



Article Robust Variable Selection with Exponential Squared Loss for the Spatial Durbin Model

Zhongyang Liu, Yunquan Song * and Yi Cheng

College of Science, China University of Petroleum, Qingdao 266580, China * Correspondence: statistics99@163.com; Tel.: +86-0532-8698-2921

Abstract: With the continuous application of spatial dependent data in various fields, spatial econometric models have attracted more and more attention. In this paper, a robust variable selection method based on exponential squared loss and adaptive lasso is proposed for the spatial Durbin model. Under mild conditions, we establish the asymptotic and "Oracle" properties of the proposed estimator. However, in model solving, nonconvex and nondifferentiable programming problems bring challenges to solving algorithms. To solve this problem effectively, we design a BCD algorithm and give a DC decomposition of the exponential squared loss. Numerical simulation results show that the method is more robust and accurate than existing variable selection methods when noise is present. In addition, we also apply the model to the 1978 housing price dataset in the Baltimore area.

Keywords: spatial Durbin model; exponential squared loss; robust variable selection

1. Introduction

In recent years, spatial section data have been widely used in geography, politics, environment, and other fields. Therefore, spatial econometrics, a model initially used in the economic area, has also attracted much attention. Anselin (1988) [1] divides spatial econometric models into spatial error models, spatial hysteresis models, and spatial Durbin models (SDM). Among them, the spatial Durbin model is represented as $y = \rho W y + X \beta +$ $WX\delta + \varepsilon$. The spatial Dubin model considers the influence of the independent variable and the dependent variable of the spatial lag on the dependent variable simultaneously and can more easily estimate the unbiased coefficient. At the same time, the spatial Dubin model can also calculate spatial spillover effects based on panel data. In spatial regression analysis, the influence of regional locations on observations is expressed employing spatial weight matrix W, and the appropriate setting of spatial weight matrix is an essential basis for spatial econometric analysis. There are two main ways to select a spatial weight matrix: The first method is to select a spatial weight matrix from an optional set of spatial weight matrices. Kelejian (2008) [2] uses GMM estimation to select an actual spatial weight matrix. A non-nested J test method is proposed to test a set of alternative models with different spatial weight matrices for the empty SAR model. The second type of method estimates the weight matrix by averaging different spatial weight matrices. Zhang and Yu (2018) [3] propose a model averaging process to reduce estimation error. This approach overcomes the difficulty that the actual spatial weight matrix is not in the candidate matrix.

In the field of classical linear regression, much work has contributed to variable selection. Among them, the most popular method is to add penalty functions to the model for variable selection. These punishment regression methods have a unified theoretical framework, such as most minor absolute shrinkage and selection operators (lasso, Tibshirani, 1996) [4], smoothly clipped absolute deviation (SCAD, Fan and Li, 2001) [5], and adaptive lasso (Zou, 2006) [6]. Since SDM has spatial autocorrelation, the above variable selection method can be directly applied to the SDM model.

Due to noise and outliers, the classical variable selection methods in regression models often face the problem of instability, so many scholars have proposed some robust variable



Citation: Liu, Z.; Song, Y.; Cheng, Y. Robust Variable Selection with Exponential Squared Loss for the Spatial Durbin Model. *Entropy* **2023**, 25, 249. https://doi.org/10.3390/ e25020249

Academic Editor: Yuehua Wu

Received: 29 November 2022 Revised: 25 January 2023 Accepted: 28 January 2023 Published: 30 January 2023



Copyright: © 2023 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/). selection algorithms. The Huber loss function was widely used in early studies, but this function has some limitations in efficiency and solution. Wang et al. (2013) [7] proposed a robust parameter estimation method based on the exponential squared loss function, which is widely used in boosting algorithms (Friedman et al., 2000) [8]. When γ is small, the loss of experience caused by a larger |t| value is close to 1; therefore, the loss function is robust to parameter estimation. Wang et al. (2013) [7] also point out that this method is more robust than other robust variable selection methodsm such as Huber estimation, quantile regression estimation (Koenker and Bassett, 1978) [9], and compound quantile regression estimation (Zou and Yuan, 2008) [10], and proposed the selection method of parameter γ .

Our research focuses on variable selection for spatial Durbin models. The spatial Durbin model combines the spatial interaction of dependent and explanatory variables, but only a few researchers use and study this model. Beer and Rield (2011) [11] used the maximum likelihood estimation to estimate the parameters of the spatial Durbin model. They used the Monte Carlo method to analyze the characteristics of the estimator. Mustaqim (2018) [12] discusses instrumental variable efficiency in simultaneous spatial Durbin models. Estimation methods are 2SLS and GMM-S2SLS. The analysis results show that the GMM-S2SLS method produces less bias than the 2SLS method. Zhu, Yanli (2020) [13] proposed parameter estimation of the spatial Durbin model based on Markov Chain Monte Carlo (MCMC). Wei, Lili (2021) [14] proposed a within-group spatial two-stage least squares estimator. However, the existing variable selection methods are affected by outliers in limited samples and are not robust enough. Therefore, it is imperative to study a robust variable selection method.

Considering robustness, we combine parameter penalty with exponential square loss and assume that the errors of the model are independent and identically distributed. For the parameter penalty method, we use adaptive lasso. We applied the robust selection method based on the exponential squared loss variable to the spatial autoregressive model and achieved satisfactory results [15]. The spatial autoregressive model is one of the special cases of the spatial Durbin model. In this paper, we aim to study the application of the robust selection method based on the exponential squared loss variable in the spatial Durbin model.

A robust variable selection method for the spatial Durbin model based on adaptive lasso penalty and exponential square loss function is proposed in this paper. This method cannot only estimate the regression coefficient but also has the function of variable selection. Next, we show the framework of the paper.

- 1. We build a robust variable selection method for SDM, equipped with an exponential squared loss, resistant to the influence of outliers in the observed values and errors estimating the space weight matrix.
- 2. To solve the optimization problem of SDM, we propose a block coordinate descent (BCD) algorithm. Secondly, to solve the subproblems generated by the BCD algorithm, we design the DC decomposition of exponential square loss and construct the CCCP program. Finally, to obtain the BCD algorithm's convergence, we analyze the algorithm's convergence rate to the stagnation point under mild conditions.
- 3. We proved the "Oracle" property of the robust variable selection method and conducted numerical experiments to verify the robustness and effectiveness of the model. Numerical studies show that when there are outliers in the observed data, the method proposed in this paper is superior to the comparison method in correctly identifying zero coefficients, nonzero coefficients, and MedSE incorrectly.

The structure of this paper is as follows. Section 2 introduces the spatial Durbin model and gives the exponential square loss function based on adaptive lasso. In Section 3, we propose an effective algorithm to complete the variable selection process. In order to check the performance of the model under limited samples, we have carried out a numerical simulation in Section 4. In Section 5, we apply our model to real-world datasets. We summarize the full text in Section 6.

2. Variable Selection and Estimation

2.1. Spatial Durbin Model

The observed dependent variable $Y_i \in R^{1 \times 1}$, and the corresponding independent vector $X_i = (X_{i1}, \ldots, X_{ip})$, where the *p* is a fixed constant. Let the dependent variable vector $Y = (Y_1, \ldots, Y_n)^T$ and the independent variable matrix $X = (X_1, \ldots, X_n)^T \in R^{n \times p}$. The spatial Durbin model is as follows:

$$Y = \rho W Y + X \beta + W X \delta + \varepsilon. \tag{1}$$

where the regression coefficient vector $\beta = (\beta_1, \ldots, \beta_p)^T \in \mathbb{R}^{p \times 1}$, the spatial autocorrelation coefficient $\rho \in \mathbb{R}^{1 \times 1}$, the regression coefficient vector of exogenous variables $\delta = (\delta_1, \ldots, \delta_p)^T \in \mathbb{R}^{p \times 1}$, and the error vector $\varepsilon = (\varepsilon_1, \ldots, \varepsilon_n)^T \in \mathbb{R}^{n \times 1}$. WX is a spatial lag term that reflects the interaction of independent variables between individuals. *Wy* embodies the interaction between the strain variable *y* and its surrounding *y*. We assume that noises ε all obey $N(0, \sigma^2)$ and are independent of each other. y can be expressed as the following formula:

$$Y = (I_n - \rho W)^{-1} (X\beta + WX\delta + \varepsilon).$$
⁽²⁾

Since the maximum eigenvalue of *W* is 1 after normalization, to guarantee $(I_n - \rho W)$ reversibility, we order $|\rho| < 1$. Additionally, in this article, we ignore the endogenous nature of the model.

2.2. Variable Selection Method for SDM

Rewrite model (1) as model (3) in the following form:

$$\varepsilon_i(\theta) = Y_i - (\rho W Y_i + X_i \beta + W X_i \delta).$$
(3)

Take the variable selection for the SDM into consideration. In practical applications, the regression coefficient vector β^* is usually sparse. At the same time, sparse solutions can find useful dimensions and reduce redundancy, as well as improve the accuracy and robustness of regression prediction (Fan and Li, 2001 [5]; Tibshirani, 1996 [4]). Applying the penalized method to variable selection is natural, which can select essential variables and estimate the regression coefficient. In this article, we punish the loss function using the adaptive lasso penalty function. The adaptive lasso penalty is described as follows:

$$\sum_{j=1}^{p} P(|\beta_{j}|) = \sum_{j=1}^{p} \eta_{j} |\beta_{j}|.$$
(4)

where $\eta_j = \frac{1}{|\hat{\beta}_j|}$, $\hat{\beta}_j$ is generally given by least squares estimates. Considering that the exponential square loss function has good robustness, we use it as the model's loss function in this paper. The exponential square loss expression is as follows:

$$\phi_{\gamma}(t) = 1 - \exp\left(-t^2/\gamma\right). \tag{5}$$

Here, γ is a parameter that controls the robustness of the loss function. γ limits the effect of outliers on the model but also reduces the accuracy of the model. Therefore, it is essential to choose the right γ . The method of selecting the right γ is shown in Section 2.4.

The model is constructed on the basis of the above model (3). The objective function to be solved is as follows:

$$\min_{\boldsymbol{\beta}\in\boldsymbol{R}^{p},\boldsymbol{\delta}\in\boldsymbol{R}^{p},\boldsymbol{\rho}\in[0,1]}L(\boldsymbol{\beta},\boldsymbol{\delta},\boldsymbol{\rho}) = \frac{1}{n}\sum_{i=1}^{n}\phi_{\gamma}(Y_{i}-\boldsymbol{\rho}Y_{i}-X_{i}\boldsymbol{\beta}-WX_{i}\boldsymbol{\delta}) + \lambda\sum_{j=1}^{p}\eta_{j}|\boldsymbol{\beta}| + \lambda\sum_{j=1}^{p}\sigma_{j}|\boldsymbol{\delta}|.$$
 (6)

We may as well order

$$\tilde{Y} = WY$$
,

$$ilde{X} = [X, WX],$$

 $ilde{eta} = \left[eta^T, \delta^T
ight]^T.$

We can obtain a simplified expression of (6) as follows as (7):

$$\min_{\widetilde{\beta}\in R^{2p}, \rho\in[0,1]} L(\widetilde{\beta},\rho) = \sum_{i=1}^{n} \phi_{\gamma} \Big(Y_{i} - \rho \widetilde{Y}_{i} - \widetilde{X}_{i} \widetilde{\beta} \Big) + \lambda \sum_{j=1}^{2p} \eta_{j} \Big| \widetilde{\beta}_{j} \Big|,$$
(7)

where $\lambda > 0$ is a regularization parameter. $\phi_{\gamma}(.)$ is exponential squared loss.

2.3. Oracle Properties and Large Sample Properties

In this section, we discuss the large sample properties and oracle properties of the proposed spatial Durbin model parameter estimation method.

First of all, let us make the true value of $\tilde{\boldsymbol{\beta}}$ be $\tilde{\boldsymbol{\beta}}_0 = (\tilde{\beta}_{10}, \dots, \tilde{\beta}_{2p0})^T$. Additionally, because $\tilde{\boldsymbol{\beta}}_0 = [\boldsymbol{\beta}_0^T, \boldsymbol{\delta}_0^T]^T$, where $\boldsymbol{\beta}_0 = (\boldsymbol{\beta}_{10}^T, \boldsymbol{\beta}_{20}^T)^T$, $\boldsymbol{\delta}_0 = (\boldsymbol{\delta}_{10}^T, \boldsymbol{\delta}_{20}^T)^T$. Based on the sparsity assumed above, we assume that $\boldsymbol{\beta}_{20} = \mathbf{0}, \boldsymbol{\delta}_{20} = \mathbf{0}$. So, $\tilde{\boldsymbol{\beta}}_0 = [\boldsymbol{\beta}_{10}^T, \mathbf{0}^T, \mathbf{0}^T, \mathbf{0}^T]^T$. For the convenience of expression, we make a transformation to the $\tilde{\boldsymbol{\beta}}_0$, so that $\tilde{\boldsymbol{\beta}}_0 = [\boldsymbol{\beta}_{10}^T, \boldsymbol{\delta}_{10}^T, \mathbf{0}^T, \mathbf{0}^T]^T = [\tilde{\boldsymbol{\beta}}'_{10}^T, \mathbf{0}^T]^T$. In order to adapt to this transformation in $\tilde{\boldsymbol{\beta}}_0, \tilde{X}$ needs to make a similar transformation. In the following, we all assume that \tilde{X} was transformed accordingly. For convenience, we express $\tilde{\boldsymbol{\beta}}'_{10}$ as $\tilde{\boldsymbol{\beta}}_{10}$ in the following text. Let $\hat{\boldsymbol{\beta}} = (\hat{\boldsymbol{\beta}}_1^T, \hat{\boldsymbol{\beta}}_2^T)^T$ be the resulting estimator of (4), suppose that the $\hat{\boldsymbol{\beta}}$ here has also undergone the above transformation. $I(\tilde{\boldsymbol{\beta}}, \gamma) = \frac{2}{\gamma} \int ZZ^T e^{-r^2/\gamma} (\frac{2r^2}{\gamma} - 1) dF(Z, y)$, where $r = Y - (I_n - \rho W)^{-1} \tilde{X} \tilde{\boldsymbol{\beta}} = Y - Z \tilde{\boldsymbol{\beta}}, Z = (I_n - \rho W)^{-1} \tilde{X}, a_n = \max \{ p'_{\lambda_{nj}}(|\tilde{\boldsymbol{\beta}}_{0j}|) : \tilde{\boldsymbol{\beta}}_{0j} \neq 0 \}$. Let the true value of ρ be ρ_0 . Thus, $\boldsymbol{\theta}_0 = (\rho_0, \tilde{\boldsymbol{\beta}}_0)$. For ease of presentation, let $\tilde{\boldsymbol{\beta}}_{10} = \rho$ and $\tilde{\boldsymbol{\beta}}_{1j} = \tilde{\boldsymbol{\beta}}_{1j}, j = 1, 2, \dots, s$, then denote $\tilde{\boldsymbol{\beta}}_1 = (\rho, \tilde{\boldsymbol{\beta}}_{11}, \dots, \tilde{\boldsymbol{\beta}}_{1s})^T$ and $\tilde{\boldsymbol{\beta}}_{01} = (\rho_0, \tilde{\boldsymbol{\beta}}_{01}, \dots, \tilde{\boldsymbol{\beta}}_{0s})^T$.

