

# On a Semiparametric Data-driven Nonlinear Model with Penalized Spatio-temporal Lag Interactions

Dawlah Al-Sulami<sup>a</sup>, Zhenyu Jiang<sup>b</sup>, Zudi Lu<sup>c</sup>, Jun Zhu<sup>d</sup>

<sup>a</sup>*Department of Statistics, King Abdulaziz University*

<sup>b</sup>*Statistical Sciences Research Institute, University of Southampton*

<sup>c</sup>*Statistical Sciences Research Institute and School of Mathematical Sciences, University of Southampton*

<sup>d</sup>*Department of Statistics and Department of Entomology, University of Wisconsin-Madison*

---

## Abstract

To study possibly nonlinear relationship between housing price index and consumer price index for individual states in the US, accounting for the temporal lag interactions of the housing price in a given state and spatio-temporal lag interactions between states could improve the accuracy of estimation and forecasting. There lacks, however, methodology to objectively identify and estimate such spatio-temporal lag interactions. In this paper, we propose a semiparametric data-driven nonlinear time series regression method that accounts for lag interactions across space and over time. A penalized procedure utilizing adaptive Lasso is developed for the identification and estimation of important spatio-temporal lag interactions. Theoretical properties for our proposed methodology are established under a general near epoch dependence structure and thus the results can be applied to a variety of linear and nonlinear time series processes. For illustration, we analyze the US housing price data and demonstrate substantial improvement in forecasting via the identification of nonlinear relationship between housing price index and consumer price index as well as spatio-temporal lag interactions.

**Keywords:** adaptive Lasso, asymptotic property, nonlinear regression, regularization, spatial statistics, time series

---

## 1. Introduction

Of broad practical interest are nonlinear relationships between covariate variables and a response variable with spatio-temporal effects for not necessarily Gaussian data. For example, for a

given state in the United States (US), ignoring the temporal lag interactions of the housing price return in this state and between states may result in biased estimates of possibly nonlinear effect of the consumer price index (CPI) change on the housing price index (HPI) return (see Section 5 for more detail). Such spatio-temporal lag interactions in a nonlinear regression setting with irregularly spaced time series data are not well understood in the literature (Al-Sulami et al., 2017). The purpose of this paper is to develop a semiparametric data-driven nonlinear time series regression method that objectively selects spatio-temporal lag interactions based on data.

It is well known that too many unnecessary predictors may give inefficient estimation and prediction and therefore, selecting the more important predictors among a large number of predictors is of keen interest. Variable selection procedures such as stepwise selection, Akaike’s information criterion (AIC) (Akaike, 1973) and Bayesian information criterion (BIC) (Schwarz, 1978) are extensively employed to choose the appropriate covariates in linear regression as well as the appropriate lag order in time series analysis (see, e.g., Ramanathan, 1992; McQuarrie and Tsai, 1998; Shumway and Stoffer, 2000). These traditional methods have a number of drawbacks including instability (Breiman, 1996), as the estimation and variable selection are executed in separate steps, and the stochastic error is not taken into account in the variable selection step (Fan and Li, 2001).

To overcome the limitations of the traditional variable selection methods and enhance the prediction accuracy, a variety of penalized methods have been developed for linear regression and gained popularity. See, for example, bridge regression (Frank and Friedman, 1993), least absolute shrinkage and selection operator (Lasso) (Tibshirani, 1996), and elastic net (Zou and Hastie, 2005). By applying an  $\mathcal{L}_1$ -penalty function with a single regularization parameter, Lasso produces a sparse model and yields a consistent estimator under appropriate conditions. The necessary conditions to achieve the consistency were explored by Zou (2006) and Zhao and Yu (2007). Fan and Li (2001) proposed a penalized likelihood approach with a set of penalty functions including Lasso as a special case, and showed that penalized likelihood with a smoothly clipped absolute deviation (SCAD) penalty function enhances the performance in model selection. Zou (2006) developed

adaptive Lasso by using an  $\mathcal{L}_1$ -penalty function which assigns different weights to different interactions, and demonstrated that adaptive Lasso enjoys *oracle* properties and ease of applications. Variable selection techniques have also been developed for generalized linear models, survival data, and time series data (see, e.g., Tibshirani et al., 1997; Wang et al., 2007; Van De Geer, 2008; Hsu et al., 2008; Haufe et al., 2009; Ren and Zhang, 2010).

For spatial lattice data, Zhu et al. (2010) applied adaptive Lasso to select covariates and spatial neighborhoods in a spatial linear regression model (see also, Huang et al., 2010). Reyes et al. (2012) further developed adaptive Lasso for the selection of covariates and spatio-temporal coefficients. However, the spatial and spatio-temporal neighborhood structures considered in Zhu et al. (2010) and Reyes et al. (2012) are relatively simple (see also Hallin et al., 2004; Gao et al., 2006; Lu et al., 2007, 2009). Bayesian model selection has also been considered for spatial data such as Song and De Oliveira (2012) and Hefley et al. (2017).

In spatial econometrics, a spatial weight matrix measures possibly complex spatial interactions between spatial locations or spatial units on a lattice and is crucial to construct when applying a spatial or spatio-temporal model (Anselin, 1988). The influence of the construction of a spatial weight matrix on both model testing and parameter estimation has been documented by many (see, e.g., Stetzer, 1982; Griffith and Lagona, 1998; Smith, 2009; Stakhovych and Bijmolt, 2009). In many applications, however, the spatial weight matrix is assumed to be known *a priori*, which may or may not reflect the true underlying spatial interaction. Thus considerable attention has been devoted to the estimation of the spatial weight matrix in various settings with low- or high-dimensional spatial panels in recent years. For example, estimation methods were proposed for the spatial weight matrix in low-dimensional spatial context under different structural constraints (Bhattacharjee and Holly, 2011; Bhattacharjee and Jensen-Butler, 2013; Bhattacharjee and Holly, 2013). de Souza (2012) proposed a spatial autoregressive model with exogenous covariates and estimated the spatial weight matrix by Lasso, but required *a priori* knowledge about the structure of the spatial interactions. Lam and Souza (2013) developed a spatial lag model with exogenous

covariates and estimated the spatial weight matrix together with the model parameters via adaptive Lasso, whereas Manresa (2013) considered a non-autoregressive spatial model with exogenous covariates and estimated the spatial weight matrix by a pooled Lasso technique. Ahrens and Bhat-tacharjee (2015) proposed a spatial autoregressive model with exogenous covariates and developed a two-step Lasso method to estimate the spatial weight matrix, which was demonstrated to be effective to uncover the underlying spatial dependence structure via a simulation study. All of these methods above in spatial econometrics, however, focused on the estimation of the spatial weight matrix for linear models under spatial stationarity.

In contrast to the research development above, here we construct and estimate the spatio-temporal weights for nonlinear models with non-Gaussian errors. In our previous work (Al-Sulami et al., 2017), we considered a partially nonlinear regression model for spatial time series. However, the neighborhood structure was pre-specified and thus can be quite subjective. Here, we develop a penalized method to simultaneously identify and estimate the spatio-temporal lag interactions in the setting of a semiparametric data-driven nonlinear model (see, e.g., Zhang et al., 2003; Lu et al., 2008, 2009). Although the theory for nonparametric or nonlinear estimation is well established, here we consider the adaptive Lasso for selecting lagged variables in time series and space-time model under a general near epoch dependence structure, which can be applied to a wide range of linear and nonlinear time series processes (c.f., Lu and Linton, 2007; Li et al., 2012).

The rest of this paper is organized as follows. In Sections 2 and 3, we describe our proposed model and estimation procedure. In Section 4, the asymptotic properties for the estimation procedure are established. For illustration, we apply our method to explore possible nonlinear relationship between housing price and consumer price in selected states of the US in Section 5. Concluding remarks are given in Section 6. Technical details including proofs of theorems are relegated to Appendixes as web-based supplementary materials.

## 2. Model

Let  $Y_t(\mathbf{s}_0)$  and  $X_t(\mathbf{s}_0)$  denote two time series which are observed at discrete time points  $t = 1, \dots, T$ , at a given spatial location  $\mathbf{s}_0 = (u_0, v_0) \in \mathbb{R}^2$ , where  $u_0$  and  $v_0$  are the  $x$  and  $y$  coordinates, respectively, representing a spatial unit on a lattice. In the US housing price data example,  $u_0$  and  $v_0$  are the latitude and longitude of the centroid of a given state. Further,  $Y_t(\mathbf{s}_0)$  is a univariate response variable and  $X_t(\mathbf{s}_0)$  is the covariate vector of dimension  $d$ , where  $X_t$  can be the same for different  $\mathbf{s}_0$  while  $d$  should be small to avoid the curse of dimensionality. A possibly nonlinear relationship between the response  $Y_t(\mathbf{s}_0)$  and the covariate  $X_t(\mathbf{s}_0)$  is of interest.

For a given spatial unit represented by the spatial location  $\mathbf{s}_0$ , we may have  $N$  other spatial units represented by spatial locations on a possibly irregular lattice, say  $\mathbf{s}_k := (u_k, v_k) \in \mathbb{R}^2$  for  $k = 1, \dots, N$ , over which  $Y_t(\mathbf{s}_k)$  are observed at the  $T$  time points  $t = 1, \dots, T$  affecting  $Y_t(\mathbf{s}_0)$ . At a given spatial location  $\mathbf{s}_0$ , we consider a semiparametric nonlinear regression time series model written as,

$$Y_t(\mathbf{s}_0) = g_0(X_t(\mathbf{s}_0)) + \sum_{i=1}^p \sum_{k=1}^N \lambda_{0k,i} Y_{t-i}(\mathbf{s}_k) + \sum_{l=1}^q \alpha_{0,l} Y_{t-l}(\mathbf{s}_0) + \varepsilon_t(\mathbf{s}_0) \quad (1)$$

where  $t = r + 1, \dots, T$  with  $r = \max(p, q)$  and  $g_0(\cdot)$  is an unknown function characterizing the relationship between  $X_t(\mathbf{s}_0)$  and  $Y_t(\mathbf{s}_0)$ . Further,  $\lambda_{0k,i}$  are unknown coefficients representing the spatio-temporal lag interactions of orders  $p$  between spatial units, while  $\alpha_{0,l}$  are unknown coefficients representing the temporal lag interactions of order  $q$  for the given spatial unit  $\mathbf{s}_0$ . Finally,  $\varepsilon_t(\mathbf{s}_0)$  are i.i.d. errors, not necessarily Gaussian, with mean  $E\varepsilon_t(\mathbf{s}_0) = 0$  and variance  $E\varepsilon_t^2(\mathbf{s}_0) = \sigma_0^2$ .

Model (1) can be viewed as an extension of the model in Al-Sulami et al. (2017) at a given spatial location  $\mathbf{s}_0$ . Unlike Al-Sulami et al. (2017), however, here the lag effects between two spatial units  $\lambda_{0k,i}$  are assumed to be unknown and depend on the time lag  $i$  in the characterization of spatio-temporal lag interactions. In fact, when  $\lambda_{0k,i} = \lambda_i(\mathbf{s}_0)w_{0k}$  with unknown  $\lambda_i(\mathbf{s}_0)$  but

known  $w_{0k}$  such that  $\sum_{k=1}^N w_{0k} = 1$ , our model (1) reduces to the partially nonlinear model of Al-Sulami et al. (2017) at the spatial location  $\mathbf{s}_0$ .

### 3. Estimation

How well the unknown function  $g_0(\cdot)$  is estimated relies on the selection and estimation of the unknown spatio-temporal lag interactions  $\lambda_{0k,i}$  and  $\alpha_{0,l}$ . A challenging issue here is that the unknown vector  $\boldsymbol{\lambda}(\mathbf{s}_0)$ , consisting of all the spatio-temporal lag interactions  $\lambda_{0k,i}$ , has  $Np$  unknown components, which can be quite large. For example, in the US housing price data example, we have  $N = 50$ . With say  $p = 6$  time lags, the dimension of  $\boldsymbol{\lambda}(\mathbf{s}_0)$  is 300. In addition, for different time lags  $i$ , the number of non-zero interactions (i.e., number of  $\mathbf{s}_k$ 's with  $\lambda_{0k,i} \neq 0$ ) for  $\mathbf{s}_0$  can be different. Thus, the estimation of model (1) is more challenging than that of Al-Sulami et al. (2017) with pre-specified spatial weights at a fixed location.

