	2.1942	0.0338	no	yes	Yes	no	Yes
141	2.6726	2.7283	0.0163	no	Yes	yes	no	Yes
142	2.1805	2.2079	0.0174	Yes	yes	no	no	Yes
153	-1.2693	-1.2718	0.2234	No	No	no	no	Yes
155	-1.2334	-1.2359	0.2472	No	No	no	no	Yes

The IOs identified by $ISRs - P_{osi}$ and $ESRs - P_{osi}$ both have large Studentized residuals and large potentials as can be observed in Table 4. Figure 7 shows the outliers in X, y and both X and y directions. The CD_{si} detected the largest Studentized residual with a high potential as IO. The H_{si1}^2 identified two observations with large Studentized values and high potential values. The H_{si2}^2 detected all suspected IOs, but with many having both small values of Studentized residuals and potential values.

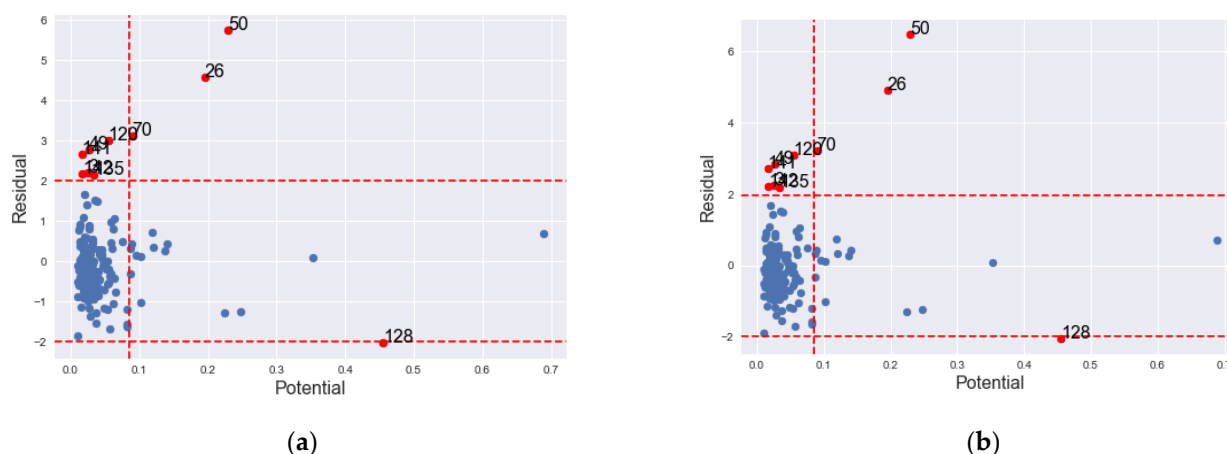


Figure 7. Graph of IO classification according to GLP, BLP and vertical outlier in the State of Georgia, USA COVID-19 data. (a) $ISRs - P_{osi}$; (b) $ESRs - P_{osi}$.

While examining the outlyingness of the classified counties, we find that county 50 is clearly an IO since it has both large Studentized residual and a large potential value. It is outside the threshold value of its neighbourhood.

Four of the counties classified by $ISRs - P_{osi}$ and $ESRs - P_{osi}$ (i.e., 26, 50, 70 and 128) are classified as vertical outliers while the counties 3, 49, 120, 135, 141 and 142 have large potential values and Studentized values greater than their threshold values and are classified as BLPs and hence IOs.

Besides the counties classified by $ISRs - P_{osi}$ and $ESRs - P_{osi}$, all the other counties detected by H_{si2}^2 have their Studentized difference residuals below their neighbourhood threshold. Though their potential values are mostly large, their prediction Studentized residuals are small in both ISRs and ESRs.

5.3. Example 3

In example 3, the life expectancy of the counties of the US was measured by population density (X_1), fair/poor health status (X_2), obesity (X_3), population in rural area (X_4), inactivity rate (X_5), population of smokers (X_6), population of black people (X_7), population of Asians (X_8) and population of Hawaiians (X_9). The data were obtained from the Kaggle website (<https://www.kaggle.com/johnjdavisiv/us-counties-covid19-weather-sociohealth-data>; accessed on 13 December 2020).