We prove the asymptotic and oracle properties of the proposed penalty estimators. Before we can prove it, we need the following hypothesis.

Assumption 1. $\Sigma = E(ZZ^T)$ is positive definite and $E||Z||^3 < \infty$.

Assumption 2. The matrix $I_n - \rho W$ is nonsingular with $|\rho| < 1$.

Assumption 3. *The row and column sums of the matrices* W_n *and* $I - \rho W_n$ *are bounded uniformly in absolute value.*

Assumption 4. For matrix $G_n = W(I - \rho W)^{-1}$, there exists a constant $\tilde{\lambda}_c$ such that $\tilde{\lambda}_c I_n - G_n G_n^T$ is positive semidefinite for all n.

Assumption 5. $1/\min_{s+1 \le j \le p} \lambda_j = o_p$ (1). Additionally, with probability 1, $\liminf_{n \to \infty} \liminf_{t \to 0^+} \left\{ \min_{s+1 \le j \le p} \frac{p'_{\lambda_j(t)}}{\lambda_j} \right\} > 0.$

Assumption 6. $\sqrt{n}a_n = o_p(1), b_n = o_P(1).$

Assumption 7. $(\gamma_n - \gamma_0) = o_p(1)$ for some $\gamma_0 > 0$.

Assumption 8. There are constants C_1 and C_2 such that, when $\theta_1, \theta_2 > C_1\lambda_j \left| p_{\lambda_j}''(\theta_1) - p_{\lambda_j}''(\theta_2) \right| \le C_2 |\theta_1 - \theta_2|$, for j = 0, 1, ..., p.

For our proposed estimator, we give the following sample properties. The following theorem gives the consistency and "oracle" property of the proposed estimator.

Theorem 1. If Assumptions 1 - 8 are true, then there is a local maximizer $\hat{\boldsymbol{\theta}}$ such that $\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0\| = O_p (n^{-1/2} + a_n)$.

Theorem 2. (Oracle Property). Suppose that Assumptions 1 - 8 hold, and $I(\tilde{\boldsymbol{\beta}}_0, \gamma_0)$ is negative definite. If $\gamma_n - \gamma_0 = o_p(1)$ for some $\gamma_0 > 0$, $\hat{\theta} = (\hat{\rho}, \hat{\beta}_1^T, \hat{\beta}_2^T)^T$ must satisfy:

- (1) sparsity, that is, $\hat{\beta}_{n2} = \mathbf{0}$ with probability 1;
- (2) *asymptotic normality:*

$$\begin{split} \sqrt{n} \big(I_1(\tilde{\boldsymbol{\beta}}_{01}, \gamma_0) + \Sigma_1 \big) \Big\{ \Big(\hat{\boldsymbol{\beta}}_{n1} - \tilde{\boldsymbol{\beta}}_{01} \Big) + \big(I_1(\tilde{\boldsymbol{\beta}}_{01}, \gamma_0) + \Sigma_1 \big)^{-1} \Delta \Big\} &\to N(\boldsymbol{0}, \Sigma_2) \\ where \ \hat{\boldsymbol{\beta}}_{n1} &= \Big(\hat{\rho}, \hat{\boldsymbol{\beta}}_{11}, \dots, \hat{\boldsymbol{\beta}}_{1s} \Big)^T, \text{ and } \tilde{\boldsymbol{\beta}}_{01} &= \big(\rho_0, \tilde{\boldsymbol{\beta}}_{01}, \dots, \tilde{\boldsymbol{\beta}}_{0s} \big)^T, \\ \Sigma_1 &= \operatorname{diag} \Big\{ p_{\lambda_1}''(|\tilde{\boldsymbol{\beta}}_{01}|), \dots, p_{\lambda_s}''(|\tilde{\boldsymbol{\beta}}_{0s}|) \Big\} \\ \Sigma_2 &= \operatorname{cov} \Big(\exp \Big(-r^2/\gamma_0 \Big) \frac{2r}{\gamma_0} Z_{i1} \Big) \\ \Delta &= \Big(p_{\lambda_1}'(|\tilde{\boldsymbol{\beta}}_{01}|) \operatorname{sign}(\tilde{\boldsymbol{\beta}}_{01}), \dots, p_{\lambda_s}'(|\tilde{\boldsymbol{\beta}}_{0s}|) \times \operatorname{sign}(\tilde{\boldsymbol{\beta}}_{0s}) \Big)^T \\ I_1(\tilde{\boldsymbol{\beta}}_{01}, \gamma_0) \\ &= \frac{2}{\gamma_0} E \Big[\exp \Big(-r^2/\gamma_0 \Big) \Big(\frac{2r^2}{r_0} - 1 \Big) \Big] \times \Big(E Z_{i1} Z_{i1}^T \Big). \end{split}$$

The detailed proofs of Theorem 1 and Theorem 2 are shown in the Appendixes A and B.

2.4. The Selection of Parameter γ

Parameter γ can control the robustness and efficiency of the robust variable selection method. Wang et al. (2013) [7] proposed a parameter selection method based on normal regression. In this paper, we extend the selection method of parameter γ to the spatial Durbin model. The specific process is as follows:

Step 1. Initialize $\hat{\rho} = \rho^{(0)}$ and $\hat{\beta} = \tilde{\beta}^{(0)}$. Set $\rho^{(0)} = \frac{1}{2}$, $\tilde{\beta}^{(0)}$ a robust estimator. Rewrite the model $Y = \rho WY + X\tilde{\beta} + WX\delta + \epsilon$ as $Y^* = X^*\tilde{\beta}^* + \epsilon$, where $Y^* = Y - \rho WY$, $X^* = [X WX]$, $\tilde{\beta}^* = [\tilde{\beta}, \delta]^T$.

Step 2. Find the pseudo-outlier set of the sample: Let $D_n = \{ (X_1^*, Y_1^*), \dots, (X_n^*, Y_n^*) \}$. Calculate $r_i(\hat{\beta}) = Y_i^* - X_i^* \hat{\beta}, i = 1, \dots, n$ and $S_n = 1.4826 \times \text{median}_i | r_i(\hat{\beta}^*) - \text{median}_j(r_j(\hat{\beta}^*)) |$. Then, there exist the pseudo-outlier set $D_m = \{ (X_i, Y_i) : |r_i(\hat{\beta}^*)| \ge 2.5S_n \}$, set $m = \sharp \{ 1 \le i \le n : |r_i(\hat{\beta}^*)| \ge 2.5S_n \}$, and $D_{n-m} = D_n \setminus D_m$.

Step 3. Select the tuning parameter γ_n : construct $\hat{V}(\gamma) = {\{\hat{I}(\hat{\beta}^*)\}}^{-1} \tilde{\Sigma}_2 {\{\hat{I}(\hat{\beta}^*)\}}^{-1}$, in which

$$\hat{I}(\hat{\beta}^*) = \frac{2}{\gamma} \left\{ \frac{1}{n} \sum_{i=1}^n \exp\left(-r_i^2(\hat{\beta}^*)/\gamma\right) \left(\frac{2r_i(\hat{\beta}^*)}{\gamma} - 1\right) \right\} \cdot \left(\frac{1}{n} \sum_{i=1}^n X_i X_i^T\right),$$

$$\tilde{\Sigma}_2 = \operatorname{Cov}\left\{ \exp\left(-r_1^2(\hat{\beta}^*)/\gamma\right) \frac{2r_1(\hat{\beta}^*)}{\gamma} X_1, \dots, \exp\left(-r_n^2(\hat{\beta}^*)/\gamma\right) \frac{2r_n(\hat{\beta}^*)}{\gamma} X_n \right\}.$$

Next, let γ_n be the minimizer of det($\hat{V}(\gamma)$) in the set $G = \{\gamma : \zeta(\gamma) \in (0, 1]\}$, where $\zeta(\cdot)$ enjoys the common definition with that in Wang et al. (2013) [7].

Step 4. Update $\hat{\rho}$ and $\hat{\beta}$ as the optimal solution of $\min_{\tilde{\beta} \in R^p, \rho \in [0,1]} \frac{1}{n} \sum_{i=1}^n \phi_{\gamma} (Y_i - \rho \tilde{Y}_i - \tilde{X}_i \tilde{\beta})$, where $\tilde{Y} = WY$, $\tilde{X} = [X WX]$, $\tilde{\beta} = [\beta, \delta]^T$. Go to Step 2 until convergence.

In the above process, the initial step requires an initial value $\tilde{\beta}^{(0)}$. In practice, the estimate of LAD loss is usually used as $\tilde{\beta}^{(0)}$.

2.5. The Selection of Parameter λ and η_i

We order $\lambda_i = \lambda \cdot \eta_i$, in which λ and η_i are from model (7). Usually, researchers use cross-validation, AIC, and BIC criteria to select λ_i . In this paper, considering the complexity of computation and the consistency of variable selection, we adopt the method of Wang, Li, and Tsai (2007) [16] to consider regularization parameters by minimizing the BIC-type objective function. The BIC-type objective function is as follows:

$$\sum_{t=1}^{n} \left[1 - \exp\left(Y_i - \rho \widetilde{Y}_i - \widetilde{X}_i \widetilde{\beta}\right)^2 / \gamma \right] + n \sum_{j=1}^{2p} \lambda_i |\widetilde{\beta}_j| - \sum_{j=1}^{2p} \log(0.5n\lambda_j) \log(n),$$

The selection method of parameter γ is given above. This makes $\lambda_i = \log(n)/n|\theta_i|$. In practice, let $\theta_i = \hat{\theta}_i$, where $\hat{\theta}_i$ is the exponential square loss estimator without penalty term. Note that this choice satisfies the condition $\hat{\lambda}_i \rightarrow 0$ for $i \leq d_0$, and $\hat{\lambda}_i \rightarrow \infty$ for $i > d_0$, with d_0 as the number of nonzero value in the θ_0 . Therefore, the final estimator can ensure the consistency of variable selection.

3. Algorithm for Model Solving

In this section, we focus on designing algorithms to solve model (7). This optimization problem has two optimization variables, $\tilde{\beta} \in R^{2p}$ and $\rho \in [0, 1]$. So, the block coordinate descent algorithm becomes our first choice. However, the subproblems used to solve $\tilde{\beta}$ are nonconvex functions and are not differentiable, and the convergence of the block coordinate drop algorithm is difficult to guarantee. In this case, we used bump decomposition and CCCP algorithms to deal with it. Finally, regarding the processing of penalty terms in the optimization model, we use the ISTA algorithm. This is reflected below.

3.1. Block Coordinate Descent Algorithm Frame

We present the framework of the block coordinate descent algorithm in Algorithm A. Next, we need to solve subproblems (8) and (9).

3.2. Solving the Subproblem (8)

Subproblem (8) minimizes the univariate function at the interval [0,1], so it can be solved using a golden section algorithm based on parabolic interpolation. For more information about the algorithm, see Forsythe et al. (1977) [17]. It is not repeated in this article.

3.3. Solving the Subproblem (9)

For subproblem (9), by observation, we can see that the penalty term part of the optimization model is the convex function, and the loss function part ϕ_{γ} can also be decomposed into the difference between the two convex functions, that is, the DC function. So, subproblem (8) is DC programming. We can construct corresponding algorithms to solve the problem.

Algorithm 1: Block coordinate descent algorithm

1. Set initial value for $^{0} \in \mathbb{R}^{2p}$ and $\rho^{0} \in (0,1)$;

- 2. *repeat* { For $k = 0, 1, 2, \dots$ }
- 3. Solve the subproblem about ρ with initial point ρ^k :

ĥ

$$p^{k+1} \leftarrow \min_{\rho \in [0,1]} L\left(\tilde{\beta}^k, \rho\right)$$
 (8)

4. Solve the subproblem with initial value $\tilde{\beta}^k$,

$$\beta^{k+1} \leftarrow \min_{\tilde{\beta} \in \mathbb{R}^{2p}} L\left(\tilde{\beta}, \rho^{k+1}\right) \tag{9}$$

to obtain a solution $\tilde{\beta}^{k+1}$, ensuring that $L(\tilde{\beta}^k, \rho^{k+1}) - L(\tilde{\beta}^{k+1}, \rho^{k+1}) \leq 0$, and $\tilde{\beta}^{k+1}$ is a stationary point of $L(\tilde{\beta}, \rho^{k+1})$. 5. *until* convergence.

We can first perform a DC decomposition of the loss function $\phi_{\gamma}(t) = 1 - \exp(-t^2/\gamma)$. Suppose there are two convex functions F(t) and G(t), make $F(t) - G(t) = \phi_{\gamma}(t)$. Because $F(t) = G(t) + \phi_{\gamma}(t)$ is a convex function, $F''(t) = \phi_{\gamma}''(t) + G''(t) > 0$, $\forall t \in \mathbb{R}$. We may as well order $G''(t) = \frac{2}{\gamma}\frac{2}{\gamma}t^2$. So, we can make $G(t) = \frac{1}{3\gamma^2}t^4$, $F(t) = G(t) + \phi_{\gamma}(t) = 1 - \exp(-t^2/\gamma) + \frac{1}{3\gamma^2}t^4$. It can be verified that both F(t) and G(t) are convex functions.