Before giving the estimation detail, we introduce some notation. At a given spatial location  $\mathbf{s}_0$ , let  $\boldsymbol{\eta}(\mathbf{s}_0) = (\boldsymbol{\lambda}(\mathbf{s}_0)', \boldsymbol{\alpha}(\mathbf{s}_0)')'$ , with  $\boldsymbol{\lambda}(\mathbf{s}_0) = (\boldsymbol{\lambda}_1(\mathbf{s}_0)', \boldsymbol{\lambda}_2(\mathbf{s}_0)', \dots, \boldsymbol{\lambda}_N(\mathbf{s}_0)')'$ ,  $\boldsymbol{\lambda}_k(\mathbf{s}_0) = (\lambda_{0k,1}, \lambda_{0k,2}, \dots, \lambda_{0k,p})'$ , and  $\boldsymbol{\alpha}(\mathbf{s}_0) = (\alpha_{0,1}, \alpha_{0,2}, \dots, \alpha_{0,q})'$ . Moreover, let  $\mathbf{Z}_t(\mathbf{s}_0) = (\mathbf{Z}_{1,t}(\mathbf{s}_0)', \mathbf{Z}_{2,t}(\mathbf{s}_0)')'$ , where  $\mathbf{Z}_{1,t}(\mathbf{s}_0) = (\mathbf{Y}_{t-1}(\mathbf{s}_1)', \mathbf{Y}_{t-1}(\mathbf{s}_2)', \dots, \mathbf{Y}_{t-1}(\mathbf{s}_N)')'$ , with  $\mathbf{Y}_{t-1}(\mathbf{s}_k) = (Y_{t-1}(\mathbf{s}_k), Y_{t-2}(\mathbf{s}_k), \dots, Y_{t-p}(\mathbf{s}_k))'$ , and  $\mathbf{Z}_{2,t}(\mathbf{s}_0) = (Y_{t-1}(\mathbf{s}_0), Y_{t-2}(\mathbf{s}_0), \dots, Y_{t-q}(\mathbf{s}_0))'$ . Then model (1) can be written more succinctly as

$$Y_t(\mathbf{s}_0) = g_0(X_t(\mathbf{s}_0)) + \mathbf{Z}_{1,t}(\mathbf{s}_0)' \boldsymbol{\lambda}(\mathbf{s}_0) + \mathbf{Z}_{2,t}(\mathbf{s}_0)' \boldsymbol{\alpha}(\mathbf{s}_0) + \varepsilon_t(\mathbf{s}_0). \quad (2)$$

By taking conditional expectation of the terms in (2), we have

$$\begin{aligned} g_0(X_t(\mathbf{s}_0)) &= E[Y_t(\mathbf{s}_0)|X_t(\mathbf{s}_0)] - E[\mathbf{Z}_{1,t}(\mathbf{s}_0)|X_t(\mathbf{s}_0)]' \boldsymbol{\lambda}(\mathbf{s}_0) - E[\mathbf{Z}_{2,t}(\mathbf{s}_0)|X_t(\mathbf{s}_0)]' \boldsymbol{\alpha}(\mathbf{s}_0) \\ &\equiv g_1(X_t(\mathbf{s}_0), \mathbf{s}_0) - \mathbf{g}_{21}(X_t(\mathbf{s}_0), \mathbf{s}_0)' \boldsymbol{\lambda}(\mathbf{s}_0) - \mathbf{g}_{22}(X_t(\mathbf{s}_0), \mathbf{s}_0)' \boldsymbol{\alpha}(\mathbf{s}_0), \end{aligned} \quad (3)$$

where  $g_1(X_t(\mathbf{s}_0), \mathbf{s}_0) \in \mathbb{R}^1$ ,  $\mathbf{g}_{21}(X_t(\mathbf{s}_0), \mathbf{s}_0) \in \mathbb{R}^{Np}$  and  $\mathbf{g}_{22}(X_t(\mathbf{s}_0), \mathbf{s}_0) \in \mathbb{R}^q$  denote the condi-

121 tional means of  $E[Y_t(\mathbf{s}_0)|X_t(\mathbf{s}_0)]$ ,  $E[\mathbf{Z}_{1,t}(\mathbf{s}_0)|X_t(\mathbf{s}_0)]$  and  $E[\mathbf{Z}_{2,t}(\mathbf{s}_0)|X_t(\mathbf{s}_0)]$ , respectively, at the  
 122 spatial location  $\mathbf{s}_0$ .

### 123 3.1. Estimation of $g_0(X_t(\mathbf{s}_0))$

To estimate  $g_0(X_t(\mathbf{s}_0))$ , we first consider estimating the three conditional means in (3),  $g_1(x, \mathbf{s}_0)$ ,  $\mathbf{g}_{21}(x, \mathbf{s}_0)$ , and  $\mathbf{g}_{22}(x, \mathbf{s}_0)$  at  $X_t(\mathbf{s}_0) = x$ . Define  $a_0 = a_0(x, \mathbf{s}_0) = g_1(x, \mathbf{s}_0)$  and let  $a_1 = a_1(x, \mathbf{s}_0) = \dot{g}_1(x, \mathbf{s}_0)$  denote the first-order derivative of  $g_1$  with respect to  $x$ . By local linear fitting (Fan and Gijbels, 1996; Hallin et al., 2004), we obtain the estimators  $\hat{a}_0 = \hat{g}_1(x, \mathbf{s}_0)$  and  $\hat{a}_1$  by

$$(\hat{a}_0, \hat{a}_1)' = \underset{(a_0, a_1)' \in \mathbb{R}^2}{\operatorname{argmin}} \sum_{t=r+1}^T \{Y_t(\mathbf{s}_0) - a_0 - a_1(X_t(\mathbf{s}_0) - x)\}^2 K_b(X_t(\mathbf{s}_0) - x), \quad (4)$$

124 where  $T_0 = T - r$  is the effective sample size with  $r = \max\{p, q\}$ ,  $b = b_{T_0}$  is a bandwidth tending  
 125 to zero as  $T \rightarrow \infty$ ,  $K(\cdot)$  is a bounded kernel function, and  $K_b(\cdot) = b^{-1}K(\cdot/b)$ .

Let  $\mathbf{A}(x)$  be a  $T_0 \times 2$  matrix with the  $(t - r)$ th-row  $(1, b^{-1}(X_t(\mathbf{s}_0) - x))$  for  $t = r + 1, \dots, T$ , and let  $\mathbf{B}(x)$  be a  $T_0 \times T_0$  diagonal matrix with the  $t$ th diagonal element  $K_b(X_t(\mathbf{s}_0) - x)$  for  $t = r + 1, \dots, T$ . Let  $\mathbf{Y} = (Y_{r+1}(\mathbf{s}_0), \dots, Y_T(\mathbf{s}_0))'$  denote a  $T_0$ -dimensional vector of responses. Then the local linear estimator is given by

$$(\hat{a}_0, b\hat{a}_1)' = \mathbf{U}_{T_0}^{-1} \mathbf{V}_{T_0},$$

126 where  $\mathbf{U}_{T_0} = \mathbf{A}(x)' \mathbf{B}(x) \mathbf{A}(x) = \begin{pmatrix} u_{T_0,00} & u_{T_0,01} \\ u_{T_0,10} & u_{T_0,11} \end{pmatrix}$  is a  $2 \times 2$  matrix and  $\mathbf{V}_{T_0} = \mathbf{A}(x)' \mathbf{B}(x) \mathbf{Y} =$   
 127  $(v_{T_0,0}, v_{T_0,1})'$  is a  $2 \times 1$  vector. By denoting  $\left(\frac{X_t(\mathbf{s}_0) - x}{b}\right)^0 = 1$ , we have

$$\begin{aligned} u_{T_0,jk} &= (T_0 b)^{-1} \sum_{t=r+1}^T \left(\frac{X_t(\mathbf{s}_0) - x}{b}\right)^j \left(\frac{X_t(\mathbf{s}_0) - x}{b}\right)^k K\left(\frac{X_t(\mathbf{s}_0) - x}{b}\right), \quad j, k = 0, 1 \\ v_{T_0,j} &= (T_0 b)^{-1} \sum_{t=r+1}^T Y_t(\mathbf{s}_0) \left(\frac{X_t(\mathbf{s}_0) - x}{b}\right)^j K\left(\frac{X_t(\mathbf{s}_0) - x}{b}\right), \quad j = 0, 1. \end{aligned}$$

Thus, with  $\mathbf{e}_1 = (1, 0)' \in \mathbb{R}^2$ , the local linear estimator of  $g_1(x, \mathbf{s}_0)$  is

$$\hat{g}_1(x, \mathbf{s}_0) = \hat{a}_0 = \mathbf{e}_1' \mathbf{U}_{T_0}^{-1} \mathbf{V}_{T_0}. \quad (5)$$

Similarly, we estimate  $\mathbf{g}_{21}(X_t(\mathbf{s}_0), \mathbf{s}_0) \in \mathbb{R}^{Np}$  with

$$\mathbf{g}_{21}(X_t(\mathbf{s}_0), \mathbf{s}_0) = \left( (g_{21}^{1,k}(X_t(\mathbf{s}_0), \mathbf{s}_0), \dots, g_{21}^{p,k}(X_t(\mathbf{s}_0), \mathbf{s}_0))' : k = 1, \dots, N \right),$$

where  $g_{21}^{i,k}(X_t(\mathbf{s}_0), \mathbf{s}_0) = E[Y_{t-i}(\mathbf{s}_k) | X_t(\mathbf{s}_0)]$ . Let  $\mathbf{Z}_1^{i,k} = (Y_{(r+1)-i}(\mathbf{s}_k), \dots, Y_{T-i}(\mathbf{s}_k))'$  be a  $T_0 \times 1$  vector and  $\mathbf{R}_{1T_0}^{i,k} = \mathbf{A}(x)' \mathbf{B}(x) \mathbf{Z}_1^{i,k} = (\mathbf{r}_{1T_0,0}^{i,k}, \mathbf{r}_{1T_0,1}^{i,k})'$  a  $2 \times 1$  vector, where

$$\mathbf{r}_{1T_0,j}^{i,k} = (T_0 b)^{-1} \sum_{t=r+1}^T Y_{t-i}(\mathbf{s}_k) \left( \frac{X_t(\mathbf{s}_0) - x}{b} \right)^j K \left( \frac{X_t(\mathbf{s}_0) - x}{b} \right); j = 0, 1.$$

Then the local linear estimator of  $g_{21}^{i,k}(x, \mathbf{s}_0)$  is

$$\hat{g}_{21}^{i,k}(x, \mathbf{s}_0) = \mathbf{e}_1' \mathbf{U}_{T_0}^{-1} \mathbf{R}_{1T_0}^{i,k}. \quad (6)$$

Therefore, the local linear estimator of the unknown vector  $\mathbf{g}_{21}(x, \mathbf{s}_0)$  is given by

$$\hat{\mathbf{g}}_{21}(x, \mathbf{s}_0) = \mathbf{e}_1' \left( \mathbf{U}_{T_0}^{-1} \mathbf{R}_{1T_0}^{1,1}, \dots, \mathbf{U}_{T_0}^{-1} \mathbf{R}_{1T_0}^{p,1}, \mathbf{U}_{T_0}^{-1} \mathbf{R}_{1T_0}^{1,2}, \dots, \mathbf{U}_{T_0}^{-1} \mathbf{R}_{1T_0}^{p,2}, \dots, \mathbf{U}_{T_0}^{-1} \mathbf{R}_{1T_0}^{1,N}, \dots, \mathbf{U}_{T_0}^{-1} \mathbf{R}_{1T_0}^{p,N} \right)'.$$

Further, we estimate  $\mathbf{g}_{22}(X_t(\mathbf{s}_0), \mathbf{s}_0) = (g_{22}^1(X_t(\mathbf{s}_0), \mathbf{s}_0), \dots, g_{22}^q(X_t(\mathbf{s}_0), \mathbf{s}_0))'$  with

$$g_{22}^l(X_t(\mathbf{s}_0), \mathbf{s}_0) = E[Y_{t-l}(\mathbf{s}_0) | X_t(\mathbf{s}_0)],$$

for  $l = 1, \dots, q$ . Let  $\mathbf{Z}_2^l = (Y_{(r+1)-l}(\mathbf{s}_0), \dots, Y_{T-l}(\mathbf{s}_0))'$  be a  $T_0 \times 1$  vector and let  $\mathbf{R}_{2T_0}^l =$



$\mathbf{A}(x)' \mathbf{B}(x) \mathbf{Z}_2^l = (\mathbf{r}_{2T_0,0}^l, \mathbf{r}_{2T_0,1}^l)'$  be a  $2 \times 1$  vector, where

$$\mathbf{r}_{2T_0,j}^l = (T_0 b)^{-1} \sum_{t=r+1}^T Y_{t-l}(\mathbf{s}_0) \left( \frac{X_t(\mathbf{s}_0) - x}{b} \right)^j K \left( \frac{X_t(\mathbf{s}_0) - x}{b} \right), j = 0, 1.$$

Then the local linear estimator of  $g_{22}^l(x, \mathbf{s}_0)$  is

$$\hat{g}_{22}^l(x, \mathbf{s}_0) = \mathbf{e}_1' \mathbf{U}_{T_0}^{-1} \mathbf{R}_{2T_0}^l. \quad (7)$$

128 Therefore, the local linear estimator of the unknown vector  $\mathbf{g}_{22}(x, \mathbf{s}_0)$  is given by

$$\hat{\mathbf{g}}_{22}(x, \mathbf{s}_0) = (\hat{g}_{22}^1(x, \mathbf{s}_0), \dots, \hat{g}_{22}^q(x, \mathbf{s}_0))' = (\mathbf{e}_1' \mathbf{U}_{T_0}^{-1} \mathbf{R}_{2T_0}^1, \dots, \mathbf{e}_1' \mathbf{U}_{T_0}^{-1} \mathbf{R}_{2T_0}^q)'. \quad (8)$$

Finally, by (3), (5)–(7), and given both  $\boldsymbol{\lambda}(\mathbf{s}_0)$  and  $\boldsymbol{\alpha}(\mathbf{s}_0)$ , the unknown function  $g_0(X_t(\mathbf{s}_0))$  can be estimated by

$$\hat{g}_0(X_t(\mathbf{s}_0)) = \hat{g}_1(X_t(\mathbf{s}_0), \mathbf{s}_0) - \hat{\mathbf{g}}_{21}(X_t(\mathbf{s}_0), \mathbf{s}_0)' \boldsymbol{\lambda}(\mathbf{s}_0) - \hat{\mathbf{g}}_{22}(X_t(\mathbf{s}_0), \mathbf{s}_0)' \boldsymbol{\alpha}(\mathbf{s}_0). \quad (8)$$

### 129 3.2. Estimation of Parameters

Replacing  $g_0(X_t(\mathbf{s}_0))$  in (2) by its estimator (8), we re-write model (2) as

$$\hat{Y}_t(\mathbf{s}_0) = \hat{\mathbf{Z}}_t(\mathbf{s}_0)' \boldsymbol{\eta}(\mathbf{s}_0) + \varepsilon_t(\mathbf{s}_0),$$

130 where  $\hat{Y}_t(\mathbf{s}_0) = Y_t(\mathbf{s}_0) - \hat{g}_1(X_t(\mathbf{s}_0), \mathbf{s}_0)$ ,  $\hat{\mathbf{Z}}_t(\mathbf{s}_0) = \mathbf{Z}_t(\mathbf{s}_0) - \hat{\mathbf{g}}_2(X_t(\mathbf{s}_0), \mathbf{s}_0)$  and  $\boldsymbol{\eta}(\mathbf{s}_0) = (\boldsymbol{\lambda}(\mathbf{s}_0)', \boldsymbol{\alpha}(\mathbf{s}_0)')'$ ,

131  $\mathbf{Z}_t(\mathbf{s}_0) = (\mathbf{Z}_{1,t}(\mathbf{s}_0)', \mathbf{Z}_{2,t}(\mathbf{s}_0)')'$  and  $\hat{\mathbf{g}}_2(X_t(\mathbf{s}_0), \mathbf{s}_0) = (\hat{\mathbf{g}}_{21}(X_t(\mathbf{s}_0), \mathbf{s}_0)', \hat{\mathbf{g}}_{22}(X_t(\mathbf{s}_0), \mathbf{s}_0)')'$ .