The spatial error model (SEM) had the lowest AIC value and was fitted to the data. The model was significant at the 5% level with a significant Moran's I of 0.2160. X_1 and X_4 were not significant at the 5%. All the other estimates were significant at the 5% level. The fitted model is given by:

$$\hat{y} = \hat{\beta}_0 + \sum_{i=1}^9 \hat{\beta}_i X_i + \hat{\lambda} W \xi$$

where $\hat{\lambda} = 0.4343$, $\hat{\beta}_0 = 88.4885$, $\hat{\beta}_1 = 0.0000$, $\hat{\beta}_2 = -0.0954$, $\hat{\beta}_3 = -0.0377$, $\hat{\beta}_4 = 0.0040$, $\hat{\beta}_5 = -0.0630$, $\hat{\beta}_6 = -0.3892$, $\hat{\beta}_7 = -0.0113$, $\hat{\beta}_8 = 0.1437$, $\hat{\beta}_9 = -0.2016$. Counties with fair/poor health facility had a 0.1 lower life expectancy for an increase in the population. Counties with a larger number of obese people had a decrease in life expectancy of 0.03. Similarly, those counties with a large number of people with inactivity had a life expectancy decreased by 0.06, and counties with a larger number of smokers had a life expectancy decreased by 0.04 per increase in the population. Counties with a higher number of black people and Hawaiians had a life expectancy decreased by 0.01 and 0.2, respectively, while those with a higher number of Asians had an increased rate of 0.14 in population.

The $ISRs - P_{osi}$ classified 139 counties as IOs, while $ESRs - P_{osi}$ classified eight more counties, making a total of 147. H_{si1}^2 and H_{si2}^2 have classified 24 and 324 counties as IOs, respectively. CD_{si} classified no county as IO.

6. Conclusions

In this article, we demonstrated the application of influential observations (IOs) detection techniques from the classical regression to the spatial regression model. Measures that contained spatial information in the spatial autoregression in the dependent variables and residuals were obtained. We also evaluated the performance of some methods employed in classical regression to their spatial counterparts. Though the methods work well in classical regression models, they are mostly prone to either masking or swamping in spatial applications. This is attributable to the local nature of spatial outliers. Hence, we proposed new $ISRs - P_{osi}$ and $ESRs - P_{osi}$ plots to classify observations into four categories: regular observations, vertical outliers, good leverage points and bad leverage points, whereby IOs are those observations which fall in the vertical and bad leverage point categories. Interestingly, the proposed $ESRs - P_{osi}$ diagnostic plot was very successful in classifying observations into the correct categories followed by the $ISRs - P_{osi}$, as demonstrated by the results obtained from a simulation study and real data examples. Thus, the newly established $ESRs - P_{osi}$ plot can be a suitable alternative to identify IOs in the spatial regression model.

Author Contributions: Conceptualization, A.M.B. and H.M.; methodology, A.M.B.; software, A.M.B.; validation, A.M.B. and H.M.; formal analysis, A.M.B. and N.H.A.R.; investigation, A.M.B. and N.H.A.R.; resources, A.M.B., H.M. and M.B.A.; data curation, A.M.B.; writing—original draft preparation, A.M.B.; writing—review and editing, A.M.B., H.M., M.B.A. and N.H.A.R.; visualization, A.M.B. and M.B.A.; supervision, H.M.; project administration, H.M.; funding acquisition, H.M. All authors have read and agreed to the published version of the manuscript.

Funding: This article was partially supported by the Fundamental Research Grant Scheme (FRGS) under the Ministry of Higher Education, Malaysia with project number FRGS/1/2019/STG06/UPM/01/1

Data Availability Statement: Data are available online. Data for Example 1 are available in page 332 of [19]. Data for Example 2 are available online, website link (<http://dph.georgia.gov/covid-19-daily-status-report>; accessed on 30 June 2020 and <http://www.countyheathrankings.org>; accessed on 30 June 2020). Data for Example 3 are available online, website link (<https://www.kaggle.com/johnjdavisiv/us-counties-covid19-weather-sociohealth-data>; accessed on 13 June 2020)

Conflicts of Interest: The authors declare no conflict of interest.