The DC decomposition of $\phi_{\gamma}(t)$ is as follows:

ŀ

$$F(t) = 1 - \exp\left(-t^2/\gamma\right) + \frac{1}{3\gamma^2}t^4,$$
 (10)

$$G(t) = \frac{1}{3\gamma^2} t^4,\tag{11}$$

$$\phi_{\gamma}(t) = F(t) - G(t). \tag{12}$$

We can use the CCCP algorithm to solve the problem after DC decomposition. Next, define the following two functions:

$$J_{\text{vex}}(\tilde{\beta}) = \frac{1}{n} \sum_{i=1}^{n} F\left(Y_i - \rho^{k+1} \langle w_i, Y \rangle - X_i \tilde{\beta}\right) + \lambda \sum_{j=1}^{2p} P\left(\left|\tilde{\beta}_j\right|\right),$$
(13)

$$J_{\text{cav}}(\tilde{\beta}) = \frac{1}{n} \sum_{i=1}^{n} G\Big(Y_i - \rho^{k+1} \langle w_i, Y \rangle - X_i \tilde{\beta}\Big).$$
(14)

 w_i is the *i* th row of the weight matrix *W*, and $\sum_{j=1}^{p} P(|\tilde{\beta}_j|)$ is a convex penalty with respect to $\tilde{\beta}$. Then, $J_{\text{vex}}(\cdot)$ and $J_{\text{cav}}(\cdot)$ are a convex function and concave function, respectively. So, the suboptimization problem (9) can be rewritten as

$$\min_{\tilde{\beta} \in \mathbb{R}^{2p}} L(\tilde{\beta}, \rho^{k+1}) = J_{\text{vex}}(\tilde{\beta}) + J_{\text{cav}}(\tilde{\beta}).$$
(15)

At this point, it can be found that the optimization problem (15) can be solved by the CCCP(Concave–Convex Procedure) algorithm. The CCCP algorithm framework is shown below (Algorithm 2):

Algorithm 2: The Concave–Convex Procedure (CCCP) 1. Initialize $\tilde{\beta}^0$. Set k = 0. 2. *repeat* 3. $\tilde{\beta}^{k+1} = \operatorname{argmin}_{\tilde{\beta}} J_{\operatorname{vex}}(\tilde{\beta}) + J'_{\operatorname{cav}}(\tilde{\beta}^k) \cdot \tilde{\beta}$ (16)

4. *untill* convergence of $\tilde{\beta}^k$.

It is easy to know that the optimization problem (16) is a convex optimization problem. The CCCP algorithm minimizes the problem (15) by iteratively solving a series of convex problems (16). Therefore, the solving method of subproblem (16) directly affects the iterative efficiency of the CCCP algorithm.

Observe subproblem (16): $J'_{cav}(\tilde{\beta}^k) \cdot \tilde{\beta}$ is a linear function about $\tilde{\beta}$. $J_{vex}(\tilde{\beta})$ contains the convex function $\frac{1}{n} \sum_{i=1}^{n} F(Y_i - \rho^{k+1} \langle w_i, Y \rangle - X_i \tilde{\beta})$ and penalty term $\lambda \sum_{j=1}^{2p} P(|\tilde{\beta}_j|)$ for $\tilde{\beta}$. We might as well order

$$\psi(\tilde{\beta}) = \frac{1}{n} \sum_{i=1}^{n} F\left(Y_i - \rho^{k+1} \langle w_i, Y \rangle - X_i \tilde{\beta}\right) + J_{cav}'\left(\tilde{\beta}^k\right) \cdot \tilde{\beta},\tag{17}$$

where $\psi(\tilde{\beta})$ is a convex function about $\tilde{\beta}$. So, subproblem (16) can be represented as

$$\min_{\tilde{\beta}\in \mathbb{R}^p}\psi(\tilde{\beta}) + \lambda \sum_{i=1}^p P(|\tilde{\beta}_i|).$$
(18)

Optimization problems (17) are composed of convex functions and adaptive lasso penalty terms, and we can use the ISTA algorithm to solve such problems.

For all L > 0, ISTA approximates the function $F(\beta) = \psi(\tilde{\beta}) + \lambda \sum_{i=1}^{2p} \eta_i |\tilde{\beta}_i|$ at $\tilde{\beta} = \xi$ as

$$Q_L(\tilde{\beta},\xi) = \psi(\xi) + \langle \tilde{\beta} - \xi, \nabla \psi(\xi) \rangle + \frac{L}{2} \|\tilde{\beta} - \xi\|^2 + \lambda \sum_{i=1}^{2p} \eta_i |\tilde{\beta}_i|.$$
(19)

This function has the following minimum point:

$$\Theta_{L}(\xi) = \operatorname{argmin}_{\tilde{\beta} \in R^{2p}} Q_{L}(\tilde{\beta}, \xi)$$

= $\operatorname{argmin}_{\tilde{\beta} \in R^{2p}} \left\{ \lambda \sum_{i=1}^{2p} \eta_{i} |\tilde{\beta}_{i}| + \frac{L}{2} \left\| \tilde{\beta} - \left(\xi - \frac{1}{L} \nabla \psi(\xi) \right) \right\|^{2} \right\}$
= $S_{\lambda \eta/L} \left(\xi - \frac{1}{L} \nabla \psi(\xi) \right).$ (20)

With $\eta = [\eta_1, \ldots, \eta_{2p}] \in R^{2p}$, and for $\nu = \lambda \eta / L \in R^p_+, S_\alpha : \mathbb{R}^{2p} \to \mathbb{R}^{2p}$ the vectorformed soft-thresholding operator $S_v(\tilde{\beta}) = \overline{\tilde{\beta}}, \quad \overline{\tilde{\beta}}_i = (|\tilde{\beta}_i| - v_i)_+ \operatorname{sgn}(\tilde{\beta}_i), i = 1, \ldots, 2p.$

Thus, the solution of problem (11) can be simply expressed as $\tilde{\beta}^k = \Theta_L(\tilde{\beta}^{k-1})$.

In this article, we use the FISTA algorithm with a faster convergence speed than ISTA (Beck and Teboulle, 2009) [18]. The FISTA algorithm framework with backtracking steps is given below (Algorithm 3):

Algorithm 3: FISTA with Backtracking Step for solving (17)

Require: $A, \xi, w\lambda > 0$ Ensure: solution $\tilde{\beta}$ 1: Step 0. Select $L^0 > 0, \eta > 1, \tilde{\beta}^0 \in \mathbb{R}^{2p}$ Let $\xi^1 = \tilde{\beta}^0, t^1 = 1$ 2: Step $k(k \ge 1)$. 3: Determine the smallest non-negative integer i^k which make $\bar{L} = \eta^{i^k} L^{k-1}$ satisfy 4 $F\left(\Theta_{\bar{L}}\left(\xi^k\right)\right) \le Q_{\bar{L}}\left(\Theta_{\bar{L}}\left(\xi^k\right), \xi^k\right).$ 5: Let $L^k = \eta^{i^k} L^{k-1}$ according to (19), calculate: 6: $\tilde{\beta}^k = \Theta_{L^k}\left(\xi^k\right)$ 7: $t^{k+1} = \frac{1}{2}\left[1 + \sqrt{1 + 4(t^k)^2}\right]$

8:
$$\boldsymbol{\xi}^{k+1} = \tilde{\boldsymbol{\beta}}^k + \frac{t^k - 1}{t^{k+1}} \left(\tilde{\boldsymbol{\beta}}^k - \tilde{\boldsymbol{\beta}}^{k-1} \right)$$

9: Output $\tilde{\beta} := \tilde{\beta}^k$.

So far, we completed the solution of subproblem (9).

4. Numerical Examples

We designed five numerical experiments to verify the performance and accuracy of variable selection methods under different conditions. For example, there are abnormal values in dependent variable Y and too many insignificant covariates.

Data generation will be based on model (1). We make the covariance matrix X an $n \times (q+3)$ matrix, and the X obeys the (q+3)-dimensional normal distribution, the mean value is zero, and the covariance matrix is (σ_{ij}) , where $\sigma_{ij} = 0.5^{|i-j|}$. This means that the number of samples is n, the number of significant covariates is 3, and the number of insignificant covariates is q. In the following experiments, we set n and q to $n \in \{200, 360, 500\}$ and $q \in \{5, 20, 40, 60\}$. For the spatial regression coefficient ρ , in the experiment, we set it to $\rho \in \{0.2, 0.5, 0.8\}$.

We define the spatial weights matrix as a k-diagonal matrix, i.e., a matrix with only the main diagonal and the k-1 skew diagonals around it as element 1, and the other elements as 0. In numerical experiments, we set k = 7.

The regression coefficient β is set to: $\beta = (\beta_1, \beta_2, \beta_3, 0_q)$, where $(\beta_1, \beta_2, \beta_3) = (3, 2, 1.5)$. The regression coefficient vector of exogenous variables δ is set to: $\delta = (\delta_1, \delta_2, \delta_3, 0_q)$, where $(\delta_1, \delta_2, \delta_3) = (1.5, 1.2, 1)$, and where 0_q is a zero vector and its dimension is q; this means that the number of 0 elements of β and δ that we set in the experiment is the same, both of which are q. The dependent variable Y is generated by the model (2).

For the error term, let $\varepsilon \sim \mathcal{N}(0, \sigma^2 I_n), \sigma^2$ obeys uniform distribution, and its generation interval is $[\sigma_1 - 0.1, \sigma_1 + 0.1]$ with $\sigma_1 \in \{1, 2\}$. Of course, in practice, the observation noise does not completely conform to the Gaussian distribution, and there may be abnormal values in the response. The abnormal values in the response are discussed in Section 4.3.

To reflect the excellence of this model, we also used square loss and LAD to compare with our exponential square loss. To ensure the accuracy of the experiment, we repeated each experiment 100 times. The following results are the results of MSE in the middle of 100 repeated experiments.We express the median of MSE as MedSE.

4.1. Nonregular Estimation of Normal Data

In this section, we conduct experiments on the condition that q = 5, the noise is Gaussian noise, and the penalty term is not set for the parameter estimation model. The results are shown in Table 1. Square, Exp, and LAD represent square loss, exponential square loss, and LAD loss, respectively. (1) This shows that Exp, Square, and LAD made estimates of β and δ , which are close to typical values (the means of the true values of *beta*1,

 β 2, and β 3 are 3.0, 2.0, and 1.6, then, the mean sof the true values of δ_1 , δ_2 , and δ_3 are 2.0, 1.5, and 1.0.). By comparison, the estimated value obtained by the square loss model is the best. (2) For MedSE, the square loss model also performs the best. (3) The three loss functions can give accurate estimates of the spatial autoregressive coefficients ρ .

		n = 200, 2q =	: 10	1	a = 360, 2q =	: 10	1	n = 500, 2q =	10
	Exp	Square	LAD	Exp	Square	LAD	Exp	Square	LAD
$\rho = 0.8, \sigma = 1$ β_1 β_2 β_3 δ_1 δ_2 δ_3 $\hat{\rho}$ MedSE	3.0904 2.0303 1.6422 1.242 1.5109 0.8155 0.8001 0.5994	2.6866 1.9449 1.4689 1.3382 1.3963 1.0711 0.8011 0.4158	3.2503 1.8899 1.3069 1.3582 1.0625 0.7999 0.4693	3.1335 1.9594 1.5725 1.504 1.1245 1.1101 0.7999 0.2518	$\begin{array}{c} 2.8487\\ 1.9498\\ 1.5409\\ 1.6117\\ 1.132\\ 1.0575\\ 0.8006\\ 0.2827\end{array}$	3.0801 2.0897 1.3781 1.3924 1.3174 0.9693 0.7997 0.3432	2.8084 2.0949 1.6174 1.4604 1.1786 1.0903 0.8002 0.248	2.7947 2.1354 1.7394 1.2156 1.1139 1.0092 0.7979 0.234	2.9486 2.0207 1.6394 1.3616 1.111 1.0871 0.7981 0.3086
$\rho = 0.5, \sigma = 1$ β_1 β_2 β_3 δ_1 δ_2 δ_3 $\hat{\rho}$ MedSE	3.0854 2.0058 1.6799 1.2338 1.4943 0.8849 0.5021 0.6007	3.0349 2.1532 1.3788 1.219 1.4233 1.0766 0.4999 0.388	3.0617 1.927 1.6744 1.7939 1.5202 0.5961 0.5 0.4564	3.1254 1.9556 1.5702 1.4848 1.1322 1.1036 0.5003 0.2808	$\begin{array}{c} 2.8039\\ 2.1823\\ 1.4268\\ 1.1814\\ 1.3266\\ 0.9614\\ 0.5\\ 0.2829\end{array}$	3.0451 2.2277 1.7227 1.3734 1.3411 0.8644 0.5 0.3287	$\begin{array}{c} 2.8058\\ 2.0986\\ 1.6208\\ 1.458\\ 1.186\\ 1.0966\\ 0.4998\\ 0.2452\end{array}$	3.1899 1.9975 1.6813 1.4145 0.9884 0.9671 0.4999 0.2262	3.2542 2.0256 1.303 1.6612 1.3373 0.9961 0.4999 0.2939
$\rho = 0.2, \sigma = 1$ β_1 β_2 β_3 δ_1 δ_2 δ_3 δ_1 δ_2 δ_3 δ_1 MedSE	3.0072 1.8903 1.5386 0.8622 1.4427 0.6609 0.2417 0.9037	$\begin{array}{c} 2.7008\\ 1.7081\\ 1.4297\\ 1.2241\\ 1.0584\\ 0.5715\\ 0.2419\\ 0.8134 \end{array}$	2.7579 1.93 1.571 0.9667 0.7845 1.1202 0.2519 0.9757	3.0283 1.8152 1.4646 1.2297 0.8279 1.0265 0.2271 0.5492	$\begin{array}{c} 2.9572\\ 2.081\\ 1.2788\\ 1.2184\\ 0.8104\\ 0.9351\\ 0.2365\\ 0.6286\end{array}$	2.9077 1.9718 1.3697 1.015 1.2721 0.4027 0.2437 0.7235	$\begin{array}{c} 2.635 \\ 1.8603 \\ 1.4512 \\ 1.0963 \\ 0.7684 \\ 0.7551 \\ 0.25 \\ 0.9921 \end{array}$	2.8836 2.0794 1.5486 1.184 0.889 0.8819 0.2216 0.5287	3.0321 1.4453 1.5858 1.1689 1.0247 0.9224 0.2317 0.6407
$\rho = 0.8, \sigma = 2$ β_1 β_2 β_3 δ_1 δ_2 $\delta3$ $\hat{\rho}$ MedSE	3.0727 2.1164 1.6747 0.9807 1.7814 0.7301 0.801 1.2058	$\begin{array}{c} 2.8882 \\ 1.8141 \\ 1.3996 \\ 1.5773 \\ 0.8604 \\ 0.9358 \\ 0.8045 \\ 0.7731 \end{array}$	3.0548 2.0596 1.4515 1.6723 0.9949 0.8807 0.7943 0.9502	3.2634 2.0387 1.7457 1.4112 1.0979 1.0707 0.7996 0.5131	$\begin{array}{c} 2.7795\\ 1.6757\\ 1.7597\\ 1.3493\\ 1.1257\\ 0.6032\\ 0.8042\\ 0.5493\end{array}$	3.423 2.1636 1.7635 1.5103 1.0441 1.0806 0.7978 0.6914	$\begin{array}{c} 2.6808\\ 2.1004\\ 1.5395\\ 1.4831\\ 1.1439\\ 1.0423\\ 0.7989\\ 0.5536\end{array}$	$\begin{array}{c} 3.1902 \\ 1.9817 \\ 1.3549 \\ 1.6088 \\ 1.0058 \\ 1.0058 \\ 1.0172 \\ 0.7989 \\ 0.4719 \end{array}$	3.0531 1.9528 1.2795 1.5378 1.3651 0.8286 0.7897 0.5597
$\rho = 0.5, \sigma = 2$ β_1 β_2 β_3 δ_1 δ_2 δ_3 δ_1 δ_2 δ_3 δ_1 MedSE	3.0762 2.0839 1.7169 0.9706 1.7802 0.8067 0.5044 1.2459	3.1916 2.2408 1.384 1.3618 1.1337 1.2863 0.5035 0.8065	$\begin{array}{c} 3.1159 \\ 1.5295 \\ 1.8689 \\ 1.483 \\ 1.2179 \\ 1.0691 \\ 0.5 \\ 0.9201 \end{array}$	3.2325 1.8826 1.6075 1.437 1.0509 1.2173 0.5007 0.6033	3.0528 1.9422 1.7206 1.5269 1.0885 0.8635 0.4996 0.5434	2.8093 2.1974 1.5534 1.7395 1.4445 0.901 0.4986 0.6822	$\begin{array}{c} 2.6731 \\ 2.1207 \\ 1.5565 \\ 1.4862 \\ 1.1825 \\ 1.0777 \\ 0.4994 \\ 0.5387 \end{array}$	3.0178 2.1175 1.355 1.3611 1.2947 0.9428 0.4973 0.4699	3.0432 1.9125 1.2292 1.1453 1.3262 0.9601 0.4998 0.5622
$\rho = \overline{0.2, \sigma} = 2$ β_1 β_2 β_3 δ_1 δ_2 δ_3 $\hat{\rho}$ MedSE	2.9838 1.9525 1.5477 0.5511 1.6614 0.5739 0.2624 1.5138	3.0811 1.8371 1.5198 1.5501 0.7911 0.3887 0.2341 1.1417	3.0512 2.3438 1.2998 1.1878 1.9069 0.1885 0.2342 1.5123	$\begin{array}{c} 3.1253\\ 1.7215\\ 1.4741\\ 1.1662\\ 0.6863\\ 1.1021\\ 0.235\\ 0.8143\end{array}$	2.965 1.9963 1.6448 1.1816 0.9868 0.8204 0.2307 0.816	2.7197 2.2379 1.0015 0.8059 1.2406 0.8706 0.2313 0.907	$\begin{array}{c} 2.5382 \\ 1.8893 \\ 1.4075 \\ 1.1623 \\ 0.8082 \\ 0.7875 \\ 0.2472 \\ 1.0921 \end{array}$	2.8219 1.886 1.7982 1.0545 1.0211 0.5343 0.2349 0.7422	2.7438 1.8351 1.7257 1.1654 1.3427 0.8185 0.238 0.8705