Following the idea of adaptive Lasso, we estimate the interactions  $\boldsymbol{\lambda}(\mathbf{s}_0)$  and  $\boldsymbol{\alpha}(\mathbf{s}_0)$  at the spatial

location  $\mathbf{s}_0$  by minimizing the following penalized sum of squared errors

$$Q(\boldsymbol{\eta}(\mathbf{s}_0)) = L(\boldsymbol{\eta}(\mathbf{s}_0)) + T_0 \sum_{k=1}^N \sum_{i=1}^p \gamma_i^k(\mathbf{s}_0) |\lambda_{0k,i}| + T_0 \sum_{l=1}^q \beta_l(\mathbf{s}_0) |\alpha_{0l}|, \quad (9)$$

where  $L(\boldsymbol{\eta}(\mathbf{s}_0)) = \sum_{t=r+1}^T \hat{\varepsilon}_t(\mathbf{s}_0)^2$  is the sum of squared errors with  $\hat{\varepsilon}_t(\mathbf{s}_0) = \hat{Y}_t(\mathbf{s}_0) - \hat{\mathbf{Z}}_t(\mathbf{s}_0)' \boldsymbol{\eta}(\mathbf{s}_0)$ . The last two terms in (9) are adaptive Lasso penalties with regularization parameters  $\{\gamma_i^k(\mathbf{s}_0)\}_{i=1}^p \}_{k=1}^N$  and  $\{\beta_l(\mathbf{s}_0)\}_{l=1}^q$ . Thus, the penalized parameter estimator of  $\boldsymbol{\eta}(\mathbf{s}_0)$  is

$$\hat{\boldsymbol{\eta}}(\mathbf{s}_0) = \underset{\boldsymbol{\eta} \in \mathbb{R}^{Np+q}}{\operatorname{argmin}} \{Q(\boldsymbol{\eta}(\mathbf{s}_0))\}, \quad (10)$$

132 where the objective function  $Q(\boldsymbol{\eta}(\mathbf{s}_0))$  is given by (9).

To estimate the regularization parameters  $\{\gamma_i^k(\mathbf{s}_0)\}_{i=1}^p \}_{k=1}^N$  for  $k = 1, \dots, N$  and  $\{\beta_l(\mathbf{s}_0)\}_{l=1}^q$  in (9) with  $(Np + q)$  parameters, we let

$$\gamma_i^k(\mathbf{s}_0) = \gamma(\mathbf{s}_0) \frac{\log(T_0)}{T_0 |\tilde{\lambda}_i^k(\mathbf{s}_0)|} \quad \text{and} \quad \beta_l(\mathbf{s}_0) = \beta(\mathbf{s}_0) \frac{\log(T_0)}{T_0 |\tilde{\alpha}_l(\mathbf{s}_0)|}, \quad (11)$$

for  $k = 1, \dots, N, i = 1, \dots, p, j = 1, \dots, q$ , where  $\tilde{\lambda}_i^k(\mathbf{s}_0)$  and  $\tilde{\alpha}_l(\mathbf{s}_0)$  are the initial least squares estimators obtained by minimizing the objective function  $Q(\boldsymbol{\eta}(\mathbf{s}_0))$  without penalty (Wang et al., 2007; Zhu et al., 2010). That is,  $\tilde{\lambda}_i^k(\mathbf{s}_0)$  and  $\tilde{\alpha}_l(\mathbf{s}_0)$  are the minimizers of  $\sum_{t=r+1}^T \left\{ \hat{Y}_t(\mathbf{s}_0) - \hat{\mathbf{Z}}_t(\mathbf{s}_0)' \boldsymbol{\eta}(\mathbf{s}_0) \right\}^2$ . In the case of multicollinearity,  $\tilde{\lambda}_i^k(\mathbf{s}_0)$  and  $\tilde{\alpha}_l(\mathbf{s}_0)$  could be the ridge regression estimators (Zou, 2006). By (11), we reduce from  $(Np + q)$  to just two regularization parameters,  $\gamma(\mathbf{s}_0)$  and  $\beta(\mathbf{s}_0)$ . We select  $\gamma(\mathbf{s}_0)$  and  $\beta(\mathbf{s}_0)$  by a Bayesian information criterion (BIC),

$$\text{BIC}(\gamma(\mathbf{s}_0), \beta(\mathbf{s}_0)) = \log(\hat{\sigma}^2(\mathbf{s}_0)) + \kappa \log(T_0)/T_0, \quad (12)$$

133 where  $\hat{\sigma}^2(\mathbf{s}_0) = T_0^{-1} \sum_{t=r+1}^T \left\{ \hat{Y}_t(\mathbf{s}_0) - \hat{\mathbf{Z}}_t(\mathbf{s}_0)' \hat{\boldsymbol{\eta}}(\mathbf{s}_0) \right\}^2$ ,  $\hat{\boldsymbol{\eta}}(\mathbf{s}_0)$  is the penalized parameter estimator in

(10) corresponding to the regularization parameters  $\gamma(\mathbf{s}_0)$  and  $\beta(\mathbf{s}_0)$ , and  $\kappa$  is the effective number of parameters. For all possible combinations of  $\gamma(\mathbf{s}_0)$  and  $\beta(\mathbf{s}_0)$ , we choose the combination of regularization parameters that gives the minimum value of  $\text{BIC}(\gamma(\mathbf{s}_0), \beta(\mathbf{s}_0))$ .

After estimating  $\boldsymbol{\lambda}(\mathbf{s}_0)$  and  $\boldsymbol{\alpha}(\mathbf{s}_0)$ , we obtain the final estimator of the function  $g_0(x)$  by substituting  $\hat{\boldsymbol{\lambda}}(\mathbf{s}_0)$  and  $\hat{\boldsymbol{\alpha}}(\mathbf{s}_0)$  into (8),

$$\hat{g}_0(x) = \hat{g}_1(x, \mathbf{s}_0) - \hat{\mathbf{g}}_{21}(x, \mathbf{s}_0)' \hat{\boldsymbol{\lambda}}(\mathbf{s}_0) - \hat{\mathbf{g}}_{22}(x, \mathbf{s}_0)' \hat{\boldsymbol{\alpha}}(\mathbf{s}_0). \quad (13)$$

In the methodology developed here, the order of spatio-temporal lag interactions  $p$  and the order of temporal lag interactions  $q$  need to be pre-specified in model (1). In the housing price data example (Section 5), we recommend to consider  $p = q = 6$ , which is slightly larger than what was chosen by Akaike's information criterion with correction (AICc) in Appendix of Al-Sulami et al. (2017). This recommendation is owing to the advantage of the data-driven method which enables the selection of the more important lag interactions by penalization. Alternatives to local linear fitting considered above include local constant fitting, Nadaraya-Watson kernel estimate, and spline based methods. However, local linear fitting is known to outperform local constant fitting and Nadaraya-Watson kernel estimate (see, e.g., Fan and Gijbels, 1996).

#### 4. Asymptotic Properties

In semiparametric nonlinear regression time series model (1), there are two sets of parameters  $\lambda_{0k,i}$  and  $\alpha_{0,l}$ , with  $\lambda_{0k,i}$  reflecting the interaction between spatial locations  $\mathbf{s}_k$  and  $\mathbf{s}_0$  at time lag  $i$ , while  $\alpha_{0,l}$  reflecting the temporal lag interaction for a given spatial location  $\mathbf{s}_0$  at time lag  $l$ . Denote by  $S_N(\mathbf{s}_0) = \{\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_N\}$  the set of  $N$  spatial locations that are potentially interacting with spatial location  $\mathbf{s}_0$ . Let  $A_i(\mathbf{s}_0) = \{\mathbf{s}_k \in S_N(\mathbf{s}_0) : \lambda_{0k,i} \neq 0\}$  denote the set of spatial locations that do interact with spatial location  $\mathbf{s}_0$  at time lag  $i$ , whereas  $A_i^c(\mathbf{s}_0) = \{\mathbf{s}_k \in S_N(\mathbf{s}_0) : \lambda_{0k,i} = 0\}$  the set of spatial locations that do not interact with spatial location  $\mathbf{s}_0$  at time lag  $i$ . We let  $n_i(\mathbf{s}_0) =$

154  $\#A_i(\mathbf{s}_0)$  denote the cardinality of  $A_i(\mathbf{s}_0)$  and  $n_i^*(\mathbf{s}_0) = N - n_i(\mathbf{s}_0)$  denote the cardinality of  $A_i^c(\mathbf{s}_0)$ .

155 We partition  $\boldsymbol{\lambda}(\mathbf{s}_0)$ , the  $Np$ -dimensional vector of spatio-temporal lag interactions  $\lambda_{0k,i}$ , into  
 156  $\boldsymbol{\lambda}(\mathbf{s}_0) = (\boldsymbol{\lambda}_1(\mathbf{s}_0)', \boldsymbol{\lambda}_2(\mathbf{s}_0)')'$ , where  $\boldsymbol{\lambda}_1(\mathbf{s}_0)$  is a vector of  $\lambda_{0k,i}$  with  $i, k$  such that  $\mathbf{s}_k \in A_i(\mathbf{s}_0)$ , and  
 157  $\boldsymbol{\lambda}_2(\mathbf{s}_0)$  is a vector of  $\lambda_{0k,i}$  with  $i, k$  such that  $\mathbf{s}_k \in A_i^c(\mathbf{s}_0)$ . Therefore, with  $(Np)_0 = \sum_{i=1}^p n_i(\mathbf{s}_0)$ ,  
 158  $\boldsymbol{\lambda}_1(\mathbf{s}_0)$  is an  $(Np)_0$ -dimensional vector of non-zero spatio-temporal lag interactions and  $\boldsymbol{\lambda}_2(\mathbf{s}_0)$  is  
 159 an  $(Np - (Np)_0)$ -dimensional vector of zero spatio-temporal lag interactions.

160 Similarly, we partition  $\boldsymbol{\alpha}(\mathbf{s}_0)$ , the  $q$ -dimensional vector of temporal lag interactions  $\alpha_{0,l}$ , into  
 161  $\boldsymbol{\alpha}(\mathbf{s}_0) = (\boldsymbol{\alpha}_1(\mathbf{s}_0)', \boldsymbol{\alpha}_2(\mathbf{s}_0)')'$ . Let  $B(\mathbf{s}_0) = \{1 \leq l \leq q : \alpha_{0,l} \neq 0\}$  denote a set of  $q_0 = \#B(\mathbf{s}_0)$   
 162 non-zero temporal lag interactions and  $B^c(\mathbf{s}_0) = \{1 \leq l \leq q : \alpha_{0,l} = 0\}$ , a set of  $(q - q_0)$  zero  
 163 temporal lag interactions. Then,  $\boldsymbol{\alpha}_1(\mathbf{s}_0)$  is a  $q_0$ -dimensional vector of all  $\alpha_{0,l}$  such that  $l \in B(\mathbf{s}_0)$   
 164 and  $\boldsymbol{\alpha}_2(\mathbf{s}_0)$  is a  $(q - q_0)$ -dimensional vector of all  $\alpha_{0,l}$  such that  $l \in B^c(\mathbf{s}_0)$ .

165 Further, let  $\boldsymbol{\eta}^0(\mathbf{s}_0) = (\boldsymbol{\eta}_1^0(\mathbf{s}_0)', \boldsymbol{\eta}_2^0(\mathbf{s}_0)')'$  denote an  $(Np + q)$ -dimensional vector of true interac-  
 166 tions, where  $\boldsymbol{\eta}_1^0(\mathbf{s}_0) = (\boldsymbol{\lambda}_1(\mathbf{s}_0)', \boldsymbol{\alpha}_1(\mathbf{s}_0)')'$  is an  $((Np)_0 + q_0)$ -dimensional vector of non-zero interac-  
 167 tions and  $\boldsymbol{\eta}_2^0(\mathbf{s}_0) = (\boldsymbol{\lambda}_2(\mathbf{s}_0)', \boldsymbol{\alpha}_2(\mathbf{s}_0)')'$  is an  $(Np - (Np)_0 + q - q_0)$ -dimensional vector of zero interac-  
 168 tions. For the regularization parameters, we define  $a_{T_0}^*(\mathbf{s}_0) = \max \{\gamma_i^k(\mathbf{s}_0), \beta_l(\mathbf{s}_0) : \mathbf{s}_k \in A_i(\mathbf{s}_0), l \in B(\mathbf{s}_0)\}$ ,  
 169 and  $d_{T_0}^*(\mathbf{s}_0) = \min \{\gamma_i^k(\mathbf{s}_0), \beta_l(\mathbf{s}_0) : \mathbf{s}_k \in A_i^c(\mathbf{s}_0), l \in B^c(\mathbf{s}_0)\}$ , where the maximum and minimum  
 170 are taken over  $i = 1, \dots, p$ ,  $k = 1, \dots, N$ , and  $l = 1, \dots, q$ . We let  $|\cdot|$  and  $\|\cdot\|$  denote the  $L_1$   
 171 and  $L_2$  norm, respectively. We let  $\xrightarrow{P}$  and  $\xrightarrow{D}$  denote convergence in probability and convergence  
 172 in distribution, respectively.

173 For establishing the asymptotic properties, we impose regularity conditions given in Appendix 2  
 174 and provide proofs in Appendix 3 as web-based supplementary materials. In the theorems below,  
 175 we state the assumptions made about the regularization parameters and the asymptotic results. At  
 176 a given spatial location  $\mathbf{s}_0$ , we establish the consistency, sparsity, and asymptotic normality of the  
 177 penalized parameter estimator  $\hat{\boldsymbol{\eta}}(\mathbf{s}_0)$  obtained in (10) as follows.

178 **Theorem 1.** *Suppose that the regularity conditions in Appendix 2 hold and that the regularization*

179 parameter satisfies  $a_{T_0}^*(s_0) = o_p(T_0^{-1/2})$  as  $T_0 \rightarrow \infty$ . Then, there exists a global minimizer  $\hat{\boldsymbol{\eta}}(s_0)$   
 180 of the objective function  $Q(\boldsymbol{\eta}(s_0))$  such that  $\|\hat{\boldsymbol{\eta}}(s_0) - \boldsymbol{\eta}^0(s_0)\| = O_p\left(T_0^{-1/2} + a_{T_0}^*(s_0)\right)$

181 By Theorem 1, when the regularization parameters associated with the non-zero interactions  
 182 converge to zero at a rate of  $\sqrt{T_0}$ , the penalized parameter estimator  $\hat{\boldsymbol{\eta}}(s_0)$  is a global minimizer  
 183 and  $\sqrt{T_0}$ -consistent.