References

1. Belsley, D.A.; Kuh, E.; Welsch, R.E. *Regression Diagnostics: Identifying Influential Data and Sources of Collinearity*; John Wiley & Sons: New York, NY, USA, 1980; Volume 571.
2. Rashid, A.M.; Midi, H.; Slwabi, W.D.; Arasan, J. An Efficient Estimation and Classification Methods for High Dimensional Data Using Robust Iteratively Reweighted SIMPLS Algorithm Based on Nu-Support Vector Regression. *IEEE Access* **2021**, *9*, 45955–45967, doi:10.1109/ACCESS.2021.3066172.
3. Cook, R.D. Influential Observations in Linear Regression. *J. Am. Stat. Assoc.* **1977**, *74*, 169–174.
4. Hoaglin, D.C.; Welsch, R.E. The Hat Matrix in Regression and ANOVA. *Am. Stat.* **1978**, *32*, 17, <https://doi.org/10.2307/2683469>.
5. Andrews, D.F.; Pregibon, D. Finding the Outliers That Matter. *J. R. Stat. Soc. Ser. B (Methodol.)* **1978**, *40*, 85–93.
6. Hawkins, D.M. *Identification of Outliers*; Springer: Berlin/Heidelberg, Germany, 1980; Volume 11.
7. Huber, P. *Robust Statistics*; John Wiley and Sons: New York, NY, USA, 1981.
8. Cook, R.D.; Weisberg, S. Monographs on statistics and applied probability. In *Residuals and Influence in Regression*; Chapman and Hall: New York, NY, USA, 1982; ISBN 978-0-412-24280-9.
9. Chatterjee, S.; Hadi, A.S. *Sensitivity Analysis in Linear Regression*; John Wiley & Sons: New York, NY, USA, 1988; Volume 327.
10. Hadi, A.S. A New Measure of Overall Potential Influence in Linear Regression. *Comput. Stat. Data Anal.* **1992**, *14*, 1–27.