Table 1. Nonregular estimation of normal data (q = 5).

4.2. Nonregular Estimation for High-Dimensional Data

In this subsection, we made $q \in \{20, 40, 60\}$, and the parameter estimation results of the model on normal data with huge sample dimensions are explained. The results are shown in Table 2. It can be found that the estimation of β , δ , and ρ of any model is far less effective than that of q = 5. The results of MedSE are also not satisfactory. Due to the insufficient number of samples, such results can be expected.

	n = 200, 2q = 40			n=360, 2q=80			n	n = 500, 2q = 120		
	Exp	Square	LAD	Exp	Square	LAD	Exp	Square	LAD	
$\rho = 0.8, \sigma = 1$ β_1 β_2 β_3 δ_1 δ_2 δ_3 δ_1 δ_2 δ_3 δ MedSE	$\begin{array}{c} 2.9991 \\ 1.87 \\ 1.6641 \\ 1.5449 \\ 1.1133 \\ 1.104 \\ 0.7841 \\ 1.0389 \end{array}$	2.9355 1.8534 1.7286 1.4123 1.2998 1.0165 0.7976 1.1975	2.8898 2.148 1.3997 1.2135 1.4747 0.8516 0.7812 2.4913	3 2.2471 1.7933 1.1743 1.193 1.5647 0.7814 3.9024	3.1579 2.097 1.4782 0.9747 1.006 0.6918 0.7626 1.3519	3.273 1.7728 1.3566 1.6914 0.6772 0.8134 0.7698 3.216	$\begin{array}{c} 2.1076 \\ 1.4959 \\ 1.2695 \\ 0.8435 \\ 1.5524 \\ 0.9871 \\ 0.7785 \\ 3.5719 \end{array}$	2.8197 2.1326 1.5688 1.4414 1.0102 0.8072 0.7921 1.4983	3.0519 2.4031 1.1838 1.6309 1.3473 0.3043 0.7578 3.4753	
$\rho = 0.5, sigma = 1$ β_1 β_2 β_3 δ_1 δ_2 δ_3 δ_1 δ_2 δ_3 δ MedSE	$\begin{array}{c} 2.9674\\ 1.9594\\ 1.7014\\ 1.401\\ 1.3144\\ 1.1186\\ 0.4984\\ 0.7268\end{array}$	3.0461 1.956 1.5528 1.4171 1.3097 1.0907 0.499 0.8361	3.0203 2.2012 1.4707 1.8907 1.1209 0.9799 0.5 1.0131	$\begin{array}{c} 2.8981 \\ 2.2222 \\ 1.698 \\ 1.4537 \\ 1.3471 \\ 0.9439 \\ 0.5008 \\ 0.846 \end{array}$	3.3433 1.8716 1.5486 1.3244 1.0562 0.977 0.5004 0.8735	2.8355 2.041 1.533 1.4937 1.4014 1.1517 0.5 1.0286	3.0614 2.1271 1.4004 1.4151 1.2555 1.3021 0.4997 1.152	3.0561 2.0035 1.8345 1.6572 1.1235 1.2148 0.5007 0.9011	$\begin{array}{c} 2.7106\\ 2.1597\\ 1.5845\\ 1.516\\ 0.9392\\ 1.0304\\ 0.5\\ 1.1143\end{array}$	
$\rho = 0.2, \sigma = 1$ β_1 β_2 β_3 δ_1 δ_2 δ_3 δ_1 MedSE	$\begin{array}{c} 2.6774\\ 1.8758\\ 1.6073\\ 0.6648\\ 1.3401\\ 0.8527\\ 0.2731\\ 1.6545\end{array}$	$\begin{array}{c} 2.7337\\ 1.5391\\ 1.1691\\ 0.1578\\ 0.5349\\ 0.2299\\ 0.368\\ 3.3194 \end{array}$	$\begin{array}{c} 2.2269\\ 2.0274\\ 1.5513\\ -0.522\\ 0.2749\\ 1.1821\\ 0.424\\ 4.2853\end{array}$	$\begin{array}{c} 2.7222\\ 1.9726\\ 1.5515\\ 1.3422\\ 1.0115\\ 0.4795\\ 0.309\\ 2.3481\end{array}$	$\begin{array}{c} 2.6372 \\ 1.4435 \\ 1.4517 \\ 0.5994 \\ 0.1372 \\ 0.6847 \\ 0.3721 \\ 3.2567 \end{array}$	$\begin{array}{c} 2.5639\\ 1.4924\\ 1.9718\\ 0.0136\\ 0.3757\\ -0.254\\ 0.4449\\ 4.9894 \end{array}$	$\begin{array}{c} 2.7679\\ 1.7517\\ 1.3588\\ 0.4778\\ -0.186\\ 0.4057\\ 0.429\\ 5.0342\end{array}$	$\begin{array}{c} 2.625\\ 1.8328\\ 1.4107\\ 0.5541\\ -0.095\\ 0.9201\\ 0.4203\\ 3.9322 \end{array}$	$\begin{array}{c} 2.3669\\ 1.7334\\ 1.5919\\ 0.5164\\ -0.533\\ -0.205\\ 0.4601\\ 4.9488\end{array}$	
$\rho = 0.8, \sigma = 2$ β_1 β_2 β_3 δ_1 δ_2 δ_3 $\hat{\rho}$ MedSE	3.0412 1.8198 1.5695 1.3762 1.4344 0.9648 0.7775 1.9747	$\begin{array}{c} 2.9279\\ 1.731\\ 1.6793\\ 1.8133\\ 1.0001\\ 1.2588\\ 0.7886\\ 1.9176\end{array}$	3.2539 1.3838 1.9783 0.6013 1.841 1.7202 0.7847 3.2383	3.0001 2.2004 1.818 1.0067 1.3286 1.4715 0.7808 4.2531	2.9425 2.009 1.7579 1.2925 1.3243 1.123 0.7977 2.1049	3.3778 1.6173 2.1066 1.0236 1.0053 1.021 0.7196 3.9547	$\begin{array}{c} 3.0661 \\ 2.0961 \\ 1.3409 \\ 1.2358 \\ 1.0631 \\ 1.629 \\ 0.787 \\ 2.641 \end{array}$	$\begin{array}{c} 2.9154 \\ 1.7923 \\ 1.8222 \\ 1.3486 \\ 1.0136 \\ 0.9606 \\ 0.7941 \\ 2.139 \end{array}$	2.9285 1.9155 1.9951 1.1922 1.148 0.9799 0.7652 4.1654	
$\rho = 0.5, \sigma = 2$ β_1 β_2 β_3 δ_1 δ_2 δ_3 δ MedSE	2.9924 1.937 1.6366 1.1814 1.6792 1.0417 0.4982 1.7761	$\begin{array}{c} 3.127\\ 1.7325\\ 1.9284\\ 1.3155\\ 1.4335\\ 0.9891\\ 0.5005\\ 1.6992\end{array}$	3.3996 2.2317 1.3111 1.5799 1.1793 1.1182 0.5 2.0438	$\begin{array}{c} 2.9071\\ 2.1674\\ 1.7516\\ 1.3022\\ 1.5096\\ 0.8649\\ 0.5013\\ 1.9405\end{array}$	$\begin{array}{c} 2.9556\\ 2.1939\\ 1.5356\\ 1.3435\\ 1.4914\\ 0.8226\\ 0.5015\\ 1.7439\end{array}$	3.2707 1.9265 1.6939 1.781 0.8587 1.0424 0.5 2.1404	3.1224 2.1034 1.329 1.3229 1.1117 1.6105 0.4992 2.3691	2.9192 2.0763 1.2335 1.2476 1.3711 0.9617 0.4996 1.8227	2.7059 2.0959 1.5695 1.5013 0.9883 1.3937 0.5 2.2753	
$\rho = \overline{0.8, \sigma} = 2$ β_1 β_2 β_3 δ_1 δ_2 δ_3 $\hat{\rho}$ MedSE	$\begin{array}{c} 2.7126\\ 1.8389\\ 1.5575\\ 0.5004\\ 1.6485\\ 0.8129\\ 0.2716\\ 2.1955\end{array}$	$\begin{array}{c} 2.463\\ 1.7703\\ 1.302\\ 0.375\\ 0.5337\\ 0.676\\ 0.3598\\ 3.5415 \end{array}$	3.0218 1.1025 1.5166 1.4855 -0.864 -0.697 0.5 5.0709	$\begin{array}{c} 2.7365\\ 1.9165\\ 1.5931\\ 1.1763\\ 1.0869\\ 0.4247\\ 0.311\\ 2.8997\end{array}$	$\begin{array}{c} 2.6766\\ 1.625\\ 1.1868\\ 0.9432\\ 0.1892\\ -0.641\\ 0.355\\ 3.8662\end{array}$	$\begin{array}{c} 2.3112\\ 1.8649\\ 1.3592\\ -0.186\\ -0.044\\ 0.3875\\ 0.4821\\ 5.2612 \end{array}$	$\begin{array}{c} 2.819 \\ 1.7335 \\ 1.299 \\ 0.4655 \\ -0.227 \\ 0.6011 \\ 0.4261 \\ 5.3588 \end{array}$	$\begin{array}{c} 2.7498\\ 1.7638\\ 1.4037\\ 0.406\\ 0.1407\\ -0.289\\ 0.4147\\ 4.256\end{array}$	$\begin{array}{c} 2.8593 \\ 1.7436 \\ 1.2125 \\ 0.846 \\ -1.145 \\ 0.6946 \\ 0.4469 \\ 5.583 \end{array}$	

Table 2. Nonregular estimation for high-dimensional data.