184 **Theorem 2.** Suppose that the regularity conditions in Appendix 2 hold and that the regulariza-  
 185 tion parameter satisfies  $\sqrt{T_0}d_{T_0}^*(s_0) \rightarrow \infty$  as  $T_0 \rightarrow \infty$  and  $\|\hat{\boldsymbol{\eta}}(s_0) - \boldsymbol{\eta}^0(s_0)\| = O_p\left(T_0^{-1/2}\right)$ .  
 186 Then,  $\hat{\boldsymbol{\eta}}_2(s_0) = \mathbf{0}$  with probability tending to one. That is, as  $T_0 \rightarrow \infty$ ,  $P\left(\hat{\boldsymbol{\lambda}}_2(s_0) = \mathbf{0}\right) \rightarrow$   
 187 1 and  $P\left(\hat{\boldsymbol{\alpha}}_2(s_0) = \mathbf{0}\right) \rightarrow 1$ .

188 Theorem 2 shows the sparsity of the penalized parameter estimator  $\hat{\boldsymbol{\eta}}(s_0)$ , which is  $\sqrt{T_0}$ -  
 189 consistent. That is, with the regularization parameters in (11), with probability tending to one,  
 190 the zero interactions are estimated to be 0.

**Theorem 3.** Suppose that the regularity conditions in Appendix 2 hold,  $\sqrt{T_0}a_{T_0}^*(s_0) \rightarrow 0$  and  
 $\sqrt{T_0}d_{T_0}^*(s_0) \rightarrow \infty$ . Then,

$$\sqrt{T_0}(\hat{\boldsymbol{\eta}}_1(s_0) - \boldsymbol{\eta}_1^0(s_0)) \xrightarrow{D} N(\mathbf{0}, \boldsymbol{\Psi}_1(s_0)),$$

191 where  $\boldsymbol{\Psi}_1(s_0) = \sigma_0^2 \boldsymbol{\Sigma}_1^{-1}(s_0)$  and  $\boldsymbol{\Sigma}_1(s_0) = \begin{pmatrix} \boldsymbol{\Sigma}_{\lambda_1}(s_0) & \boldsymbol{\Sigma}_{\lambda_1 \alpha_1}(s_0) \\ \boldsymbol{\Sigma}_{\alpha_1 \lambda_1}(s_0) & \boldsymbol{\Sigma}_{\alpha_1}(s_0) \end{pmatrix}$ , with  $\boldsymbol{\Sigma}_{\lambda_1}(s_0)$ ,  $\boldsymbol{\Sigma}_{\alpha_1}(s_0)$   
 192 and  $\boldsymbol{\Sigma}_{\lambda_1 \alpha_1}(s_0)$  defined in Appendix 1.

193 Theorem 3 establishes a central limit theorem for the penalized parameter estimator of the  
 194 non-zero interactions.

**Theorem 4.** Suppose that the regularity conditions in Theorem 3 hold and the bandwidth  $b$  satisfies  
 conditions (C7)(ii,iii) for  $x$  in the support of  $X(s_0)$  at the spatial location  $s_0$ . Then, as  $T \rightarrow \infty$ , we

have

$$\sqrt{T_0 b} \left[ \left\{ \hat{g}_0(x) - g_0(x) \right\} - (1/2)b^2 B_0(x, \mathbf{s}_0) \right] \xrightarrow{D} N(\mathbf{0}, \Gamma(x, \mathbf{s}_0)),$$

where  $B_0(x, \mathbf{s}_0) = \frac{\partial^2 g_0(x)}{\partial x^2} \int u^2 K(u) du$ ,  $\Gamma(x, \mathbf{s}_0) = \frac{\sigma^2(x, \mathbf{s}_0)}{p(x, \mathbf{s}_0)} \int K^2(u) du$ ,  $p(x, \mathbf{s}_0)$  is the marginal density function of  $X_t(\mathbf{s}_0)$ , and

$$\sigma^2(x, \mathbf{s}_0) = \text{Var} [(Y_t(\mathbf{s}_0) - \mathbf{Z}_{1,t}(\mathbf{s}_0)' \boldsymbol{\lambda}(\mathbf{s}_0) - \mathbf{Z}_{2,t}(\mathbf{s}_0)' \boldsymbol{\alpha}(\mathbf{s}_0)) | X_t(\mathbf{s}_0) = x].$$

195

196 Theorem 4 establishes the asymptotic properties of  $\hat{g}_0(x)$ , which is the nonparametric estimator  
197 of the unknown possibly nonlinear function  $g(x)$ .

## 198 5. Data Example: US Housing Price Index

199 We now return to the study of US housing price index (HPI), which is known to fluctuate  
200 between boom and recession periods. The consumer price index (CPI) is an important economic  
201 factor that may impact housing prices negatively or positively. The methodology developed in Sec-  
202 tions 2–4 can be applied to identify important spatio-temporal lag interactions as well as possibly  
203 nonlinear relationship between the HPI and the CPI.

### 204 5.1. Data and Exploratory Analysis

205 The HPI data are obtained from the Federal Home Loan Mortgage Corporation ([www.freddiemac.com](http://www.freddiemac.com)).  
206 This time series comprises 453 monthly observations of HPI from January 1975 to September 2012  
207 for each of the 50 states and the District of Columbia (DC). The HPI time series for DC and three  
208 other states (Hawaii, Texas, and Washington) are plotted in Figure 1(a) and are clearly nonstation-  
209 ary in time. These three states are chosen for illustration, because their geographical locations are  
210 outlying and thus it is not straightforward to decide which states are considered to their neighbors  
211 by the standard spatial econometric methods (see, e.g., Anselin, 1988). Further, time series of

the geometric return of the HPI are plotted in Figure 1(b), defined as  $Y_t(\mathbf{s}_0) = \log P_t(\mathbf{s}_0)/P_{t-1}(\mathbf{s}_0)$  where  $P_t(\mathbf{s}_0)$  is the HPI in month  $t$  at a given state  $\mathbf{s}_0$ . Unlike the HPI, the geometric returns of the HPI appear to be stationary in time. The box plots of the geometric returns of the HPI for the 50 states and DC are shown in Figure 2 and the temporal averages are mapped across states in Figure 3.

The CPI, on the other hand, is an average change of the prices of products over a time period and is an indicator of inflation in a country. Here we consider the CPI for all urban consumers (CPI-U), which represents about 80% of the American population, published by the Bureau of Labor Statistics ([www.bls.gov](http://www.bls.gov)). In particular, we focus on the CPI-U of primary residence (CPI-UR) and the time series of monthly CPI-UR in Figure 4(a) shows nonstationarity. Thus we consider a monthly increment of the CPI-UR, defined as  $x_t = \text{CPI}_t - \text{CPI}_{t-1}$ , where  $\text{CPI}_t$  is the CPI-UR in month  $t$ . Figure 4(b) plots the monthly increment CPI-UR indicating approximate stationarity, while Figure 4(c) plots its kernel density estimate in comparison to a Gaussian density estimate indicating that the distribution may not be Gaussian.

## 5.2. Model Fitting

For illustration, we evaluate possibly nonlinear relationships between the geometric return of the HPI and the increment of the CPI-UR in DC, Hawaii, Texas, and Washington (Figure 1), each with possible interactions with the other  $N = 50$  states (or district). We fit the semiparametric nonlinear regression time series model (1) to the data such that the monthly geometric return of the HPI (multiplied by 100)  $Y_t(\mathbf{s}_0)$  is the response variable for a given state  $\mathbf{s}_0$  and the monthly increment of CPI-UR  $X_t(\mathbf{s}_0) = x_{t-1}$  is the covariate that is the same for all the states for a given month. We use the temporal lag 1 of the increment of CPI (i.e.,  $x_{t-1}$ , not  $x_t$ ) here to avoid potential endogeneity of CPI in modeling HPI (c.f., Kuang et al., 2015; Panagiotidis and Printzis, 2016) and also for the purpose of forecasting. As mis-specification of the spatio-temporal lag interactions in model (1) may bias the estimate of the relationship between the geometric return of the HPI and

the monthly increment of the CPI-UR, we consider model (1) with temporal lag orders  $p = q = 6$ , which are chosen slightly larger than the orders  $p = q = 5$  selected by AICc in Al-Sulami et al. (2017). The results are similar for orders greater than  $p = q = 6$ , which is not unexpected as interactions for larger temporal lags are much smaller (see Figures 5 and A1–A6).

Thus, for a given state  $\mathbf{s}_0$  and  $N = 50$ , the following forecasting model can be easily obtained based on model (1):

$$Y_t(\mathbf{s}_0) = g_0(x_{t-1}) + \sum_{i=1}^6 \sum_{k=1}^N \lambda_{0k,i} Y_{t-i}(\mathbf{s}_k) + \sum_{l=1}^6 \alpha_{0,l} Y_{t-l}(\mathbf{s}_0) + \varepsilon_t(\mathbf{s}_0), \quad (14)$$

for  $t = 1, \dots, T (= 452)$  and  $k = 1, \dots, N (= 50)$ . By (14), we have the fitted response variable

$$\hat{Y}_t(\mathbf{s}_0) = \sum_{i=1}^6 \sum_{k=1}^N \lambda_{0k,i} \hat{Y}_{t-i}(\mathbf{s}_k) + \sum_{l=1}^6 \alpha_{0,l} \hat{Y}_{t-l}(\mathbf{s}_0), \quad (15)$$

where the conditional means  $\hat{Y}_t(\mathbf{s}_0) = Y_t(\mathbf{s}_0) - \hat{E}[Y_t(\mathbf{s}_0)|x_{t-1}]$ ,  $\hat{Y}_{t-i}(\mathbf{s}_k) = Y_{t-i}(\mathbf{s}_k) - \hat{E}[Y_{t-i}(\mathbf{s}_k)|x_{t-1}]$ , and  $\hat{Y}_{t-l}(\mathbf{s}_0) = Y_{t-l}(\mathbf{s}_0) - \hat{E}[Y_{t-l}(\mathbf{s}_0)|x_{t-1}]$  are estimated by local linear regression with bandwidth selected by cross-validation (CV). Then, we obtain the penalized parameter estimate of  $\lambda_{0k,i}$  and  $\alpha_{0,l}$  by minimizing the following penalized sum of squared errors with a single regularization parameter,

$$Q(\boldsymbol{\eta}(\mathbf{s}_0)) = \sum_{t=7}^T \left[ \hat{Y}_t(\mathbf{s}_0) - \hat{\mathbf{Z}}_t(\mathbf{s}_0)' \boldsymbol{\eta}(\mathbf{s}_0) \right]^2 + T_0 \sum_{i=1}^6 \sum_{k=0}^N \gamma_i^k(\mathbf{s}_0) |\lambda_{0k,i}|,$$

where  $\boldsymbol{\eta}(\mathbf{s}_0) = (\lambda_{01,1}, \dots, \lambda_{0N,1}, \dots, \lambda_{01,6}, \dots, \lambda_{0N,6}, \lambda_{00,1}, \dots, \lambda_{00,6})'$  with  $\lambda_{00,i}$  denoting  $\alpha_{0,i}$ . Further,  $\gamma_i^k(\mathbf{s}_0) = \gamma(\mathbf{s}_0) \frac{\log(T_0)}{T_0 |\lambda_{0k,i}|}$ , where  $\tilde{\lambda}_{0k,i}$  are the initial estimates of  $\lambda_{0k,i}$  by the least squares for model (15) and  $\gamma(\mathbf{s}_0)$  is selected by the BIC. For simplicity, here we consider only one tuning parameter  $\gamma(\mathbf{s}_0)$  for both spatio-temporal and temporal interactions.

For DC, the estimates of the temporal lag interactions  $\alpha_{0,l}$  and the unknown function  $g(\cdot)$  are plotted in Figure 5(a) and (b), respectively. Further, Figure 5(a) shows that the estimates of the tem-



poral lag interactions for DC  $\alpha_{0,\ell}$  decrease in magnitude over the time lags and appear relatively strong for time lags 1–4 but weak for time lags 5–6. It is also interesting that the temporal lag interactions oscillate around 0, which may indicate a self-adjustment or reversion of the housing price change along time. As commented by a referee, this is also probably due to the weakly correlated errors that are compensating each other at adjacent lags. In Figure 5(b), the estimated function  $g(\cdot)$  is nearly flat, implying a weak or no relationship between the geometric return of HPI return and monthly increment of CPI-UR for DC, after taking into account the spatio-temporal lag interactions. We have also considered a residual analysis for model diagnostics and in particular, we plot the autocorrelation functions to evaluate the assumption of independence for the innovations. For example, we find that compared with the strong autocorrelation of the original return series  $Y_t$ , the time series of residuals of (14), denoted by  $e_t$ , is largely uncorrelated for DC (Figure 6).

For the other three states of (Hawaii, Texas, and Washington), we have presented the estimates for  $\alpha_{0,\ell}$  and  $g(\cdot)$  in Figure 5. Figures A1–A6 in Appendix 5 plot the estimated spatio-temporal lag interactions for all 50 states and DC. Figure 5 shows that the patterns of the temporal lag interactions for Hawaii, Texas, and Washington are similar to that for the DC. However, the relationships between HPI and CPI-UR for the three states differ from that for DC, with obvious nonlinearity. For Texas, the relationship is positive when the monthly increment of CPI-UR is, approximately, below 0.3 and above 0.75, relatively flat when the monthly increment of CPI-UR is between about 0.3 and 0.5, and negative when the monthly increment of CPI-UR is between about 0.5 and 0.75. For Washington, the relationship between HPI and CPI-UR is positive when the CPI-UR is, approximately, below 0.4 and above 0.8, and turns negative when the CPI-UR is between 0.4 and 0.8. For Hawaii, the relationship between HPI and CPI-UR starts off negative when the CPI-UR is below 0.2, turns positive when the CPI-UR is between 0.2 and 0.75, and turns negative again when the CPI-UR is above 0.75. These are potentially interesting patterns about housing prices and investments in different parts of the country. In particular, important spatio-temporal lag interactions could occur between states that are not necessarily close in geographical distance.