11. Habshah, M.; Norazan, M.R.; Rahmatullah Imon, A.H.M. The Performance of Diagnostic-Robust Generalized Potentials for the Identification of Multiple High Leverage Points in Linear Regression. *J. Appl. Stat.* **2009**, *36*, 507–520, <https://doi.org/10.1080/02664760802553463>.
12. Meloun, M.; Militký, J. *Statistical Data Analysis: A Practical Guide*; Woodhead Publishing Limited: Sawston, Cambridge, 2011;.
13. Puterman, M.L. Leverage and Influence in Autocorrelated Regression Models. *J. R. Stat. Soc. Ser. C (Appl. Stat.)* **1988**, *37*, 76–86.
14. Schall, R.; Dunne, T.T. A Unified Approach to Outliers in the General Linear Model. *Sankhyā Indian J. Stat. Ser. B* **1988**, *50*, 157–167.
15. Martin, R.J. Leverage, Influence and Residuals in Regression Models When Observations Are Correlated. *Commun. Stat.-Theory Methods* **1992**, *21*, 1183–1212.
16. Shi, L.; Chen, G. Influence Measures for General Linear Models with Correlated Errors. *Am. Stat.* **2009**, *63*, 40–42.
17. Cerioli, A.; Riani, M. Robust Transformations and Outlier Detection with Autocorrelated Data. In *From Data and Information Analysis to Knowledge Engineering*; Spiliopoulou, M., Kruse, R., Borgelt, C., Nürnberger, A., Gaul, W., Eds.; Studies in Classification, Data Analysis, and Knowledge Organization; Springer: Berlin/Heidelberg, Germany, 2006; pp. 262–269, ISBN 978-3-540-31313-7.
18. Christensen, R.; Johnson, W.; Pearson, L.M. Prediction Diagnostics for Spatial Linear Models. *Biometrika* **1992**, *79*, 583–591, <https://doi.org/10.1093/biomet/79.3.583>.
19. Haining, R. Diagnostics for Regression Modeling in Spatial Econometrics*. *J. Reg. Sci.* **1994**, *34*, 325–341, <https://doi.org/10.1111/j.1467-9787.1994.tb00870.x>.
20. Haining, R.P.; Haining, R. *Spatial Data Analysis: Theory and Practice*; Cambridge University Press: Cambridge, England, 2003.
21. Dai, X.; Jin, L.; Shi, A.; Shi, L. Outlier Detection and Accommodation in General Spatial Models. *Stat. Methods Appl.* **2016**, *25*, 453–475, <https://doi.org/10.1007/s10260-015-0348-1>.
22. Singh, A.K.; Lalitha, S. A Novel Spatial Outlier Detection Technique. *Commun. Stat. Theory Methods* **2018**, *47*, 247–257.
23. Anselin, L. Some Robust Approaches to Testing and Estimation in Spatial Econometrics. *Reg. Sci. Urban Econ.* **1990**, *20*, 141–163, [https://doi.org/10.1016/0166-0462\(90\)90001-J](https://doi.org/10.1016/0166-0462(90)90001-J).
24. Cerioli, A.; Riani, M. Robust Methods for the Analysis of Spatially Autocorrelated Data. *Stat. Methods Appl.* **2002**, *11*, 335–358, <https://doi.org/10.1007/BF02509831>.
25. Yildirim, V.; Mert Kantar, Y. Robust Estimation Approach for Spatial Error Model. *J. Stat. Comput. Simul.* **2020**, *90*, 1618–1638, <https://doi.org/10.1080/00949655.2020.1740223>.
26. Kou, Y.; Lu, C.-T. Outlier Detection, Spatial. In *Encyclopedia of GIS*; Springer: Boston, MA, USA, 2008; 1539–1546.
27. Hadi, A.S.; Imon, A.H.M.R. Identification of Multiple Outliers in Spatial Data. *Int. J. Stat. Sci.* **2018**, *16*, 87–96.
28. Hadi, A.S.; Simonoff, J.S. Procedures for the Identification of Multiple Outliers in Linear Models. *J. Am. Stat. Assoc.* **1993**, *88*, 1264–1272, <https://doi.org/10.1080/01621459.1993.10476407>.
29. Aggarwal, C.C. Spatial Outlier Detection. In *Outlier Analysis*; Springer: New York, NY, USA, 2013; pp. 313–341, ISBN 978-1-4614-6395-5.
30. Anselin, L. Local Indicators of Spatial Association—LISA. *Geogr. Anal.* **1995**, *27*, 93–115.
31. Politis, D.; Romano, J.; Wolf, M. Bootstrap Sampling Distributions. In *Subsampling*; Springer: New York, NY, USA, 1999; Available online: https://link.springer.com/chapter/10.1007/978-1-4612-1554-7_1 (accessed on 16 September 2021).
32. Heagerty, P.J.; Lumley, T. Window Subsampling of Estimating Functions with Application to Regression Models. *J. Am. Stat. Assoc.* **2000**, *95*, 197–211, <https://doi.org/10.1080/01621459.2000.10473914>.
33. Anselin, L. Exploratory Spatial Data Analysis and Geographic Information Systems. *New Tools Spat. Anal.* **1994**, *17*, 45–54.
34. Cressie, N.A.C. *Statistics for Spatial Data*, Rev. ed.; Wiley series in probability and mathematical statistics; Wiley: New York, NY, USA, 1993; ISBN 978-0-471-00255-0.
35. Anselin, L. *Spatial Econometrics: Methods and Models*; Studies in Operational Regional Science; Springer: Dordrecht, The Netherlands, 1988; Volume 4, ISBN 978-90-481-8311-1.
36. LeSage, J.P. *The Theory and Practice of Spatial Econometrics*; University of Toledo: Toledo, OH, USA, 1999; Volume 28.
37. Olver, P.J.; Shakiban, C.; Shakiban, C. *Applied Linear Algebra*; Springer: Berlin/Heidelberg, Germany, 2006; Volume 1.
38. Horn, R.A.; Johnson, C.R. *Matrix Analysis*, 2nd ed.; Cambridge University Press: Cambridge, NY, USA, 2012; ISBN 978-0-521-83940-2.
39. Liesen, J.; Mehrmann, V. *Linear Algebra*; Springer Undergraduate Mathematics Series; Springer International Publishing: Cham, Germany, 2015; ISBN 978-3-319-24344-3.
40. Shekhar, S.; Lu, C.-T.; Zhang, P. A Unified Approach to Detecting Spatial Outliers. *GeoInformatica* **2003**, *7*, 139–166.
41. Billor, N.; Hadi, A.S.; Velleman, P.F. BACON: Blocked Adaptive Computationally Efficient Outlier Nominators. *Comput. Stat. Data Anal.* **2000**, *34*, 279–298.
42. Imon, A. Identifying Multiple High Leverage Points in Linear Regression. *J. Stat. Stud.* **2002**, *3*, 207–218.
43. Rahmatullah Imon, A.H.M. Identifying Multiple Influential Observations in Linear Regression. *J. Appl. Stat.* **2005**, *32*, 929–946, <https://doi.org/10.1080/02664760500163599>.
44. Midi, H.; Mohammed, A. The Identification of Good and Bad High Leverage Points in Multiple Linear Regression Model. *Math. Methods Syst. Sci. Eng.* **2015**, *147*–158.
45. Bagheri, A.; Midi, H. Diagnostic Plot for the Identification of High Leverage Collinearity-Influential Observations. *Sort: Stat. Oper. Res. Trans.* **2015**, *39*, 51–70.

-
46. Dowd, P. The variogram and kriging: Robust and resistant estimators. In *Geostatistics for Natural Resources Characterization*; Springer: Berlin/Heidelberg, Germany, 1984; pp. 91–106.