4.3. Nonregular Estimation of Data with Outliers in Dependent Variable y

In this subsection, we make the error term ϵ obey the mixed Gaussian distribution $(1 - \xi_1) \cdot \mathcal{N}(0, 1) + \xi_1 \cdot \mathcal{N}(10, 6^2)$, where $\xi_1 \in \{0.01, 0.05\}$. In this case, the observed y will have many outliers. We illustrate the results (Table 3) of the estimated coefficients of β and δ when the observations of y have outliers. (1) For MedSE, unlike the results in Table 1, where y has no outliers, in almost all tests in Table 3, exponential square loss performed the best. (2) By comparison, the estimated values of β and δ obtained by the exponential square loss is also the best. Therefore, we can conclude that when y has outliers, the SDM based on exponential square loss has good robustness.

	n=200, 2q=10		1	n=360, 2q=10			n = 500, 2q = 10		
	Exp	Square	LAD	Exp	Square	LAD	Exp	Square	LAD
$\rho = 0.8, \sigma = 1, \xi = 0.01$ β_1 β_2 β_3 δ_1 δ_2 δ_3 δ_1 MedSE	3.053 2.213 1.577 1.341 1.311 0.876 0.801 0.609	$\begin{array}{c} 3.333 \\ 1.48 \\ 1.579 \\ 1.983 \\ 0.959 \\ 1.224 \\ 0.791 \\ 1.866 \end{array}$	2.873 1.958 2.046 0.811 1.736 0.644 0.798 1.009	3.02 2.125 1.57 1.444 1.127 1.182 0.8 0.398	$\begin{array}{c} 3.237\\ 2.048\\ 1.857\\ 0.966\\ 1.113\\ 1.284\\ 0.8\\ 1.392\end{array}$	2.754 1.836 1.908 1.7 1.199 0.752 0.786 0.755	$\begin{array}{c} 2.882\\ 2.126\\ 1.604\\ 1.464\\ 0.962\\ 1.119\\ 0.794\\ 0.405\end{array}$	$\begin{array}{c} 3.102 \\ 1.948 \\ 1.677 \\ 1.165 \\ 0.906 \\ 0.382 \\ 0.756 \\ 1.268 \end{array}$	2.827 1.754 1.757 1.344 1.722 0.954 0.799 0.623
$\rho = 0.5, \sigma = 1, \xi = 0.01$ β_1 β_2 β_3 δ_1 δ_2 δ_3 δ MedSE	3.035 2.217 1.561 1.143 1.395 0.929 0.5 0.738	$\begin{array}{c} 3.371 \\ 2.19 \\ 1.875 \\ 1.37 \\ 0.621 \\ 0.918 \\ 0.497 \\ 1.341 \end{array}$	3.049 2.259 1.646 1.623 1.193 0.768 0.5 0.84	3.036 2.094 1.555 1.448 1.136 1.119 0.5 0.347	$\begin{array}{c} 3.123 \\ 1.672 \\ 2.008 \\ 1.432 \\ 1.642 \\ 0.742 \\ 0.499 \\ 1.097 \end{array}$	3.179 1.981 1.582 1.227 0.954 1.213 0.499 0.627	2.874 2.14 1.699 1.492 1.075 1.186 0.501 0.341	3.031 1.687 1.474 1.464 1.472 1.161 0.5 0.959	2.972 2.249 1.436 1.264 1.079 1.138 0.5 0.511
$ \begin{array}{l} \rho = 0.2, \sigma = 1, \xi = 0.01 \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \delta_1 \\ \delta_2 \\ \delta_3 \\ \hat{\rho} \\ \text{MedSE} \end{array} $	3.032 2.085 1.52 0.882 1.667 0.723 0.213 1.091	2.817 2.497 1.224 1.381 1.117 0.49 0.207 1.509	2.936 1.757 1.654 1.323 0.84 1.022 0.222 1.065	$\begin{array}{c} 2.961 \\ 1.946 \\ 1.441 \\ 1.216 \\ 0.879 \\ 0.996 \\ 0.225 \\ 0.556 \end{array}$	3.077 1.646 1.817 0.883 0.785 1.207 0.23 1.076	3.17 1.802 1.329 1.478 0.688 0.72 0.223 0.807	$\begin{array}{c} 2.7\\ 1.869\\ 1.506\\ 1.151\\ 0.631\\ 0.809\\ 0.256\\ 1.061\end{array}$	$\begin{array}{c} 3.333 \\ 1.9 \\ 1.614 \\ 1.654 \\ 1.323 \\ 0.618 \\ 0.188 \\ 1.052 \end{array}$	2.752 1.853 1.479 1.018 1.034 1.096 0.232 0.625
$ \begin{split} \rho &= 0.8, \sigma = 1, \xi = 0.05 \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \delta_1 \\ \delta_2 \\ \delta_3 \\ \hat{\rho} \\ \text{MedSE} \end{split} $	2.845 2.228 1.85 1.41 0.863 0.869 0.796 0.984	2.064 3.473 1.076 2.957 0.258 3.323 0.788 4.778	3.28 1.405 2.271 0.344 2.166 0.661 0.782 2.45	$\begin{array}{c} 2.914\\ 2.148\\ 1.703\\ 1.361\\ 1.05\\ 1.255\\ 0.799\\ 0.716\end{array}$	3.493 2.501 0.18 0.522 -0.3 2.546 0.794 3.707	3.116 1.688 1.86 0.856 2.015 0.566 0.793 1.366	2.935 2.109 1.675 1.791 0.899 0.95 0.788 0.835	3.857 2.336 2.19 1.181 0.782 1.568 0.77 3.077	3.369 1.564 1.487 1.062 2.052 1.107 0.789 1.048
$\rho = 0.5, \sigma = 1, \xi = 0.05$ β_1 β_2 β_3 δ_1 δ_2 δ_3 δ MedSE	2.894 2.169 1.766 1.215 1.283 0.833 0.497 0.627	3.636 1.293 0.053 -0.05 1.593 0.489 0.473 3.878	$\begin{array}{c} 2.978\\ 2.131\\ 1.419\\ 0.866\\ 2.262\\ 0.595\\ 0.5\\ 1.536\end{array}$	$\begin{array}{c} 2.8\\ 2.177\\ 1.572\\ 1.436\\ 1.02\\ 1.068\\ 0.499\\ 0.506\end{array}$	3.655 2.57 1.919 1.79 0.449 -0.27 0.494 2.989	2.922 2.54 1.607 1.238 1.381 0.837 0.501 0.996	$\begin{array}{c} 2.882\\ 2.149\\ 1.705\\ 1.729\\ 1.01\\ 1.042\\ 0.499\\ 0.799\end{array}$	3.027 1.971 2.184 0.58 1.381 -0.14 0.499 2.467	3.152 2.333 1.673 0.811 1.316 1.032 0.5 0.803
$\rho = \overline{0.2, \sigma = 1, \xi = 0.05}$ β_1 β_2 β_3 δ_1 δ_2 δ_3 $\hat{\rho}$ MedSE	3.111 2.294 1.791 1.687 1.536 0.805 0.133 0.975	2.373 3.955 0.154 2.849 1.503 1.732 0.034 3.461	3.24 2.314 1.293 1.613 0.8 1.143 0.183 1.326	2.597 2.195 1.335 1.542 0.98 0.884 0.189 1.074	3.516 2.032 1.097 2.57 1.2 0.344 0.088 2.428	2.797 1.871 1.344 1.161 1.125 0.76 0.226 0.822	$\begin{array}{c} 2.784\\ 2.013\\ 1.589\\ 1.418\\ 0.8\\ 0.934\\ 0.221\\ 0.81\end{array}$	$3.617 \\ 1.844 \\ 1.296 \\ 1.035 \\ 0.872 \\ 1.147 \\ 0.174 \\ 2.047$	2.885 2.128 1.467 1.627 1.222 1.145 0.189 0.67

Table 3. Nonregular estimation of data with outliers in dependent variable y.

4.4. Nonregular Estimation of Data with Noise in Spatial Weight Matrix

In this section, we simulate the presence of noise in the spatial weight matrix. We added a minor disturbance term ϵ' to each nonzero element in the spatial weight matrix W, where $\epsilon' \sim (1 - \xi_2) \cdot \mathcal{N}(0, 0.001) + \xi_2 \cdot \mathcal{N}(0, 1), \xi_2 \in \{0.01, 0.03, 0.05\}$, and all the simulated data are generated with $\rho = 0.5, \sigma = 1$. The test results are shown in Table 4. Compared with normal data (Table 1), the MedSE value increased. Additionally, for each loss function, the estimation of β , δ , and ρ also worsens. When the weight matrix has noise, the exponential square loss and LAD loss have good performance. Compared with the square loss, they have more accurate estimates of the parameters and smaller MedSE values. However, it cannot be denied that LAD loss performs better than exponential square loss.

		n = 200, 2q =	: 10	1	n = 360, 2q =	= 10	1	n = 500, 2q =	: 10
	Exp	Square	LAD	Exp	Square	LAD	Exp	Square	LAD
$\rho = 0.5, \sigma = 1, \xi = 0.01$									
β_1	3.125	3.143	2.909	3.286	2.614	3.142	2.82	2.138	2.895
β_2	1.692	1.89	1.934	1.3	2.39	2.025	2.07	2.367	1.826
β ₃	1.919	1.633	1.597	1.68	0.716	1.622	1.651	2.761	1.473
δ_1	1.167	0.737	1.612	0.997	1.4	1.418	1.365	0.86	1.452
δ_2	1.422	1.235	1.059	-0.28	0.584	1.318	1.247	0.485	1.191
δ_3	0.898	0.038	1.076	2.083	1.16	1.273	0.978	0.799	1.164
$\hat{\rho}$	0.501	0.492	0.496	0.501	0.477	0.5	0.5	0.486	0.5
MedSE	0.636	2.596	0.562	2.657	2.623	0.411	0.275	2.591	0.341
$\rho = 0.5, \sigma = 1, \xi = 0.03$									
β_1	2.941	1.728	2.955	1.38	2.193	3.15	2.83	2.295	2.963
β_2	1.143	2.575	2.278	0.298	2.322	1.952	1.88	2.292	2.002
β ₃	1.771	2.299	1.386	1.195	0.383	1.267	1.867	1.809	1.619
δ_1	1.008	-0.63	1.607	0.019	0.211	1.156	0.396	2.348	1.218
δ_2	0.605	2.475	1.326	0.255	0.283	0.961	0.974	1.317	0.994
δ_3	0.579	0.146	0.922	0.58	0.826	0.921	0.881	0.699	1.136
ρ̂	0.503	0.468	0.495	0.503	0.449	0.499	0.494	0.45	0.498
MedSE	1.561	3.925	0.819	3.645	3.922	0.547	1.227	4.972	0.454
$\rho = 0.5, \sigma = 1, \xi = 0.05$									
β_1	3.02	1.981	3.046	3.183	2.054	2.849	2.893	2.507	3.187
β_2	1.479	0.857	2.072	1.259	0.636	2.253	1.911	0.865	2.265
β_3	1.978	0.753	0.897	1.721	1.649	1.299	1.827	2.039	1.362
δ_1	0.837	0.645	1.349	0.785	0.67	1.362	0.61	0.934	0.962
δ_2	1.557	-0.56	1.26	0.551	-0.68	1.275	0.813	1.384	1.135
δ_3	-0.23	-0.13	0.379	1.105	0.509	0.536	1.067	-1.56	1.24
$\hat{\rho}$	0.502	0.431	0.489	0.504	0.431	0.493	0.493	0.459	0.491
MedSE	1.922	5.079	1.207	2.191	4.588	0.805	1.034	5.232	0.69

Table 4. Nonregular Estimation of Data with Noise in Spatial Weight Matrix.

4.5. Estimation with Adaptive-l1 Regularizer

We add adaptive L1 regularization to the loss function in this section and conduct experiments. We also record the average number of zero coefficients correctly selected by the model as "Correct" and the average number of nonzero coefficients incorrectly judged by the model as "Incorrect".

Table 5 shows the results of adaptive lasso regularization on normal data with q = 5. The results show that, under almost all test results, the SDM model with exponential square loss and adaptive lasso cannot only identify more true zero coefficients ("Correct" with exponential square loss model is almost twice as much as that with square loss and LAD loss model) and nearly zero 'Incorrect' numbers but also has the best MedSE and accurate estimation of $\hat{\rho}$.

Table 6 shows the results of adaptive lasso regularization on normal data with $q \in \{20, 40, 60\}$. The results show that when there are too many insignificant covariates, the accuracy of the results of the Square loss model with adaptive lasso and the LAD loss model with adaptive lasso decreases significantly. However, the model with adaptive lasso and exponential square loss is still accurate. It can identify more true "Correct" numbers and nearly zero "Incorrect" numbers and has the best MedSE and precise estimation of $\hat{\rho}$.

Table 7 shows the results of estimation with adaptive-l1 regularization when the observations of y have outliers. The results show that, in almost all test results, the exponential square loss model with adaptive L1 has identified more true zero coefficients ("Correct") and, in most cases, has lower MedSE. Compared with the model without regularization term (Table 3), the model with adaptive L1 has a better effect. In the test, the exponential square loss model with adaptive L1 identified at least 8 zero coefficients and, in most cases, determined 10 zero coefficients. For MedSE, the exponential square loss model with adaptive L1 identified at least 8 zero coefficients and, in most cases, determined 10 zero coefficients. For MedSE, the exponential square loss model with adaptive L1 has the smallest MedSE in all cases, except that when n = 500 and 2q = 10, the MedSE in some cases is slightly larger than the LAD loss model with adaptive L1. This shows that the SDM using exponential square loss and adaptive lasso has excellent variable selection ability and strong robustness when the Y observation has outliers.

	n=200, 2q=10			1	i = 360, 2q =	10	n	a = 500, 2q =	10
	Exp	Square	LAD	Exp	Square	LAD	Exp	Square	LAD
$\rho = 0.8, \sigma = 1$ Correct Incorrect $\hat{\rho}$ MedSE	10	5.23	5.78	10	5.53	5.61	10	5.61	5.64
	0	0	0	0	0	0	0	0	0
	0.8008	0.8035	0.8011	0.7999	0.8014	0.7982	0.7997	0.7995	0.801
	0.3747	0.3887	0.4697	0.1468	0.2843	0.3259	0.1316	0.2374	0.2944
$ \begin{array}{l} \rho = 0.5, \sigma = 1 \\ \text{Correct} \\ \text{Incorrect} \\ \hat{\rho} \\ \text{MedSE} \end{array} $	10	5.27	5.61	10	5.47	5.74	10	5.69	5.75
	0	0	0	0	0	0	0	0	0
	0.5013	0.5008	0.502	0.5001	0.5003	0.5005	0.4999	0.4997	0.5005
	0.3575	0.3514	0.4354	0.1342	0.2751	0.3161	0.1207	0.2217	0.2699
$\rho = 0.2, \sigma = 1$ Correct Incorrect $\hat{\rho}$ MedSE	10	5.42	5.52	9.98	5.42	5.36	9	5.3	5.46
	0	0.05	0.14	0	0.01	0.03	0	0	0
	0.2351	0.2375	0.2508	0.2257	0.231	0.2426	0.245	0.2265	0.2407
	0.7905	0.8637	1.0335	0.4758	0.6328	0.7898	0.8372	0.5565	0.6443
$\rho = 0.8, \sigma = 2$ Correct Incorrect $\hat{\rho}$ MedSE	10	5.34	5.12	10	5.1	5.42	9	5.17	5.18
	0	0	0	6	0	0	0	0	0
	0.8017	0.7992	0.8087	0.5033	0.7988	0.8018	0.8002	0.8036	0.7986
	0.8202	0.7826	0.9687	4.5219	0.5524	0.6503	0.2753	0.4729	0.5452
$\begin{array}{l} \rho = 0.5, \sigma = 2\\ \text{Correct}\\ \text{Incorrect}\\ \hat{\rho}\\ \text{MedSE} \end{array}$	10	5.34	5.21	10	5.4	5.27	9	5.27	5.23
	0	0	0	0	0	0	0	0	0
	0.5034	0.5014	0.4998	0.5005	0.5001	0.5001	0.4997	0.4988	0.4998
	0.813	0.75	0.9107	0.3039	0.5554	0.6583	0.2723	0.4426	0.5261
$\rho = 0.2, \sigma = 2$ Correct Incorrect $\hat{\rho}$ MedSE	8	5.17	5.11	9	5.51	5.29	9	5.31	5.27
	0	0.03	0.23	0	0.01	0.02	0	0	0
	0.2601	0.2382	0.2301	0.2359	0.2241	0.2535	0.2432	0.246	0.246
	1.3903	1.0942	1.3318	0.6905	0.7508	1.0301	0.8508	0.7031	0.8261

Table 5. Estimation with adaptive-l1 regularizer on normal data (q = 5).