### 5.3. Prediction Performance

To further evaluate the model fitting in Section 5.2, we consider prediction and comparison with alternative approaches. We partition the data into two parts. The first part has the first  $T = 402$  observations and is used for model fitting. The second part has the last 50 observations and we consider one-step ahead prediction for testing. The performance of the prediction is compared between two different forms of the function  $g$  in terms of mean squared prediction error (MSPE). In the first form,  $g$  is assumed to be a linear function,  $g_L(x_{t-1}) = a_0(\mathbf{s}_0) + a_1(\mathbf{s}_0)x_{t-1}$ , where the linear coefficients  $a_0(\mathbf{s}_0)$  and  $a_1(\mathbf{s}_0)$  are for a given location  $\mathbf{s}_0$ . In the second form, we apply our nonparametric specification of  $g$  in model (14) and denote it as  $g_{NP}$ . The MSPE values computed based on the two forms of  $g$  and the testing data are provided in the first two columns of Table 1. The results suggest that a flexible, possibly nonlinear, form for  $g$  that relates the geometric return of the HPI to the monthly increment of the CPI-UR is helpful and improves the accuracy of prediction for the three states and DC.

Further, to evaluate the penalized approach that enables the simultaneous selection and estimation of the spatio-temporal lag interactions, we compare our methodology with the least squares approach without regularization, which we denote as  $g_{LS}$  in the third column of Table 1. The MSPE values under least squares are much larger, demonstrating that the penalization has helped to improve prediction accuracy.

In summary, our data analysis based on the proposed model (1) is fully data-driven in the specification of spatio-temporal interactions between spatial locations, which allows spatial weights to vary over different temporal lags. Like purely nonparametric approach, our method can help to extract useful information and knowledge for improvement of model forecasting, but itself may not always perform the best in forecasting when compared with other more prior information imposed spatial time series models. However, we do find that for the DC, the forecasting based on the proposed model with refined optimal tuning parameter achieves an MSPE as small as 0.06905, which outperforms that of Al-Sulami et al. (2017) with pre-specified inverse distance based spatial

weights, which has an MSPE of 0.09416. The data-driven based spatial weight matrix is more adaptive to the data and hence likely more robust in forecasting.

## 6. Conclusions and Discussion

In this paper, we have proposed a semiparametric data-driven nonlinear method that allows potentially nonlinear relationship between a response variable and a covariate, as well as spatial interactions that could vary by temporal lag. In economic applications, spatial weights are usually specified *a priori* between two spatial locations or units that are neighbors of each other according to a neighborhood structure, which can be somewhat subjective. In contrast, our penalized estimation method identifies the important lag interactions across space and over time, providing a data-driven way to determine the spatio-temporal weights more objectively. If it is desirable to consider only interactions between neighboring spatial units guided by theory or for ease of interpretation, our method could be modified by setting the interactions to zero between spatial units that are not neighbors of each other.

We have applied our methodology to analyze a US housing price data set focusing on DC, Hawaii, Texas, and Washington. The data analysis has revealed nonlinear relationships between housing price and consumer price index for the three states but not DC. The prediction based on our method can be more accurate than the more standard approaches that assume linearity or without penalization for the US housing price data. Our approach here is different from the time simultaneous perspective of linear spatial autoregressive models (see, e.g., Ahrens and Bhattacharjee, 2015; Qu and Lee, 2015). Extending our methodology to spatial autoregressive models with nonlinear covariate structure would be interesting for understanding simultaneous spatial interactions, although it would not be suitable for forecasting.

Further, our proposed methodology works well for a given spatial location in relation to other spatial locations on a lattice, but may not be optimal in identifying the whole network. It would be interesting to extend the methodology to simultaneously examine all the locations on the lat-

tice, although we may face the challenge of estimating a much larger number of parameters. For example, with  $N = 51$  and  $p = q = 6$ , there will be 15606 parameters. Application of dimension reduction techniques of semi-data-driven estimation with prior information may be desirable. We leave these for future research.

Finally, we have established the asymptotic properties of the estimation procedure. Theoretical justifications for the nonparametric and nonlinear estimation, as well as the adaptive LASSO for selecting lagged variables in time series and space-time model are well studied under general near epoch dependence structure (c.f., Appendix 2). Thus, the results obtained here can be applied to a wide range of linear and nonlinear time series processes in statistics and econometrics (c.f., Li et al., 2012; Lu and Linton, 2007).

## Acknowledgments

We thank the co-editors and two reviewers for their constructive comments that have improved the content and presentation of this paper. Lu's research was partially supported by an EU's Marie Curie Career Integration Grant, which is acknowledged.

## References

- Ahrens, A. and A. Bhattacharjee (2015). Two-step Lasso estimation of the spatial weights matrix. *Econometrics* 3(1), 128–155.
- Akaike, H. (1973). Information theory and an extension of the maximum likelihood principle. In *Second International Symposium on Information Theory*, Volume 1, pp. 267–281. Springer Verlag.
- Al-Sulami, D., Z. Jiang, Z. Lu, and J. Zhu (2017). Estimation for semiparametric nonlinear regression of irregularly located spatial time-series data. *Econometrics and Statistics* 2, 22–35.

- Anselin, L. (1988). *Spatial Econometrics: Methods and Models*. Springer Science & Business Media.
- Bhattacharjee, A. and S. Holly (2011). Structural interactions in spatial panels. *Empirical Economics* 40(1), 69–94.
- Bhattacharjee, A. and S. Holly (2013). Understanding interactions in social networks and committees. *Spatial Economic Analysis* 8(1), 23–53.
- Bhattacharjee, A. and C. Jensen-Butler (2013). Estimation of the spatial weights matrix under structural constraints. *Regional Science and Urban Economics* 43(4), 617–634.
- Breiman, L. (1996). Heuristics of instability and stabilization in model selection. *The Annals of Statistics* 24(6), 2350–2383.
- de Souza, P. C. L. (2012). Estimating networks: Lasso for spatial weights. *The 34th Meeting of the Brazilian Econometric Society*.
- Fan, J. and I. Gijbels (1996). *Local Polynomial Modelling and its Applications*. Chapman & Hall.
- Fan, J. and R. Li (2001). Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association* 96(456), 1348–1360.
- Fan, J. and Q. Yao (2003). *Nonlinear Time Series: Nonparametric and Parametric Methods*. Springer Verlag.
- Frank, I. and J. Friedman (1993). A statistical view of some chemometrics regression tools. *Technometrics*, 109–135.
- Gao, J. (2007). *Nonlinear Time Series: Semiparametric and Nonparametric Methods*. Chapman & Hall.

367 Gao, J., Z. Lu, and D. Tjøstheim (2006). Estimation in semiparametric spatial regression. *The*  
368 *Annals of Statistics* 34(3), 1395–1435.

369 Griffith, D. A. and F. Lagona (1998). On the quality of likelihood-based estimators in spatial  
370 autoregressive models when the data dependence structure is misspecified. *Journal of Statistical*  
371 *Planning and Inference* 69(1), 153–174.

372 Hallin, M., Z. Lu, and L. Tran (2004). Local linear spatial regression. *The Annals of Statis-*  
373 *tics* 32(6), 2469–2500.

374 Haufe, S., G. Nolte, K. Mueller, and N. Kraemer (2009). Sparse causal discovery in multivariate  
375 time series. *Arxiv preprint arXiv:0901.2234*.

376 Hefley, T., M. Hooten, E. Hanks, R. Russell, and D. Walsh (2017). The Bayesian group lasso  
377 for confounded spatial data. *Journal of Agricultural, Biological and Environmental Statistics*,  
378 42–59.

379 Hsu, N., H. Hung, and Y. Chang (2008). Subset selection for vector autoregressive processes using  
380 Lasso. *Computational Statistics & Data Analysis* 52(7), 3645–3657.

381 Huang, H., N. Hsu, D. Theobald, and F. Breidt (2010). Spatial LASSO with applications to GIS  
382 model selection. *Journal of Computational and Graphical Statistics* 19(4), 963–983.

383 Kuang, W., P. Liu, et al. (2015). Inflation and house prices: theory and evidence from 35 major  
384 cities in china. *International Real Estate Review* 18(2), 217–240.

385 Lam, C. and P. C. Souza (2013). Regularization for spatial panel time series using the adaptive  
386 lasso. Technical report, Mimeo.

387 Li, D., Z. Lu, and O. Linton (2012). Local linear fitting under near epoch dependence: Uniform  
388 consistency with convergence rates. *Econometric Theory*, 28(5), 935–958.

- 389 Lu, Z. and O. Linton (2007). Local linear fitting under near epoch dependence. *Econometric*  
390 *Theory* 23(1), 37.
- 391 Lu, Z., A. Lundervold, D. Tjøstheim, and Q. Yao (2007). Exploring spatial nonlinearity using  
392 additive approximation. *Bernoulli* 13(2), 447–472.
- 393 Lu, Z., D. Steinskog, D. Tjøstheim, and Q. Yao (2009). Adaptively varying-coefficient spatiotem-  
394 poral models. *Journal of the Royal Statistical Society: Series B (statistical methodology)* 71(4),  
395 859–880.
- 396 Lu, Z., D. Tjøstheim, and Q. Yao (2008). Spatial smoothing, nugget effect and infill asymptotics.  
397 *Statistics & Probability Letters* 78(18), 3145–3151.
- 398 Manresa, E. (2013). Estimating the structure of social interactions using panel data. *Unpublished*  
399 *Manuscript. CEMFI, Madrid.*
- 400 McQuarrie, A. and C. Tsai (1998). *Regression and Time Series Model Selection*. World Scientific.
- 401 Panagiotidis, T. and P. Printzis (2016). On the macroeconomic determinants of the housing market  
402 in greece: A vecm approach. *International Economics and Economic Policy* 13(3), 387–409.
- 403 Qu, X. and L.-f. Lee (2015). Estimating a spatial autoregressive model with an endogenous spatial  
404 weight matrix. *Journal of Econometrics* 184(2), 209–232.
- 405 Ramanathan, R. (1992). *Introductory Econometrics with Applications*. Harcourt Brace Colleges.
- 406 Ren, Y. and X. Zhang (2010). Subset selection for vector autoregressive processes via adaptive  
407 Lasso. *Statistics & Probability Letters* 80(23-24), 1705–1712.
- 408 Reyes, P., J. Zhu, and B. Aukema (2012). Selection of spatial-temporal lattice models: Assessing  
409 the impact of climate conditions on a mountain pine beetle outbreak. *Journal of Agricultural,*  
410 *Biological, and Environmental Statistics* 17, 508–525.

- 411 Schwarz, G. (1978). Estimating the dimension of a model. *The Annals of Statistics* 6(2), 461–464.
- 412 Shumway, R. and D. Stoffer (2000). *Time Series Analysis and its Applications*. Springer Verlag.
- 413 Smith, T. E. (2009). Estimation bias in spatial models with strongly connected weight matrices.  
414 *Geographical Analysis* 41(3), 307–332.
- 415 Song, J. J. and V. De Oliveira (2012). Bayesian model selection in spatial lattice models. *Statistical*  
416 *Methodology* 9(1-2), 228–238.
- 417 Stakhovych, S. and T. H. Bijmolt (2009). Specification of spatial models: A simulation study on  
418 weights matrices. *Papers in Regional Science* 88(2), 389–408.
- 419 Stetzer, F. (1982). Specifying weights in spatial forecasting models: the results of some experi-  
420 ments. *Environment and Planning A* 14(5), 571–584.
- 421 Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal*  
422 *Statistical Society: Series B (Methodological)*, 267–288.
- 423 Tibshirani, R. et al. (1997). The Lasso method for variable selection in the Cox model. *Statistics*  
424 *in medicine* 16(4), 385–395.
- 425 Van De Geer, S. (2008). High-dimensional generalized linear models and the lasso. *The Annals of*  
426 *Statistics* 36(2), 614–645.
- 427 Wang, H., G. Li, and C. Tsai (2007). Regression coefficient and autoregressive order shrinkage and  
428 selection via the Lasso. *Journal of the Royal Statistical Society: Series B (Statistical Methodol-*  
429 *ogy)* 69(1), 63–78.
- 430 Zhang, W., Q. Yao, H. Tong, and N. Stenseth (2003). Smoothing for spatiotemporal models and  
431 its application to modeling muskrat-mink interaction. *Biometrics* 59(4), 813–821.



- 432 Zhao, P. and B. Yu (2007). On model selection consistency of Lasso. *Journal of Machine Learning*  
433 *Research* 7(2), 2541.
- 434 Zhu, J., H. Huang, and P. Reyes (2010). On selection of spatial linear models for lattice data.  
435 *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 72(3), 389–402.
- 436 Zou, H. (2006). The adaptive lasso and its oracle properties. *Journal of the American Statistical*  
437 *Association* 101(476), 1418–1429.
- 438 Zou, H. and T. Hastie (2005). Regularization and variable selection via the elastic net. *Journal of*  
439 *the Royal Statistical Society: Series B (Statistical Methodology)* 67(2), 301–320.

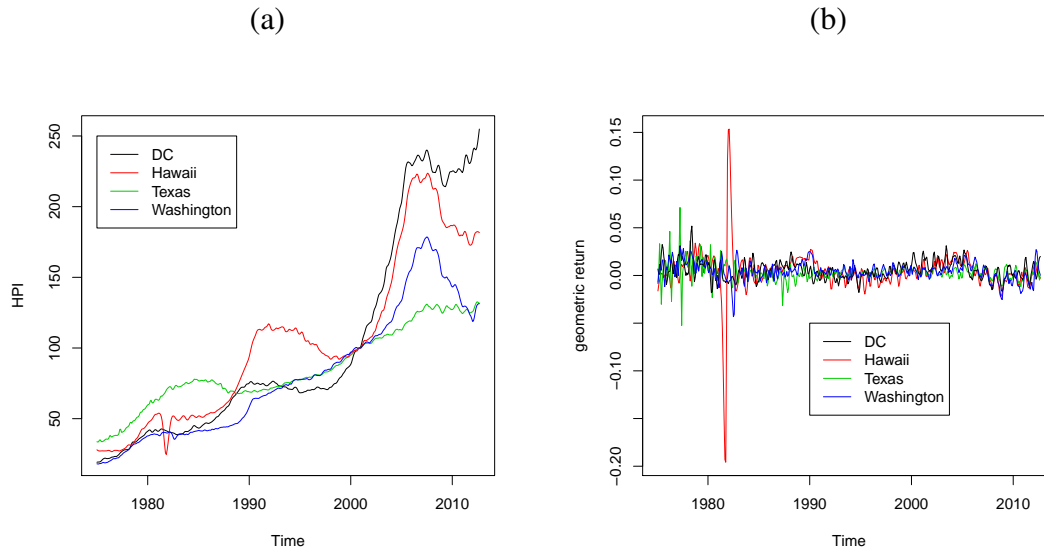


FIGURE 1: Time series plots of (a) the monthly housing price index (HPI) and (b) its geometric return for District of Columbia (DC), Hawaii, Texas and Washington from January 1975 to September 2012.

TABLE 1: Mean squared prediction error (MSPE) based on the last 50 months of housing price index (HPI) geometric returns assuming linearity  $g_L$ , nonlinearity  $g_{NP}$ , and using ordinary least squares  $g_{LS}$  without penalization.

	$g_L$	$g_{NP}$	$g_{LS}$
DC	0.090	0.069	1.574
Hawaii	0.154	0.138	5.540
Texas	0.288	0.282	1.237
Washington	0.261	0.227	0.972

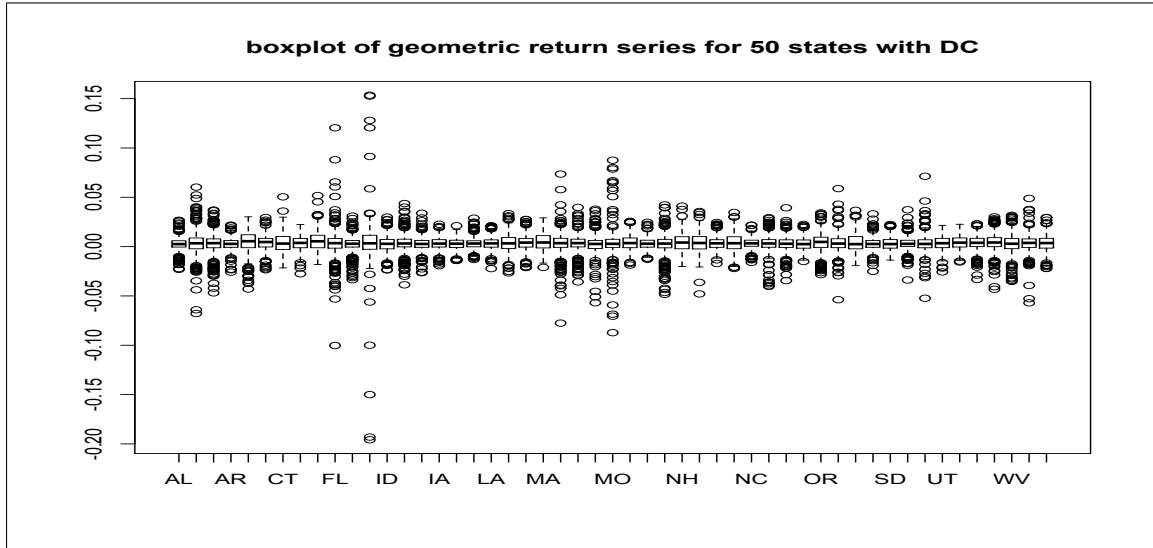


FIGURE 2: Boxplots of the geometric returns of housing price index (HPI) from January 1975 to September 2012 for each of the 50 US states and one district. The states are placed along the  $x$ -axis in the order of AL: Alabama; AK: Alaska; AZ: Arizona; AR: Arkansas; CA: California; CO: Colorado; CT: Connecticut; DE: Delaware; DC: District of Columbia; FL: Florida; GA: Georgia; HI: Hawaii; ID: Idaho; IL: Illinois; IN: Indiana; IA: Iowa; KS: Kansas; KY: Kentucky; LA: Louisiana; ME: Maine; MD: Maryland; MA: Massachusetts; MI: Michigan; MN: Minnesota; MS: Mississippi; MO: Missouri; MT: Montana; NE: Nebraska; NV: Nevada; NH: New Hampshire; NJ: New Jersey; NM: New Mexico; NY: New York; NC: North Carolina; ND: North Dakota; OH: Ohio; OK: Oklahoma; OR: Oregon; PA: Pennsylvania; RI: Rhode Island; SC: South Carolina; SD: South Dakota; TN: Tennessee; TX: Texas; UT: Utah; VT: Vermont; VA: Virginia; WA: Washington; WV: West Virginia; WI: Wisconsin; WY: Wyoming.

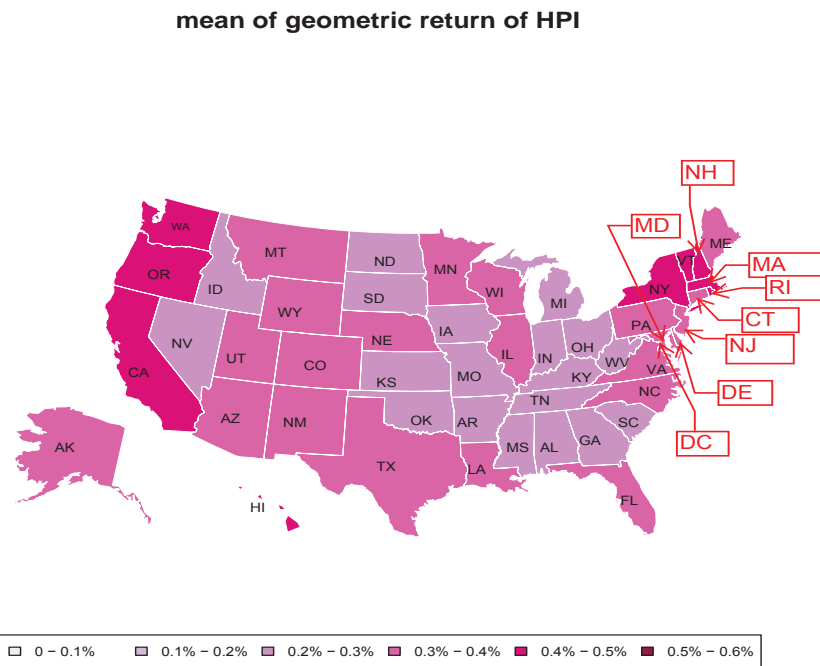


FIGURE 3: Map of the mean geometric return of the housing price index (HPI) averaged from January 1975 to September 2012 across the 50 states and Washington DC.

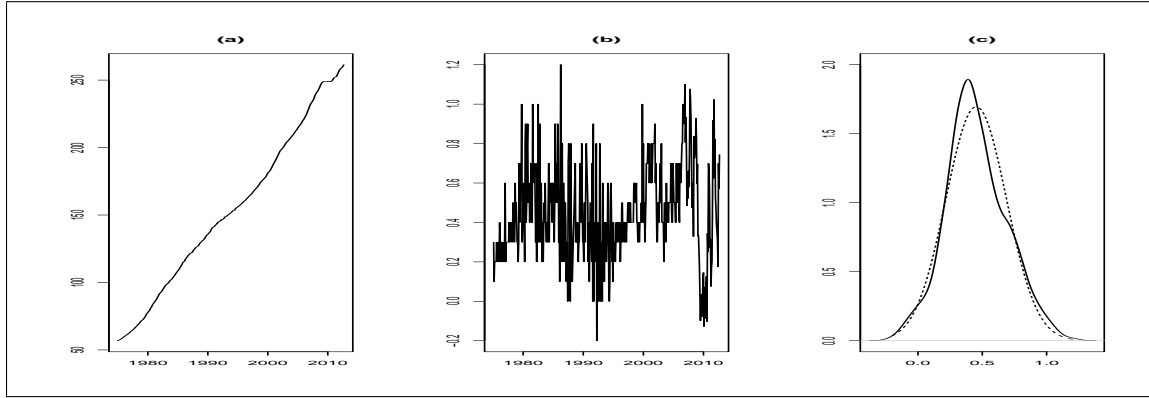


FIGURE 4: Plots of consumer price index (CPI): (a) original time series; (b) monthly increments; and (c) kernel density estimate (solid curve) compared with Gaussian density estimate with the same mean and variance (dashed curve).

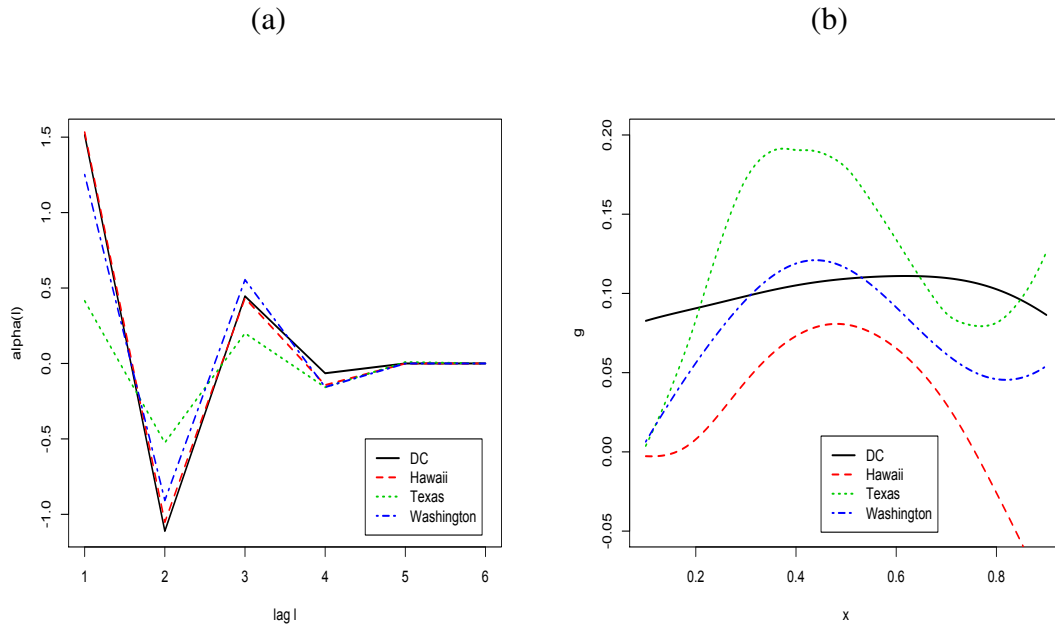


FIGURE 5: Plots of (a) the estimated temporal lag interactions  $\hat{\alpha}_\ell$  for  $\ell = 1, \dots, 6$ ; (b) the estimated function  $\hat{g}_0$  for DC, Hawaii, Texas and Washington.

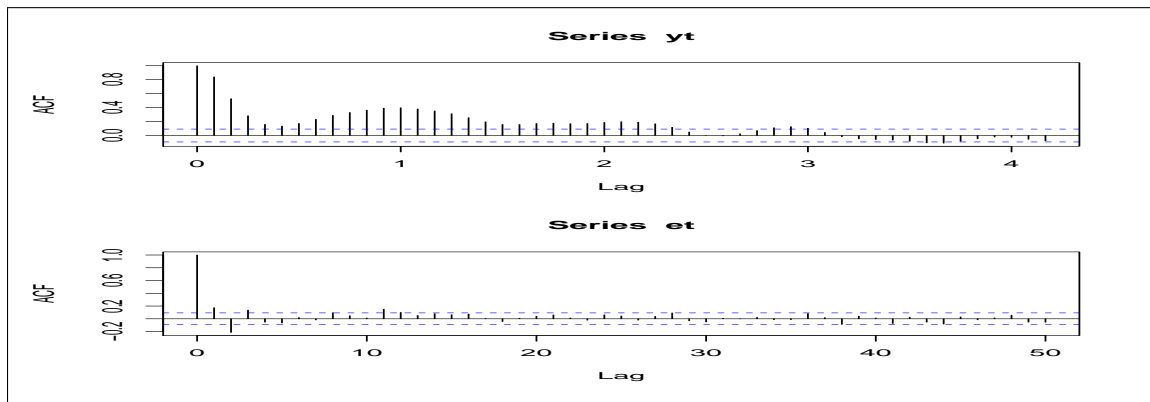


FIGURE 6: The autocorrelation function (ACF) for the original return of housing price index  $Y_t$  and the residual  $e_t$  for Washington DC.

## Web-based Supplementary Materials

### Appendix to “On a Semiparametric Data-driven Nonlinear Model with Penalized Spatio-temporal Lag Interactions”

In this appendix, we present additional notation in Section 1, regularity conditions in Section 2, proofs for Theorems 1–4 in Section 3, a simulation study in Section 4, and additional results on the estimated spatio-temporal lag interactions for all the US states in Section 5.

#### 1. Notation

Recall  $\mathbf{Z}_t(\mathbf{s}_0) = (\mathbf{Y}_{t-1}(\mathbf{s}_0)', \mathbf{Y}_{t-2}(\mathbf{s}_0)', \dots, \mathbf{Y}_{t-p}(\mathbf{s}_0)', Y_{t-1}(\mathbf{s}_0), Y_{t-2}(\mathbf{s}_0), \dots, Y_{t-q}(\mathbf{s}_0))'$ , with  $\mathbf{Y}_{t-i}(\mathbf{s}_0) = (Y_{t-i}(\mathbf{s}_1), Y_{t-i}(\mathbf{s}_2), \dots, Y_{t-i}(\mathbf{s}_N))'$ . Define  $\mathbf{Z}_t^*(\mathbf{s}_0) = \mathbf{Z}_t(\mathbf{s}_0) - E[\mathbf{Z}_t(\mathbf{s}_0)|X_t(\mathbf{s}_0)]$  and  $\mathbf{R}_t(\mathbf{s}_0) = \varepsilon_t(\mathbf{s}_0)\mathbf{Z}_t^*(\mathbf{s}_0)$ . Then, let  $\Sigma(\mathbf{s}_0) = E[\mathbf{Z}_t^*(\mathbf{s}_0)\mathbf{Z}_t^{*'}(\mathbf{s}_0)]$  denote an  $(Np + q) \times (Np + q)$  matrix,

$$\Gamma(\mathbf{s}_0) = \sum_{k=0}^{\infty} E[\mathbf{R}_t(\mathbf{s}_0)\mathbf{R}_{t+k}(\mathbf{s}_0)'], \quad \text{and } \Psi(\mathbf{s}_0) = \Sigma^{-1}(\mathbf{s}_0)\Gamma(\mathbf{s}_0)(\Sigma^{-1}(\mathbf{s}_0))'.$$

By Condition (C4)(i) in Section 2, we have  $\Gamma(\mathbf{s}_0) = E[\mathbf{R}_t(\mathbf{s}_0)\mathbf{R}_t(\mathbf{s}_0)'] = \sigma_0^2\Sigma(\mathbf{s}_0)$  and hence  $\Psi(\mathbf{s}_0) = \sigma_0^2\Sigma^{-1}(\mathbf{s}_0)$ .