 Table 6. Estimation with adaptive-l1 regularizer on normal data of high dimension.

	i	n = 200, 2q =	40	1	n = 360, 2q =	80	n	a = 500, 2q =	120
	Exp	Square	LAD	Exp	Square	LAD	Exp	Square	LAD
$ \begin{array}{l} \rho = 0.8, \sigma = 1 \\ \text{Correct} \\ \text{Incorrect} \\ \hat{\rho} \\ \text{MedSE} \end{array} $	40	21.32	21.19	80	42.42	42.63	119.01	65.11	64.21
	0	0.02	0.04	0	0	0.03	0	0.02	0.07
	0.7991	0.8011	0.769	0.8	0.788	0.775	0.7995	0.773	0.773
	0.1818	1.091	1.746	0.1553	1.348	1.969	0.2672	1.478	1.992
$\begin{array}{l} \rho=0.5,\sigma=1\\ \text{Correct}\\ \text{Incorrect}\\ \hat{\rho}\\ \text{MedSE} \end{array}$	$40 \\ 0 \\ 0.4984 \\ 0.1826$	21.61 0 0.5018 0.8458	22.4 0 0.5 0.767	80 0 0.5003 0.1489	43.52 0 0.5 0.867	45.49 0 0.5 0.788	119.99 0 0.5005 0.2289	66.78 0 0.5 0.915	69.73 0 0.5 0.809
$\rho = 0.2, \sigma = 1$ Correct Incorrect $\hat{\rho}$ MedSE	39.99	20.74	21.06	73.99	41.98	42.53	109.99	62.69	63.91
	0	0.65	0.89	0	0.72	1.15	0.99	0.72	1.14
	0.2206	0.3554	0.375	0.2476	0.3644	0.437	0.3424	0.371	0.431
	0.4032	2.9396	3.381	0.8237	3.4479	3.853	2.6921	3.691	3.975
$\begin{array}{l} \rho = 0.8, \sigma = 2\\ \text{Correct}\\ \text{Incorrect}\\ \hat{\rho}\\ \text{MedSE} \end{array}$	38	20.31	20.9	77.98	41.06	42.05	117.99	62.58	62.78
	0	0.02	0.05	0	0.01	0.01	0	0	0.16
	0.7962	0.7944	0.785	0.8002	0.7937	0.778	0.7982	0.793	0.753
	0.4685	1.7795	1.959	0.3963	2.0106	2.263	0.7239	2.123	2.766
$\begin{array}{l} \rho=0.5,\sigma=2\\ \text{Correct}\\ \text{Incorrect}\\ \hat{\rho}\\ \text{MedSE} \end{array}$	38.02	20.51	21.13	76.99	41.76	43.31	118.01	63.37	65.18
	0	0	0	0	0	0	0	0	0
	0.4963	0.5009	0.5	0.5006	0.4987	0.5	0.5008	0.499	0.5
	0.4416	1.6128	1.464	0.3987	1.7193	1.522	0.6623	1.776	1.571
$\begin{array}{l} \rho=0.8,\sigma=2\\ \text{Correct}\\ \text{Incorrect}\\ \hat{\rho}\\ \text{MedSE} \end{array}$	38.99	20.73	20.87	75	41.22	42.04	115	62.43	63.03
	0	0.81	1.22	0	0.59	1.05	0	0.8	1.1
	0.219	0.3583	0.461	0.2383	0.3523	0.434	0.2962	0.413	0.462
	0.5759	3.656	3.804	0.7593	3.6017	3.887	1.8711	4.392	4.039

	i	n = 200, 2q =	10	t	i = 360, 2q =	10	1	n = 500, 2q =	10
	Exp	Square	LAD	Exp	Square	LAD	Exp	Square	LAD
$\rho = 0.8, \sigma = 1, \xi = 0.01$ Correct Incorrect $\hat{\rho}$ MedSE	$10 \\ 0 \\ 0.8016 \\ 0.4001$	5.3 0.23 0.7759 1.8977	5.47 0.01 0.781 0.5016	10 0 0.7997 0.2229	5.05 0.07 0.7978 1.6261	5.45 0 0.7949 0.3415	9.8 0 0.7957 0.2547	5 0.09 0.7606 1.4235	5.5 0 0.7991 0.287
$ \begin{array}{l} \rho = 0.5, \sigma = 1, \xi = 0.01 \\ \text{Correct} \\ \text{Incorrect} \\ \hat{\rho} \\ \text{MedSE} \end{array} $	10 0 0.4999 0.5384	5.23 0.1 0.5003 1.4093	5.53 0 0.5024 0.4443	10 0 0.4997 0.1554	5.11 0.02 0.4967 1.2247	5.68 0 0.4987 0.3282	10 0 0.5004 0.1962	5.03 0.01 0.4981 1.1521	5.74 0 0.4999 0.2811
$\rho = 0.2, \sigma = 1, \xi = 0.01$ Correct Incorrect $\hat{\rho}$ MedSE	9.9 0 0.2096 0.7234	5.15 0.17 0.2569 1.5711	5.43 0.14 0.2543 1.0847	10 0 0.2236 0.4201	5.11 0.02 0.2091 1.1447	5.44 0.03 0.2362 0.8001	9.1 0 0.2513 0.857	5.33 0.01 0.2344 0.9537	5.62 0 0.2358 0.6616
$\rho = 0.8, \sigma = 1, \xi = 0.05$ Correct Incorrect $\hat{\rho}$ MedSE	9 0.2 0.797 0.9265	5.56 0.73 0.7892 4.6649	0.7973 5.34 0 0.4994	10 0 0.7998 0.3901	5.33 0.5 0.7954 3.7479	0.7993 5.43 0 0.3404	8.2 0 0.7857 0.5873	5.23 0.42 0.7892 2.9429	0.7991 5.73 0 0.2881
$\begin{array}{l} \rho = 0.5, \sigma = 1, \xi = 0.05\\ \text{Correct}\\ \text{Incorrect}\\ \hat{\rho}\\ \text{MedSE} \end{array}$	10 0.1 0.4969 0.475	5.31 0.45 0.4999 3.8028	0.5 5.34 0 0.4602	9.8 0 0.4991 0.2339	5.24 0.23 0.4974 2.9176	0.4999 5.87 0 0.332	8,2 0 0.4994 0.4473	5.13 0.18 0.4961 2.4682	0.5002 5.76 0 0.2825
$ ho = 0.8, \sigma = 1, \xi = 0.05$ Correct Incorrect $\hat{ ho}$ MedSE	$10 \\ 0 \\ 0.1366 \\ 0.4858$	5.15 0.25 0.1648 3.1858	0.2815 5.34 0.25 1.0803	8.2 0 0.1762 0.2833	5.06 0.09 0.1159 2.4497	0.2357 5.47 0.02 0.7324	9.1 0 0.2152 0.5134	5.04 0.03 0.1538 2.1076	0.2364 5.34 0 0.6375

Table 7. Estimation with adaptive-l1 regularization when the observations of y have outliers.

Table 8 shows the results of adaptive lasso regularization for data that q = 5, rho = 0.5, and spatial weight matrix has noise. For all test results, the exponential square loss with adaptive L1 identifies more zero coefficients than other models ('Correct'). Compared with the results of the model without regularization term (Table 4), the model with adaptive L1 has a better effect. However, for MedSE, when n = 200, 2q = 10, the exponential square loss with adaptive L1 is the best; when n = 500, 2q = 10, the LAD loss with adaptive L1 is the best; when n = 360, 2q = 10, the LAD loss with adaptive L1 is the best; when n = 360, 2q = 10, the LAD loss with adaptive L1 and the exponential square loss with adaptive L1 have little difference. However, since the exponential square loss with adaptive L1 can identify more nonzero coefficients, we believe that the exponential square loss with adaptive L1 is better than the LAD loss with adaptive L1. The results show that when the spatial weight matrix has estimation error, the SDM with exponential square loss and adaptive lasso has excellent variable selection ability and robustness.

Table 8. Estimation with adaptive-l1 regularization with noisy weighting matrix W.

	i	n = 200, 2q =	10	i	n = 360, 2q =	10	n = 500, 2q = 10		
	Exp	Square	LAD	Exp	Square	LAD	Exp	Square	LAD
$\rho = 0.5, \sigma = 1, \xi = 0.01$ Correct Incorrect $\hat{\rho}$ MedSE	8.1 0 0.4991 1.1925	5.24 0.34 0.4974 2.4503	5.4 0 0.5 0.54	10 0 0.5005 0.1897	5.04 0.41 0.489 2.5552	5.45 0 0.4989 0.3621	10 0 0.5 0.1557	5.11 0.24 0.4835 2.2579	5.54 0 0.5001 0.2916
$\rho = 0.5, \sigma = 1, \xi = 0.03$ Correct Incorrect $\hat{\rho}$ MedSE	9.8 0 0.5012 0.7048	5.54 0.87 0.4644 4.6858	5.4 0 0.4996 0.6433	6.3 1.1 0.4982 1.724	5.61 0.84 0.4665 3.5277	5.3 0 0.4998 0.4495	6.1 1.9 0.4974 3.0656	5.14 0.72 0.476 3.7178	5.76 0 0.4976 0.3789
$\rho = 0.5, \sigma = 1, \xi = 0.05$ Correct Incorrect $\hat{\rho}$ MedSE	9.96 0.04 0.5019 0.832	5.2 1.21 0.4653 4.9396	5.43 0 0.4993 0.806	7.02 0 0.4993 1.3822	5.37 1.2 0.4758 4.3962	5.38 0 0.4953 0.5848	6.02 1 0.4974 2.2125	5.29 1.04 0.4491 4.0998	5.77 0 0.495 0.4734

5. Application of Practical Examples

In this part, we apply the model to actual data to verify the accuracy and efficiency of variable selection and parameter estimation.

We selected a dataset with 211 observations. The dataset describes house sales in the Baltimore area in 1978 and contains home prices and other relevant features. Original data were made available by Robin Dubin [19], Weatherhead School of Management, Case Western Research University, Cleveland, OH. The characteristics of this data are described in Table 9. We mainly study the relationship between price and several other variables. We let the dependent variable be log(PRICE), and the independent variables are NROOM, DWELL, NBATH, PATIO, FIREPL, AC, BMENT, NSTOR, GAR, AGE, CITCOU, LOTSZ, and SQFT.

Variable	Description
STATION	ID variable
PRICE	sales price of house iin \$1000 (MLS)
NROOM	the number of rooms
DWELL	1 if detached unit, 0 otherwise
NBATH	the number of bathrooms
PATIO	1 if patio, 0 otherwise
FIREPL	1 if fireplace, 0 otherwise
AC	1 if air conditioning, 0 otherwise
BMENT	1 if basement, 0 otherwise
NSTOR	number of stories
GAR	number of car spaces in garage $(0 = no garage)$
AGE	age of dwelling in years
CITCOU	1 if dwelling is in Baltimore County, 0 otherwise
LOTSZ	lot size in hundreds of square feet
SQFT	interior living space in hundreds of square feet
X	x coordinate on the Maryland grid
Y	y coordinate on the Maryland grid

 Table 9. Variable description.

We set the spatial weight matrix W by geographic location relationship. The geographic location can be determined by features X and Y. The expression for w_{ij} looks like this:

$$w_{ij} = \frac{1}{(X_i - X_j)^2 + (Y_i - Y_j)^2}.$$
(21)

In addition, we normalize the spatial weights matrix.

Table 10 shows the variable selection results of SDM for square loss, exponential square loss, and LAD loss with adaptive lasso and no penalty. To make variable selection results more intuitive, we designed Table 11. In Table 11, if the model believes that the independent variable has a positive effect on the dependent variable, we mark it as "+"; if the model believes that the independent variable is negatively correlated to the dependent variable, we mark it as "-"; and if the model considers the independent variable not to affect the dependent variable (the absolute value of the parameter estimate is less than 0.001), we do not label it. Additionally, we let the total number of "+" features be **count** "+"; Let the total number of "-" features be **count** "-"; make the total number of all independent variables related to the dependent variable **count**. We can find that the BIC index, with or without regularization, is the lowest exponential square loss. As seen from Tables 10 and 11, our variable selection method has a smaller BIC index than other variable selection methods and selects fewer independent variables, making the model more accurate and more straightforward. This fully illustrates the excellence of the variable selection method proposed in this paper.

Next, we analyze our regression results. For the variable NROOM, the six models all think it positively correlates with the house price, so the more rooms, the higher the house price. For variables DWELL, EXP+adaptive-l1, Square+adaptive-l1, and LAD+adaptive-l1, it is not considered that they will impact house prices, while EXP+null, Square+null, and LAD+null think that they have a specific positive correlation with housing prices. The three models believe that if it is a detached unit, it will make the house price higher. For the variable NBATH, all models believe it positively correlates with the housing price. Therefore, the more bathrooms, the higher the house price. For the variables PATIO and FIREPL, the models with regularization term are considered independent of house price; however, the model without regularization term believes that it is related to house price, the regression coefficients are very small, and the signs of regression coefficients are different. Therefore, we believe that these two characteristics have little impact on house price. For the variable AC, the model, without adding the with no regularization term, thinks that it is positively related to the house price; that is, the house price with air conditioning is higher than that without air conditioning. For the variable BMENT, except for EXP+adaptive-l1, other models believe it positively relates to the house price; that is, houses with basements tend to have higher prices. For the variable CITCOU, the nonregularized model considers that it positively correlates with the house price; in Baltimore, houses in the city will be more expensive. For the spatial autocorrelation coefficient, the six models' estimated values are close to 0.5. It can be seen that the rise in house price will lead to an increase in surrounding house prices. Additionally, we can see that NROOM_W, BMENT_W, and CITCOU_W, under the estimation of the six models, all have negative regression coefficients. Therefore, we can know that the spatial regression coefficients of NROOM, BMENT, and CITCOU are negative. As a result, houses with a lot of rooms, houses with basements, and houses in urban areas can have a negative impact on house prices around them. This is also customary. After all, if all the configurations of a house perform well, people will naturally expect more from the houses around them.