With the notation introduced prior to the theorems in Section 4, let  $p_0 = \#\{1 \leq i \leq p : n_i(\mathbf{s}_0) \neq 0\}$ . For notational convenience, we suppose without loss of generality that the non-zero  $\lambda_{0k^*,i}$ 's are for those  $k^* = 1, \dots, n_i$ ,  $i = 1, \dots, p_0$ , with  $n_i = n_i(\mathbf{s}_0)$ , and the non-zero  $\alpha_{0,i}$ 's are for those  $i = 1, \dots, q_0$ . Similarly, let  $\mathbf{Z}_t^{*(1)}(\mathbf{s}_0) = \mathbf{Z}_t^{(1)}(\mathbf{s}_0) - E[\mathbf{Z}_t^{(1)}(\mathbf{s}_0)|X_t(\mathbf{s}_0)]$  and  $\mathbf{R}_t^{(1)}(\mathbf{s}_0) = \varepsilon_t(\mathbf{s}_0)\mathbf{Z}_t^{*(1)}(\mathbf{s}_0)$ , with

$$\mathbf{Z}_t^{(1)}(\mathbf{s}_0) = (Y_{t-1}(\mathbf{s}_1), \dots, Y_{t-1}(\mathbf{s}_{n_1}), \dots, Y_{t-p_0}(\mathbf{s}_1), \dots, Y_{t-p_0}(\mathbf{s}_{n_{p_0}}), Y_{t-1}(\mathbf{s}_0), Y_{t-2}(\mathbf{s}_0), \dots, Y_{t-q_0}(\mathbf{s}_0))'.$$

Then, let  $\Sigma_1(\mathbf{s}_0) = E \left[ \mathbf{Z}_t^{*(1)}(\mathbf{s}_0) \mathbf{Z}_t^{*(1)}(\mathbf{s}_0)' \right]$  denote a  $(\sum_{i=1}^{p_0} n_i + q_0) \times (\sum_{i=1}^{p_0} n_i + q_0)$  matrix and

$$\Gamma_1(\mathbf{s}_0) = \sum_{k=0}^{\infty} E \left[ \mathbf{R}_t^{(1)}(\mathbf{s}_0) \mathbf{R}_{t+k}^{(1)}(\mathbf{s}_0)' \right].$$

By Condition (C4)(i) in Section 2, we have  $\Gamma_1(\mathbf{s}_0) = E \left[ \mathbf{R}_t^{(1)}(\mathbf{s}_0) \mathbf{R}_t^{(1)}(\mathbf{s}_0)' \right] = \sigma_0^2 \Sigma_1(\mathbf{s}_0)$ .

For the sub-matrices of  $\Gamma_1(\mathbf{s}_0)$ , we let  $\mathbf{Z}_t^{*(11)}(\mathbf{s}_0) = \mathbf{Z}_t^{(11)}(\mathbf{s}_0) - E \left[ \mathbf{Z}_t^{(11)}(\mathbf{s}_0) | X_t(\mathbf{s}_0) \right]$  and  $\mathbf{R}_t^{(11)}(\mathbf{s}_0) = \varepsilon_t(\mathbf{s}_0) \mathbf{Z}_t^{*(11)}(\mathbf{s}_0)$  with  $\mathbf{Z}_t^{(11)}(\mathbf{s}_0) = (Y_{t-1}(\mathbf{s}_1), \dots, Y_{t-1}(\mathbf{s}_{n_1}), \dots, Y_{t-p_0}(\mathbf{s}_1), \dots, Y_{t-p_0}(\mathbf{s}_{n_{p_0}}))'$ . Then, let  $\Sigma_{\lambda_1}(\mathbf{s}_0) = E \left[ \mathbf{Z}_t^{*(11)}(\mathbf{s}_0) \mathbf{Z}_t^{*(11)}(\mathbf{s}_0)' \right]$  denote a  $\sum_{i=1}^{p_0} n_i \times \sum_{i=1}^{p_0} n_i$  matrix and

$$\Gamma_{\lambda_1}(\mathbf{s}_0) = \sum_{k=0}^{\infty} E \left[ \mathbf{R}_t^{(11)}(\mathbf{s}_0) \mathbf{R}_{t+k}^{(11)}(\mathbf{s}_0)' \right].$$

By Condition (C4)(i) in Section 2, we have  $\Gamma_{\lambda_1}(\mathbf{s}_0) = E \left[ \mathbf{R}_t^{(11)}(\mathbf{s}_0) \mathbf{R}_t^{(11)}(\mathbf{s}_0)' \right] = \sigma_0^2 \Sigma_{\lambda_1}(\mathbf{s}_0)$ .

Let  $\mathbf{Z}_t^{*(12)}(\mathbf{s}_0) = \mathbf{Z}_t^{(12)}(\mathbf{s}_0) - E \left[ \mathbf{Z}_t^{(12)}(\mathbf{s}_0) | X_t(\mathbf{s}_0) \right]$  and  $\mathbf{R}_t^{(12)}(\mathbf{s}_0) = \varepsilon_t(\mathbf{s}_0) \mathbf{Z}_t^{*(12)}(\mathbf{s}_0)$ , with  $\mathbf{Z}_t^{(12)}(\mathbf{s}_0) = ((Y_{t-(p_0+1)}(\mathbf{s}_{k^*}), \dots, Y_{t-p}(\mathbf{s}_{k^*})), k^* = 1, \dots, N; (Y_{t-i}(\mathbf{s}_{n_i+1}), \dots, Y_{t-i}(\mathbf{s}_N)), i = 1, \dots, p_0)'$ . Then, let  $\Sigma_{\lambda_2}(\mathbf{s}_0) = E \left[ \mathbf{Z}_t^{*(12)}(\mathbf{s}_0) \mathbf{Z}_t^{*(12)}(\mathbf{s}_0)' \right]$  denote an  $(Np - \sum_{i=1}^{p_0} n_i) \times (Np - \sum_{i=1}^{p_0} n_i)$  matrix and  $\Gamma_{\lambda_2}(\mathbf{s}_0) = \sum_{k=0}^{\infty} E \left[ \mathbf{R}_t^{(12)}(\mathbf{s}_0) \mathbf{R}_{t+k}^{(12)}(\mathbf{s}_0)' \right]$ . By Condition (C4)(i) in Section 2, we have  $\Gamma_{\lambda_2}(\mathbf{s}_0) = E \left[ \mathbf{R}_t^{(12)}(\mathbf{s}_0) \mathbf{R}_t^{(12)}(\mathbf{s}_0)' \right] = \sigma_0^2 \Sigma_{\lambda_2}(\mathbf{s}_0)$ .

Further, let  $\mathbf{Z}_t^{*(21)}(\mathbf{s}_0) = \mathbf{Z}_t^{(21)}(\mathbf{s}_0) - E \left[ \mathbf{Z}_t^{(21)}(\mathbf{s}_0) | X_t(\mathbf{s}_0) \right]$  and  $\mathbf{R}_t^{(21)}(\mathbf{s}_0) = \varepsilon_t(\mathbf{s}_0) \mathbf{Z}_t^{*(21)}(\mathbf{s}_0)$ , with  $\mathbf{Z}_t^{(21)}(\mathbf{s}_0) = (Y_{t-1}(\mathbf{s}_0), Y_{t-2}(\mathbf{s}_0), \dots, Y_{t-q_0}(\mathbf{s}_0))'$ . Then, let  $\Sigma_{\alpha_1}(\mathbf{s}_0) = E \left[ \mathbf{Z}_t^{*(21)}(\mathbf{s}_0) \mathbf{Z}_t^{*(21)}(\mathbf{s}_0)' \right]$  denote a  $q_0 \times q_0$  matrix and  $\Gamma_{\alpha_1}(\mathbf{s}_0) = \sum_{k=0}^{\infty} E \left[ \mathbf{R}_t^{(21)}(\mathbf{s}_0) \mathbf{R}_{t+k}^{(21)}(\mathbf{s}_0)' \right]$ . By Condition (C4)(i) in Section 2, we have  $\Gamma_{\alpha_1}(\mathbf{s}_0) = E \left[ \mathbf{R}_t^{(21)}(\mathbf{s}_0) \mathbf{R}_t^{(21)}(\mathbf{s}_0)' \right] = \sigma_0^2 \Sigma_{\alpha_1}(\mathbf{s}_0)$ . Also, we let  $\Sigma_{\lambda_1 \alpha_1}(\mathbf{s}_0) = E \left[ \mathbf{Z}_t^{*(11)}(\mathbf{s}_0) \mathbf{Z}_t^{*(21)}(\mathbf{s}_0)' \right]$  denote a  $\sum_{i=1}^{p_0} n_i \times q_0$  matrix.

Finally, let  $\mathbf{Z}_t^{*(22)}(\mathbf{s}_0) = \mathbf{Z}_t^{(22)}(\mathbf{s}_0) - E \left[ \mathbf{Z}_t^{(22)}(\mathbf{s}_0) | X_t(\mathbf{s}_0) \right]$  and  $\mathbf{R}_t^{(22)}(\mathbf{s}_0) = \varepsilon_t(\mathbf{s}_0) \mathbf{Z}_t^{*(22)}(\mathbf{s}_0)$ , with  $\mathbf{Z}_t^{(22)}(\mathbf{s}_0) = (Y_{t-(q_0+1)}(\mathbf{s}_0), Y_{t-(q_0+2)}(\mathbf{s}_0), \dots, Y_{t-q}(\mathbf{s}_0))'$ . Then, let  $\Sigma_{\alpha_2}(\mathbf{s}_0) = E \left[ \mathbf{Z}_t^{*(22)}(\mathbf{s}_0) \mathbf{Z}_t^{*(22)}(\mathbf{s}_0)' \right]$  denote a  $(q - q_0) \times (q - q_0)$  matrix and  $\Gamma_{\alpha_2}(\mathbf{s}_0) = \sum_{k=0}^{\infty} E \left[ \mathbf{R}_t^{(22)}(\mathbf{s}_0) \mathbf{R}_{t+k}^{(22)}(\mathbf{s}_0)' \right]$ . By Condition (C4)(i)



in Section 2, we have  $\Gamma_{\alpha_2}(\mathbf{s}_0) = E \left[ \mathbf{R}_t^{(22)}(\mathbf{s}_0) \mathbf{R}_t^{(22)}(\mathbf{s}_0)' \right] = \sigma_0^2 \Sigma_{\alpha_2}(\mathbf{s}_0)$ .

## 2. Regularity Conditions

Now we list the regularity conditions for Theorems 1–4.

(C1) (i) The covariate and response process  $\{(X_t(\mathbf{s}_0), \mathbf{Y}_t(\mathbf{s}_0))\}$ , with  $\mathbf{Y}_t(\mathbf{s}_0) = (Y_t(\mathbf{s}_1), Y_t(\mathbf{s}_2), \dots, Y_t(\mathbf{s}_N))'$ , is strictly stationary and  $\alpha$ -mixing or (more generally) near epoch dependent based on an  $\alpha$ -mixing process (c.f. Lu and Linton, 2007; Li et al., 2012) in time and  $X_t(\mathbf{s}_0)$  has a compact support with the joint probability density function  $p(x_1, x_2; \mathbf{s}_0)$  of  $X_{t_1}(\mathbf{s}_0)$  and  $X_{t_2}(\mathbf{s}_0)$  being continuous and bounded from above for all  $t_1 \neq t_2$  and  $x_1, x_2 \in \mathbb{R}$ .

(ii) The  $\alpha$ -mixing coefficient  $\alpha(\cdot)$  satisfies  $\lim_{k \rightarrow \infty} k^a \sum_{n=k}^{\infty} \{\alpha(n)\}^{\delta/(2+\delta)} = 0$  for some constant  $a > \delta/(2 + \delta)$ .

(C2) The temporal lag interactions  $\alpha_{0,\ell}$ 's in (1) satisfy the stationarity condition that all the roots of  $1 - \sum_{\ell=1}^q \alpha_{0,\ell} z^\ell = 0$  are outside the unit circle.

(C3) The functions  $g_1(x, \mathbf{s}_0) = E[Y_t(\mathbf{s}_0) | X_t(\mathbf{s}_0) = x]$  and  $\mathbf{g}_2(x, \mathbf{s}_0) = E[\mathbf{Z}_t(\mathbf{s}_0) | X_t(\mathbf{s}_0) = x]$  are continuous at all  $x$  and twice differentiable.

(C4) (i) The innovations  $\{\varepsilon_t(\mathbf{s}_0)\}$  are i.i.d. For each  $t > q$ ,  $\varepsilon_t(\mathbf{s}_0)$  is independent of  $X_t(\mathbf{s}_0)$ ,  $\{Y_{t-i}(\mathbf{s}_k)\}_{k=1}^N$  for  $i = 1, \dots, p$ , and  $Y_{t-l}(\mathbf{s}_0)$  for  $l = 1, \dots, q$ .

(ii) For some  $\delta > 0$ ,  $E[|\varepsilon_t(\mathbf{s}_0)|^{2+\delta}] < \infty$ ,  $E[|Y_t(\mathbf{s}_0)|^{2+\delta}] < \infty$ , and  $E[||\mathbf{Z}_t(\mathbf{s}_0)||^{2+\delta}] < \infty$ .

(C5) The matrix  $E[\mathbf{Z}_t^*(\mathbf{s}_0) \mathbf{Z}_t^*(\mathbf{s}_0)']$  and the covariance matrix  $\Psi(\mathbf{s}_0)$  are positive definite, where  $\mathbf{Z}_t^*(\mathbf{s}_0)$  and  $\Psi(\mathbf{s}_0)$  are defined in Section 1.