 Table 10. Variable section on real data.

		EXP		Square	LA	AD
	Adaptive-l1	Null	Adaptive-l1	Null	Adaptive-l1	Null
NROOM	0.49674002	0.20409727	0.01051881	0.1929546	0.0037123	0.2362159
DWELL	$^{-1.3922}_{10^{-17}}$	0.45980677	0.00075831	0.4703206	0.00029162	0.5097926
NBATH	0.030063578	0.36030577	0.00385469	0.3514846	0.0012552	0.4254525
PATIO	4.91478×10^{-18}	0.01072777	0.00097357	0.017285	0.00014135	-0.092123
FIREPL	$^{-9.5477}_{10^{-18}} imes$	-0.01059	0.0002972	0.0013726	0.00012271	-0.077913
AC	$^{-1.2919}_{10^{-17}} imes$	0.3021609	0.001	0.311554	0.00020267	0.3138447
BMENT	$^{-2.2645}_{10^{-17}}$	0.1187834	0.0044111	0.1235361	0.00116474	0.1317025
NSTOR	$^{-9.9947}_{10^{-17}}$ $ imes$	0.47809045	0.00359286	0.503778	0.00129079	0.4164672
GAR	$^{-2.2383}_{10^{-17}}$ ×	-0.1040092	0.00038606	-0.099652	0.00039553	-0.0844
AGE	4.72988×10^{-17}	0.01105389	0.03319136	0.0113282	0.02091404	0.0100436
CITCOU	$^{-6.0274}_{10^{-17}}$ $ imes$	0.68202393	0.00093599	0.6868509	0.00025428	0.4701451
LOTSZ	1.04997×10^{-17}	0.0011463	0.002195	0.0011849	0.00443912	4.606×10^{-5}
SQFT	$^{-7.0764}_{10^{-17}}$	-0.0362982	0.0316769	-0.037256	0.01075005	-0.034887
NROOM_W	-0.28507527	-0.1092328	-0.012775	-0.098775	-0.0059391	-0.17384
DWELL_W	1.61355×10^{-33}	-0.1051101	-0.0010685	-0.102124	-0.0005411	-0.107508
NBATH_W	$^{-7.3257 imes}_{10^{-18}}$	-0.1575366	-0.0033218	-0.159341	-0.0014154	-0.095605
PATIO_W	$^{-8.7577}_{10^{-18}} imes$	0.0642994	-8.43×10^{-5}	0.0601125	1.1236×10^{-5}	0.1341514
FIREPL_W	$^{-4.6005}_{10^{-34}}$ $^{\times}$	-0.0387668	0.00068287	-0.036572	$^{-2.587 imes}_{10^{-5}}$	-0.042569
AC_W	$^{-3.3933}_{10^{-17}}$ $ imes$	-0.2098427	-0.001	-0.223255	-0.0004486	-0.187274
BMENT_W	-0.02780421	-0.111448	-0.0074127	-0.117315	-0.0032221	-0.1317

		EXP		Square	L	AD
	Adaptive-l1	Null	Adaptive-l1	Null	Adaptive-l1	Null
NSTOR_W	$^{-2.6009}_{10^{-32}} imes$	-0.1648658	-0.0047893	-0.159493	-0.0023109	-0.178134
GAR_W	1.94949×10^{-17}	0.15495116	0.001	0.1566788	0.00017173	0.1186948
AGE_W	5.87444×10^{-33}	-0.0069188	-0.0228056	-0.007876	-0.0204497	-0.001951
CITCOU_W	-0.06178084	-0.4116914	-0.0030172	-0.426973	-0.001	-0.240381
LOTSZ_W	$^{-2.2196}_{10^{-32}} imes$	-0.0005826	-0.0034412	-0.000541	-0.0057558	-0.000449
SQFT_W	2.75575×10^{-17}	0.01072435	-0.0207167	0.0098465	-0.0112853	0.0143746
ρ MSE BIC	$\begin{array}{r} 0.498613719\\ 0.121911727\\ -304.892336\end{array}$	0.49970041 0.11475312 -317.66083	0.49571237 0.13792606 -278.85061	0.4997043 0.1149259 317.3434	0.49780529 0.14467343 -268.77299	0.499992 0.114829 -317.5214

Table 10. Cont.

Table 11. Visual representation of variable selection on real data.

		EXP		Square	LA	AD
	Adaptive-l1	Null	Adaptive-l1	Null	Adaptive-l1	Null
NROOM	+	+	+	+	+	+
DWELL		+		+		+
NBATH	+	+	+	+	+	+
PATIO		+		+		_
FIREPL		_		+		_
AC		+		+		+
BMENT		+	+	+	+	+
NSTOR		+	+	+	+	+
GAR		_		_		_
AGE		+	+	+	+	+
CITCOU		+		+		+
LOTSZ		+	+	+	+	
SQFT		_	+	_	+	_
NROOM_W	_	_	_	_	_	_
DWELL_W		_	_	_		_
NBATH_W		_	_	_	_	_
PATIO_W		+		+		+
FIREPL_W		_		_		_
AC_W		_		_		_
BMENT_W	_	_	_	_	_	_
NSTOR_W		_	_	_	_	_
GAR_W		+		+		+
AGE_W		_	_	_	_	_
CITCOU_W	-	-	-	_		-
LOTSZ_W			-		-	
SQFT_W		+	_	+	_	+
count "+"	2	13	7	14	7	11
count " $-$ "	3	12	9	11	7	13
count	5	25	16	25	14	24
BIC	-304.892336	-317.66083	-278.85061	-317.3434	-268.77299	-317.5214

6. Conclusions

This paper constructs a robust method for SDM variable selection based on adaptive lasso and exponential square loss. We established the "oracle" nature of the proposed estimators. For the nondifferentiable and nonconvex problems when the model is solved, we design the BCD algorithm, DC decomposition, and CCCP algorithm to solve them. Numerical simulations show that our method has good robustness and accuracy when there is noise in the observed data. Additionally, when the spatial weight matrix estimation is inaccurate, our method also has some robustness. In variable selection, our method is significantly better than exponential squared loss and LAD loss, and almost all zero coefficients can be identified in numerical simulations. Taking the housing price dataset of the Baltimore region in 1978 as an example, the excellence and accuracy of the variable selection method of the SDM proposed in this paper are verified. Our analysis demonstrates the difference between our robust variable selection approach and other penalty regression methods, demonstrating the importance of developing robust variable selection methods.

Author Contributions: Conceptualization, Y.S. and Z.L.; methodology, Z.L.; software, Z.L.; validation, Y.S.; formal analysis, Z.L.; investigation, Y.C.; resources, Z.L.; writing-original draft preparation, Z.L.; writing-review and editing, Z.L., Y.S. and Y.C.; supervision, Y.S.; project administration, Z.L. All authors have read and agreed to the published version of the manuscript.

Funding: The researches are supported by the National Key Research and Development Program of China (2021YFA1000102).

Institutional Review Board Statement: Not applicable.

Data Availability Statement: Not applicable.

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

Appendix A. Proof of Theorem 1

Let $\xi_n = n^{-1/2} + a_n$ and set $\|\mathbf{u}\| = C$, where \mathbf{u} is *d*-dimensional vector and *C* is a large enough constant. Similar to Fan and Li (2001), we first show that $\|\hat{\beta} - \tilde{\beta}_0\| = O_p(\xi_n)$. It suffices to show that, for any given $\epsilon > 0$, there is a large constant *C* such that, for large *n*,

$$P\left\{\sup_{|\mathbf{u}|=C}\ell_n(\boldsymbol{\theta}_0+\boldsymbol{\xi}_n\mathbf{u})<\ell_n(\boldsymbol{\theta}_0)\right\}\geq 1-\epsilon.$$
(A1)

Define $Z = (I - \rho W)^{-1} \tilde{X}^T$, $\epsilon^* = (I - /rhoW)^{-1} \epsilon_n$, and then we can represent model (1) as

$$Y = (I - \rho W)^{-1} \tilde{X}^T \tilde{\beta} + (I - \rho W)^{-1} \varepsilon = Z^T \beta + \varepsilon^*.$$
(A2)

For the optimization model (7)

$$\min_{\widetilde{\beta}\in R^{2p},\rho\in[0,1]}L(\widetilde{\beta},\rho),$$

we know that this is equivalent to

$$\max_{\widetilde{\beta}\in R^{2p},\rho\in[0,1]}-L(\widetilde{\beta},\rho)$$

which can be expressed as

$$\ell_n(\theta) = \sum_{i=1}^n \exp\{-(Y_i - Z_i\tilde{\beta})/\gamma_n\} - n\sum_{j=1}^{2p} p_{\lambda_j}(|\tilde{\beta}_j|).$$
(A3)

Let
$$D_n(\theta, \gamma) = \sum_{i=1}^n \exp\left\{-\left(Y_i - Z_i\tilde{\beta}\right)^2/\gamma\right\} \frac{2(Y_i - Z_i\tilde{\beta})}{\gamma} Z_i$$
. Since $p_{\lambda_j}(0) = 0$ for $j = 1, 2, ..., p$, we have

$$\ell_{n}(\theta_{0} + \xi_{n}\mathbf{u}) - \ell_{n}(\theta_{0}) = \sum_{i=1}^{n} \exp\{-(Y_{i} - Z_{i}(\tilde{\beta}_{0} + \xi_{n}\mathbf{u}))/\gamma_{n}\} - \sum_{i=1}^{n} \exp\{-(Y_{i} - Z_{i}\tilde{\beta}_{0})/\gamma_{n}\} - \sum_{j=1}^{2p} \{p_{\lambda_{j}}(|\tilde{\beta}_{j0} + \xi_{nu_{j}}|) - p_{\lambda_{j}}(|\tilde{\beta}_{j0}|)\}$$

$$\leq \sum_{i=1}^{n} \exp\{-(Y_{i} - Z_{i}(\tilde{\beta}_{0} + \xi_{n}\mathbf{u}))/\gamma_{n}\} - \sum_{i=1}^{n} \exp\{-(Y_{i} - Z_{i}\tilde{\beta}_{0})/\gamma_{n}\} - \sum_{j=1}^{s} \{p_{\lambda_{j}}(|\tilde{\beta}_{j0} + \xi_{nu_{j}}|) - p_{\lambda_{j}}(|\tilde{\beta}_{j0}|)\}$$

$$= S_{n}(\mathbf{u}) + K_{n}(\mathbf{u}).$$
(A4)

Note that

$$S_{n}(\|\mathbf{u}\|) = \sum_{i=1}^{n} \exp\left\{-\frac{(Y_{i} - Z_{i}(\tilde{\beta}_{0} + \xi_{n}u))^{2}}{\gamma_{n}}\right\} - \sum_{i=1}^{n} \exp\left\{-\frac{(Y_{i} - Z_{i}\tilde{\beta}_{0})^{2}}{\gamma_{n}}\right\}$$

$$= \xi_{n} \sum_{i=1}^{n} \left\{\exp\left\{-\frac{(Y_{i} - Z_{i}\tilde{\beta}_{0})^{2}}{\gamma_{n}}\right\} \frac{2(Y_{i} - Z_{i}\tilde{\beta})}{\gamma_{n}} Z_{i}^{T}\right\}^{T} \mathbf{u} - \frac{1}{2} \mathbf{u}^{T} \left\{-\frac{2}{\gamma_{n}} \int ZZ^{T} e^{-(Y - Z\tilde{\beta}_{0})^{2}/\gamma_{n}} \left(A5\right)\right\}$$

$$\times \left(\frac{2(Y - Z\tilde{\beta}_{0})^{2}}{\gamma_{n}} - 1\right) dF(Z, y) \left\{un\xi_{n}^{2} \left\{1 + o_{p}(1)\right\}\right\}$$

$$= \xi_{n} D_{n}(\beta_{0}, \gamma_{n})^{T} \mathbf{u} - \frac{1}{2} \mathbf{u}^{T} \left\{-I(\tilde{\beta}_{0}, \gamma_{n})\right\} un\xi_{n}^{2} \left\{1 + o_{p}(1)\right\}.$$
(A5)

Additionally,

$$K_{n}(\mathbf{u}) = n \sum_{j=0}^{s} \left\{ p_{\lambda_{j}}(|\tilde{\beta}_{j0} + \xi_{n}u_{j}|) - p_{\lambda_{j}}(|\tilde{\beta}_{j0}|) \right\}$$

$$= n\xi_{n} \sum_{j=0}^{s} p_{\lambda_{j}}'(|\tilde{\beta}_{j0}|) \operatorname{sign}(\tilde{\beta}_{j0})u_{j} + n\xi_{n}^{2} \sum_{j=0}^{s} p_{\lambda_{j}}''(|\tilde{\beta}_{j0}|)u_{j}^{2}\{1 + o(1)\}$$

$$\leq a_{n}n\xi_{n} \sum_{j=0}^{s} |u_{j}| + b_{n}n\xi_{n}^{2} \sum_{j=0}^{s} u_{j}^{2}\{1 + o(1)\} \leq a_{n}n\xi_{n} \sum_{j=0}^{s} |u_{j}| + 2b_{n}n\xi_{n}^{2} \|\mathbf{u}\|^{2}$$

$$\leq \sqrt{s}a_{n}n\xi_{n} \sum_{j=0}^{s} |u_{j}| + b_{n}n\xi_{n}^{2} \|\mathbf{u}\|^{2}.$$

(A6)

Since $\gamma_n - \gamma_0 = o_p(1)$, by Taylor's expansion, we have

$$\ell_n(\boldsymbol{\theta}_0 + \boldsymbol{\xi}_n \mathbf{u}) - \ell_n(\boldsymbol{\theta}_0)$$

$$\leq \boldsymbol{\xi}_n D_n(\boldsymbol{\theta}_0, \boldsymbol{\gamma}_n)^T \mathbf{u} - \frac{1}{2} \mathbf{u}^T [-I(\boldsymbol{\theta}_0, \boldsymbol{\gamma}_n)] \mathbf{u} n \boldsymbol{\xi}_n^2 \{1 + o_p(1)\} - \sqrt{s} a_n n \boldsymbol{\xi}_n \sum_{j=0}^s |\boldsymbol{u}_j| + b_n n \boldsymbol{\xi}_n^2 \|\mathbf{u}\|^2.$$
(A7)

Note that $n^{-1/2}D_n(\theta_0, \gamma_0) = O_P(1)$. So, there is $O_p(n^{1/2}\xi_n) = O_p(n\xi_n^2)$ in the last equation of (A.7). By choosing a sufficiently large *C*, the second term dominates the first term uniformly in $\|\mathbf{u}\| = C$. Since $b_n = o_p(1)$, the third term is also dominated by the second term of (A.7). Therefore, (A.1) holds by choosing a sufficiently large *C*.

Appendix B. Proof of Theorem 2

Appendix B.1. Proof of Theorem 2(i)

Here, we show the proof of the first point of Theorem 2. For this, we need only prove that, as $n \to \infty$, there is any beta1 satisfying $\tilde{\beta}_1 - \tilde{\beta}_{01} = O_p(n^{-1/2})$, and for some small $\epsilon_n = \operatorname{Cn}^{-1/2}$ and $j = s + 1, \ldots, p$, we have

$$\frac{\partial \ell_n(\tilde{\beta})}{\partial \tilde{\beta}_j} = \begin{cases} > 0, & \text{for } 0 < \tilde{\beta}_j < \epsilon_n \\ < 0, & \text{for } -\epsilon_n < \tilde{\beta}_j < 0 \end{cases}$$
(A8)

First, let us make

$$Q_n(\tilde{\beta},\gamma) = \sum_{i=1}^n \exp\left\{-\left(Y_i - Z_i^T \tilde{\beta}\right)^2 / \gamma\right\}.$$
 (A9)

Then,

$$\frac{\partial \ell_n(\tilde{\beta})}{\partial \tilde{\beta}_j} = \frac{\partial Q_n(\tilde{\beta}, \gamma_n)}{\partial \tilde{\beta}_j} - n p'_{\lambda_j}(|\tilde{\beta}_j|) \operatorname{sign}(\tilde{\beta}_j).$$
(A10)

By Taylor expansion, we can obtain

$$\frac{\partial \ell_n(\tilde{\beta})}{\partial \tilde{\beta}_j} = \frac{\partial Q_n(\tilde{\beta}_0, \gamma_n)}{\partial \tilde{\beta}_j} + \sum_{l=1}^p \frac{\partial^2 Q_n(\tilde{\beta}_0, \gamma_n)}{\partial \tilde{\beta}_j \partial \tilde{\beta}_l} (\tilde{\beta}_l - \tilde{\beta}_{l0}) + \sum_{l=1}^p \sum_{k=1}^p \frac{\partial^3 Q_n(\tilde{\beta}^*, \gamma_n)}{\partial \tilde{\beta}_j \partial \tilde{\beta}_l \partial \tilde{\beta}_k} (\tilde{\beta}_l - \tilde{\beta}_{l0}) (\tilde{\beta}_k - \tilde{\beta}_{k0})
- np'_{\lambda_j}(|\tilde{\beta}_j|) \operatorname{sign}(\tilde{\beta}_j)
= R_{11} + R_{12} + R_{13} - np'_{\lambda_j}(|\tilde{\beta}_j|) \operatorname{sign}(\tilde{\beta}_j).$$
(A11)

where $\tilde{\beta}^*$ lies between $\tilde{\beta}$ and $\tilde{\beta}_0$. Moreover, because

$$n^{-1} \frac{\partial^2 Q_n(\tilde{\beta}_0, \gamma_0)}{\partial \tilde{\beta}_j \partial \tilde{\beta}_l} = E \left\{ \frac{\partial^2 Q_n(\tilde{\beta}_0)}{\partial \tilde{\beta}_j \partial \tilde{\beta}_l} \right\} + o_p(1),$$
$$n^{-1} \frac{\partial Q_n(\tilde{\beta}_0, \gamma_0)}{\partial \tilde{\beta}_i} = O_p(n^{-1/2}).$$

So there is $R_{11} = O_p(\sqrt{n})$, $R_{12} = O_p(\sqrt{n})$, and $R_{13} = O_p(\sqrt{n})$. Additionally, because of $b_n = o_p(1)$ and $\sqrt{n}a_n = o_p(1)$, we are able to make $\tilde{\beta} - \tilde{\beta}_0 = O_p(n^{-1/2})$.

Since $1/\min_{s+1\leq j\leq d}(\sqrt{n\lambda_j}) = o_p(1)$ and $\lim_{n\to\infty} \inf\lim_{t\to 0+1} \{\min_{s+1\leq j\leq d} p_{\lambda_j}(|t|)/\lambda_j\} > 0$ with probability 1, the sign of the derivative is completely determined by that of β_j . This completes the proof of Theorem 1 (i).

Appendix B.2. Proof of Theorem 2(ii)

Here, we show the proof of the second point of Theorem 2. For brevity, let $\tilde{\beta}_{10}^* = \rho$ and $\tilde{\beta}_{1j}^* = \tilde{\beta}_{1j}, j = 1, ..., s$, then denote $\tilde{\beta}_1^* = (\rho, \tilde{\beta}_{11}, ..., \tilde{\beta}_{1s})^T$ and $\tilde{\beta}_0^* = (\rho_0, \tilde{\beta}_{10}, ..., \tilde{\beta}_{0s})^T$. We known that $\hat{\theta}$ minimizes $Q_n(\theta)$. We showed that there exists a \sqrt{n} -consistent local maximizer of $\ell_n \{ (\tilde{\beta}_1, 0) \}$. satisfying that

$$\frac{\partial \ell_n \left\{ \left(\hat{\beta}_1, 0 \right) \right\}}{\partial \tilde{\beta}_j} = 0, \quad \text{for } j = 1, \dots, s.$$

Since $\tilde{\beta}_1$ is a consistent estimator, we have

$$\frac{\partial Q_n \left\{ \left(\hat{\beta}_1, 0 \right), \gamma_n \right\}}{\partial \beta_j} - n p'_{\lambda_j} \left(\left| \hat{\beta}_j \right| \right) \operatorname{sign} \left(\hat{\beta}_j \right) \\
= \frac{\partial Q_n \left(\tilde{\beta}_0, \gamma_n \right)}{\partial \tilde{\beta}_j} + \sum_{l=1}^s \left\{ \frac{\partial^2 Q_n \left(\tilde{\beta}_0, \gamma_n \right)}{\partial \tilde{\beta}_j \partial \tilde{\beta}_l} + o_p(1) \right\} \left(\hat{\beta}_l - \tilde{\beta}_{01} \right) \\
- n \left[p'_{\lambda_j} \left(\left| \tilde{\beta}_{0j} \right| \right) \operatorname{sign} \left(\tilde{\beta}_{0j} \right) + \left\{ p''_{\lambda_j} \left(\left| \beta_{0j} \right| \right) + o_p(1) \right\} \left(\hat{\beta}_j - \tilde{\beta}_{0j} \right) \right] = 0.$$
(A12)

The above equation can be rewritten as follows:

$$\frac{\partial Q_n(\tilde{\beta}_0,\gamma_n)}{\partial \tilde{\beta}_j} = \sum_{l=1}^s \left\{ E\left\{ \frac{\partial^2 Q_n(\tilde{\beta}_0,\gamma_n)}{\partial \tilde{\beta}_j \partial \tilde{\beta}_l} \right\} + o_p(1) \right\} \left(\hat{\tilde{\beta}}_l - \tilde{\beta}_{01}\right) + n\Delta + n \left[\Sigma_1 + O_p(1)\right] \left(\hat{\tilde{\beta}}_{n1} - \tilde{\beta}_{01}\right), \tag{A13}$$

$$nI_{1}(\tilde{\beta}_{01},\gamma_{0})(\hat{\beta}_{n1}-\tilde{\beta}_{01})+n\Delta+n[\Sigma_{1}+O_{p}(1)](\hat{\beta}_{n1}-\tilde{\beta}_{01})$$

$$=n[I_{1}(\tilde{\beta}_{01},\gamma_{0})+\Sigma_{1}](\hat{\beta}_{n1}-\tilde{\beta}_{01})+n\Delta$$

$$=n[I_{1}(\tilde{\beta}_{01},\gamma_{0})+\Sigma_{1}](\hat{\beta}_{n1}-\tilde{\beta}_{01})+n[I_{1}(\tilde{\beta}_{01},\gamma_{0})+\Sigma_{1}][I_{1}(\tilde{\beta}_{01},\gamma_{0})+\Sigma_{1}]^{-1}\Delta \quad (A14)$$

$$=n[I_{1}(\tilde{\beta}_{01},\gamma_{0})+\Sigma_{1}]\{(\hat{\beta}_{n1}-\tilde{\beta}_{01})+n[I_{1}(\tilde{\beta}_{01},\gamma_{0})+\Sigma_{1}]^{-1}\Delta\}$$

$$=-\frac{\partial Q_{n}(\tilde{\beta}_{0},\gamma_{n})}{\partial \tilde{\beta}_{j}}+o_{p}(1).$$

Since $\sqrt{n}(\gamma_n - \gamma_0) = o_p(1)$, invoking Slutsky's lemma and the Lindeberg–Feller central limit theorem, we have

$$\begin{split} \sqrt{n} \big(I_1(\tilde{\beta}_{01},\gamma_0) + \Sigma_1 \big) \Big\{ \Big(\hat{\tilde{\beta}}_{n1} - \tilde{\beta}_{01} \Big) + \big(I_1(\tilde{\beta}_{01},\gamma_0) + \Sigma_1 \big)^{-1} \Delta \Big\} &\to N(\mathbf{0},\Sigma_2), \\ where \hat{\tilde{\beta}}_{n1} = \Big(\hat{\rho}, \hat{\tilde{\beta}}_{11}, \dots, \hat{\tilde{\beta}}_{15} \Big)^T, and \tilde{\beta}_{01} = \big(\rho_0, \tilde{\beta}_{01}, \dots, \tilde{\beta}_{0s} \big)^T, \\ \Sigma_1 = \operatorname{diag} \Big\{ p_{\lambda_1}''(|\tilde{\beta}_{01}|), \dots, p_{\lambda_s}''(|\tilde{\beta}_{0s}|) \Big\}, \Sigma_2 = \operatorname{cov} \Big(\exp \Big(-r^2/\gamma_0 \Big) \frac{2r}{\gamma_0} Z_{i1} \Big), \\ \Delta = \Big(p_{\lambda_1}'(|\tilde{\beta}_{01}|) \operatorname{sign}(\tilde{\beta}_{01}), \dots, p_{\lambda_s}'(|\tilde{\beta}_{0s}|) \times \operatorname{sign}(\tilde{\beta}_{0s}) \Big)^T, \\ I_1(\tilde{\beta}_{01}, \gamma_0) = \frac{2}{\gamma_0} E \Big[\exp \Big(-r^2/\gamma_0 \Big) \Big(\frac{2r^2}{\gamma_0} - 1 \Big) \Big] \times \Big(E Z_{i1} Z_{i1}^T \Big). \end{split}$$

Then, the proof of part (ii) is completed.

References

- 1. Anselin, L. Spatial Econometrics: Methods and Models; Springer Science & Business Media: Berlin/Heidelberg, Germany, 1988.
- Kelejian, H.H. A spatial J-test for model specification against a single or a set of non-nested alternatives. *Lett. Spat. Resour. Sci.* 2008, 1, 3–11. [CrossRef]
- Zhang, X.; Yu, J. Spatial weights matrix selection and model averaging for spatial autoregressive models. J. Econom. 2018, 203, 1–18 [CrossRef]
- 4. Tibshirani, R. Regression shrinkage and selection via the lasso. J. R. Stat. Soc. Ser. Methodol. 1996, 58, 267–288. [CrossRef]
- 5. Fan, J.; Li, R.. Variable selection via nonconcave penalized likelihood and its oracle properties. J. Am. Stat. Assoc. 2001, 96, 1348–1360. [CrossRef]
- 6. Zou, H. The adaptive lasso and its oracle properties. J. Am. Stat. Assoc. 2006, 101, 1418–1429. [CrossRef]
- Wang, X.; Jiang, Y.; Huang, M.; Zhang, H.: Robust variable selection with exponential squared loss. J. Am. Stat. Assoc. 2013, 108, 632–643. [CrossRef] [PubMed]
- Friedman, J.; Hastie, T.; Tibshirani, R. Additive logistic regression: a statistical view of boosting (with discussion and a rejoinder by the authors). Ann. Stat. 2000, 28, 337–407. [CrossRef]
- 9. Koenker, Roger, Bassett, Gilbert. Regression quantiles. *Econometrica* **1978**, *46*, 33–50. [CrossRef]
- 10. Zou, H.; Yuan, M. Composite quantile regression and the oracle model selection theory. Ann. Stat. 2008, 36, 1108–1126. [CrossRef]
- 11. Beer, C.; Riedl, A.. Modelling spatial externalities in panel data: The Spatial Durbin model revisited. *Pap. Reg. Sci.* 2012, *91*, 299–318. [CrossRef]
- 12. Mustaqim, Setiawan, Suhartono, Ulama, B.S.S. Efficient estimation of simultaneous equations of spatial Durbin panel data model. In *Proceedings of the AIP Conference Proceedings*; AIP Publishing LLC: Melville, NY, USA, 2018; Volume 2021, p. 060024.
- 13. Zhu, Y.; Han, X.; Chen, Y. Bayesian estimation and model selection of threshold spatial Durbin model. *Econom. Lett.* 2020, 188, 108956. [CrossRef]
- 14. Wei, L.; Zhang, C.; Su, J.J.; Yang, L. Lixiong Panel threshold spatial Durbin models with individual fixed effects. *Econom. Lett.* **2021**, 201, 109778. [CrossRef]
- Song, Y.; Liang, X.; Zhu, Y.; Lin, L. Robust variable selection with exponential squared loss for the spatial autoregressive model. *Comput. Stat. Data Anal.* 2021, 155, 107094. [CrossRef]
- 16. Wang, H., Li, G.; Tsai, C.L. Regression coefficient and autoregressive order shrinkage and selection via the lasso. J. R. Stat. Soc. Ser. (Stat. Methodol. 2007, 69, 63–78. [CrossRef]
- 17. Forsythe, George Elmer, Moler, Cleve B., Malcolm, Michael A. *Computer Methods for Mathematical Computations*; Prentice Hall: Hoboken, NJ, USA, 1977.

- 18. Beck, A.; Teboulle, M. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM J. Imaging Sci.* 2009, 2, 183–202. [CrossRef]
- 19. Dubin, R.A. Spatial autocorrelation and neighborhood quality. Reg. Sci. Urban Econ. 1992, 22, 433–452. [CrossRef]

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.