(C6) (i) The kernel function  $K(\cdot)$  is symmetric, uniformly bounded by some constant, and integrable. Further,  $\int K(u) du = 1$ ,  $\int u K(u) du = 0$ , and  $\int u^2 K(u) du < \infty$ .

(ii)  $K(u)$  is Lipschitz continuous of order 1.

(iii)  $K(u)$  has an integrable second-order radial majorant (i.e.,  $Q^K(x) = \sup_{\|y\| \geq \|x\|} [\|y\|^2 K(y)]$  is integrable).

(C7) (i) The bandwidth  $b \rightarrow 0$  in such a way that  $T_0 b \rightarrow \infty$  and  $T_0 b^4 \rightarrow 0$  as  $T_0 \rightarrow \infty$ .

(ii) There exist two sequences of positive integers,  $\phi_{T_0} \rightarrow \infty$  and  $\varphi_{T_0} \rightarrow \infty$ , such that  $\varphi_{T_0}/\phi_{T_0} \rightarrow 0$  and  $T_0 \phi_{T_0}^{-1} \alpha(\varphi_{T_0}) \rightarrow 0$  as  $T_0 \rightarrow \infty$ .

(iii) The bandwidth  $b \rightarrow 0$  in such a way that  $\varphi_{T_0} b = O(1)$  and  $b^{-\delta/(2+\delta)} \sum_{t=\varphi_{T_0}}^{\infty} \alpha(t)^{\delta/(2+\delta)} \rightarrow 0$ , and  $\log(T_0)/(T_0^{1/2} b) \rightarrow 0$ , as  $T_0 \rightarrow \infty$ .

The above regularity conditions (C1)–(C7) are fairly mild. Condition (C1) assumes that the covariate process  $X_t(s_0)$  has smooth, bounded probability density functions and is  $\alpha$ -mixing over time (Fan and Yao, 2003, pp.68), which are quite standard in nonparametric time series analysis. We actually allow the process to be near epoch dependent, which takes  $\alpha$ -mixing as a special case and is of broad interest in econometrics (see, e.g., Lu and Linton, 2007; Li et al., 2012). The boundedness of  $X_t(s_0)$  is for simplicity of proof; alternatively, we could use truncation argument for  $X_t(s_0)$  as in, for example, (Gao et al., 2006), although the proof would be more tedious. (C2) is a stationarity condition assumed about the temporal lag interactions  $\alpha_{0,\ell}$ 's in (1), whereas (C3) assumes smoothness for the functions  $g_1$  and  $g_2$ . (C4) and (C5) impose conditions on the model regarding the innovations as well as  $Y_t(s)$  and  $Z_t(s)$ , which are mild. (C6) is a standard condition imposed on the kernel function  $K(\cdot)$  for time-series based methods. Condition (C7) is a requirement for the temporal bandwidth  $b = b_{T_0}$ .

Overall, the conditions imposed on the time series in (C1)–(C7) are fairly mild and commonly used in the literature (see, e.g., Fan and Yao, 2003; Gao, 2007). Also, owing to model (1) with condition (C2),  $Y_t(s_0)$  is not necessarily  $\alpha$ -mixing, but needs to follow near epoch dependence of geometrically decreasing stability coefficient (Lu and Linton, 2007; Li et al., 2012). Hence, the theorems here may not be attainable from the asymptotic results in the literature.

### 3. Proofs

#### 3.1. Proof of Theorem 1

Let  $c_{T_0}(\mathbf{s}_0) = T_0^{-1/2} + a_{T_0}^*(\mathbf{s}_0)$ , with  $a_{T_0}^*(\mathbf{s}_0)$  defined in Section 4, and  $\boldsymbol{\omega}(\mathbf{s}_0) = (\boldsymbol{v}(\mathbf{s}_0)', \boldsymbol{\nu}(\mathbf{s}_0)')'$ , where  $\boldsymbol{v}(\mathbf{s}_0) = ((v_1^k(\mathbf{s}_0), v_2^k(\mathbf{s}_0), \dots, v_p^k(\mathbf{s}_0)), k = 1, \dots, N) \in \mathbb{R}^{Np}$  and  $\boldsymbol{\nu}(\mathbf{s}_0) = (\nu_1(\mathbf{s}_0), \nu_2(\mathbf{s}_0), \dots, \nu_q(\mathbf{s}_0)) \in \mathbb{R}^q$ .

Let  $\mathfrak{B}(\mathbf{s}_0) = \{\boldsymbol{\eta}^0(\mathbf{s}_0) + c_{T_0}(\mathbf{s}_0)\boldsymbol{\omega}(\mathbf{s}_0) : \|\boldsymbol{\omega}(\mathbf{s}_0)\| \leq \delta\}$  denote a ball centered around  $\boldsymbol{\eta}^0(\mathbf{s}_0)$  with radius  $c_{T_0}(\mathbf{s}_0)\delta$ . Recall that, with the notation introduced prior to the theorems in Section 4, we let  $p_0 = \#\{1 \leq i \leq p : n_i(\mathbf{s}_0) \neq 0\}$ , and for notational convenience we suppose without loss of generality that the non-zero  $\lambda_{0k^*,i}$ 's are for those  $k^* = 1, \dots, n_i$ ,  $i = 1, \dots, p_0$ , with  $n_i = n_i(\mathbf{s}_0)$ . Then for  $\|\boldsymbol{\omega}(\mathbf{s}_0)\| = \delta$ , it follows from (9) that

$$\begin{aligned}
& Q(\boldsymbol{\eta}^0(\mathbf{s}_0) + c_{T_0}(\mathbf{s}_0)\boldsymbol{\omega}(\mathbf{s}_0)) - Q(\boldsymbol{\eta}^0(\mathbf{s}_0)) \\
& \geq L(\boldsymbol{\eta}^0(\mathbf{s}_0) + c_{T_0}(\mathbf{s}_0)\boldsymbol{\omega}(\mathbf{s}_0)) - L(\boldsymbol{\eta}^0(\mathbf{s}_0)) \\
& \quad + T_0 \sum_{i=1}^{p_0} \sum_{k^*=1}^{n_i} \gamma_i^{k^*}(\mathbf{s}_0) (|\lambda_{0k^*,i} + c_{T_0}(\mathbf{s}_0)v_i^{k^*}(\mathbf{s}_0)| - |\lambda_{0k^*,i}|) \\
& \quad + T_0 \sum_{l=1}^{q_0} \beta_l(\mathbf{s}_0) (|\alpha_{0l}^0 + c_{T_0}(\mathbf{s}_0)\nu_l(\mathbf{s}_0)| - |\alpha_{0l}^0|) \\
& \geq L(\boldsymbol{\eta}^0(\mathbf{s}_0) + c_{T_0}(\mathbf{s}_0)\boldsymbol{\omega}(\mathbf{s}_0)) - L(\boldsymbol{\eta}^0(\mathbf{s}_0)) - T_0 c_{T_0}(\mathbf{s}_0) \sum_{i=1}^{p_0} \sum_{k^*=1}^{n_i} \gamma_i^{k^*}(\mathbf{s}_0) |v_i^{k^*}(\mathbf{s}_0)| \\
& \quad - T_0 c_{T_0}(\mathbf{s}_0) \sum_{l=1}^{q_0} \beta_l(\mathbf{s}_0) |\nu_l(\mathbf{s}_0)| \\
& \geq L(\boldsymbol{\eta}^0(\mathbf{s}_0) + c_{T_0}(\mathbf{s}_0)\boldsymbol{\omega}(\mathbf{s}_0)) - L(\boldsymbol{\eta}^0(\mathbf{s}_0)) - T_0 c_{T_0}^2(\mathbf{s}_0) \sum_{i=1}^{p_0} n_i \delta - T_0 c_{T_0}^2(\mathbf{s}_0) q_0 \delta \\
& = L(\boldsymbol{\eta}^0(\mathbf{s}_0) + c_{T_0}(\mathbf{s}_0)\boldsymbol{\omega}(\mathbf{s}_0)) - L(\boldsymbol{\eta}^0(\mathbf{s}_0)) - T_0 c_{T_0}^2(\mathbf{s}_0) \left( \sum_{i=1}^{p_0} n_i + q_0 \right) \delta, \tag{a1}
\end{aligned}$$

where the second last inequality follows from the definitions of  $a_{T_0}^*(\mathbf{s}_0)$  in Section 4 and  $c_{T_0}(\mathbf{s}_0)$  at

522 the beginning of this proof. Moreover, by a Taylor's expansion

$$\begin{aligned}
& L(\boldsymbol{\eta}^0(\mathbf{s}_0) + c_{T_0}(\mathbf{s}_0)\boldsymbol{\omega}(\mathbf{s}_0)) - L(\boldsymbol{\eta}^0(\mathbf{s}_0)) \\
&= c_{T_0}(\mathbf{s}_0)\boldsymbol{\omega}(\mathbf{s}_0)' \frac{\partial L(\boldsymbol{\eta}^0(\mathbf{s}_0))}{\partial \boldsymbol{\eta}^0(\mathbf{s}_0)} + (1/2)c_{T_0}^2(\mathbf{s}_0)\boldsymbol{\omega}(\mathbf{s}_0)' \frac{\partial^2 L(\boldsymbol{\eta}^0(\mathbf{s}_0))}{\partial \boldsymbol{\eta}^0(\mathbf{s}_0) \partial \boldsymbol{\eta}^0(\mathbf{s}_0)'} \boldsymbol{\omega}(\mathbf{s}_0) \{1 + o_p(1)\} \\
&= \sum_{t=r+1}^T \left[ -2c_{T_0}(\mathbf{s}_0)\boldsymbol{\omega}(\mathbf{s}_0)' \hat{\mathbf{Z}}_t(\mathbf{s}_0) \left\{ \hat{Y}_t(\mathbf{s}_0) - \hat{\mathbf{Z}}_t(\mathbf{s}_0)' \boldsymbol{\eta}^0(\mathbf{s}_0) \right\} \right. \\
&\quad \left. + c_{T_0}^2(\mathbf{s}_0)\boldsymbol{\omega}(\mathbf{s}_0)' \hat{\mathbf{Z}}_t(\mathbf{s}_0) \hat{\mathbf{Z}}_t(\mathbf{s}_0)' \boldsymbol{\omega}(\mathbf{s}_0) \right] \{1 + o_p(1)\} \\
&= T_0 \left[ c_{T_0}^2(\mathbf{s}_0)\boldsymbol{\omega}(\mathbf{s}_0)' \left\{ (1/T_0) \sum_{t=r+1}^T \hat{\mathbf{Z}}_t(\mathbf{s}_0) \hat{\mathbf{Z}}_t(\mathbf{s}_0)' \right\} \boldsymbol{\omega}(\mathbf{s}_0) \right. \\
&\quad \left. - 2c_{T_0}(\mathbf{s}_0)\boldsymbol{\omega}(\mathbf{s}_0)' \left\{ (1/T_0) \sum_{t=r+1}^T \hat{\mathbf{Z}}_t(\mathbf{s}_0) \left( \hat{Y}_t(\mathbf{s}_0) - \hat{\mathbf{Z}}_t(\mathbf{s}_0)' \boldsymbol{\eta}^0(\mathbf{s}_0) \right) \right\} \right] \{1 + o_p(1)\}.
\end{aligned}$$

Thus

$$\begin{aligned}
Q(\boldsymbol{\eta}^0(\mathbf{s}_0) + c_{T_0}(\mathbf{s}_0)\boldsymbol{\omega}(\mathbf{s}_0)) - Q(\boldsymbol{\eta}^0(\mathbf{s}_0)) &\geq T_0 \left[ c_{T_0}^2(\mathbf{s}_0)\boldsymbol{\omega}(\mathbf{s}_0)' \left( \frac{1}{T_0} \sum_{t=r+1}^T \hat{\mathbf{Z}}_t(\mathbf{s}_0) \hat{\mathbf{Z}}_t(\mathbf{s}_0)' \right) \boldsymbol{\omega}(\mathbf{s}_0) \right. \\
&\quad \left. - 2c_{T_0}(\mathbf{s}_0)\boldsymbol{\omega}(\mathbf{s}_0)' \left( \frac{1}{T_0} \sum_{t=r+1}^T \hat{\mathbf{Z}}_t(\mathbf{s}_0) \left( \hat{Y}_t(\mathbf{s}_0) - \hat{\mathbf{Z}}_t(\mathbf{s}_0)' \boldsymbol{\eta}^0(\mathbf{s}_0) \right) \right) \right] \{1 + o_p(1)\}
\end{aligned}$$

$$-T_0 c_{T_0}^2(\mathbf{s}_0) \left( \sum_{i=1}^{p_0} n_i + q_0 \right) \delta \equiv A_1 + A_2 + A_3$$

For  $A_1$ , since  $\hat{\mathbf{Z}}_t(\mathbf{s}_0) = \mathbf{Z}_t(\mathbf{s}_0) - \hat{E}[\mathbf{Z}_t(\mathbf{s}_0)|X_t(\mathbf{s}_0)]$ , by adding and subtracting  $E[\mathbf{Z}_t(\mathbf{s}_0)|X_t(\mathbf{s}_0)]$ , we have

$$\begin{aligned}
\hat{\mathbf{Z}}_t(\mathbf{s}_0) &= \mathbf{Z}_t(\mathbf{s}_0) - E[\mathbf{Z}_t(\mathbf{s}_0)|X_t(\mathbf{s}_0)] + \left[ E[\mathbf{Z}_t(\mathbf{s}_0)|X_t(\mathbf{s}_0)] - \hat{E}[\mathbf{Z}_t(\mathbf{s}_0)|X_t(\mathbf{s}_0)] \right] \\
&= \mathbf{Z}_t^*(\mathbf{s}_0) + \left[ E[\mathbf{Z}_t(\mathbf{s}_0)|X_t(s)] - \hat{E}[\mathbf{Z}_t(\mathbf{s}_0)|X_t(\mathbf{s}_0)] \right]. \tag{a2}
\end{aligned}$$

Moreover, for  $A_2$ , since

$$\hat{Y}_t(\mathbf{s}_0) - \hat{\mathbf{Z}}_t(\mathbf{s}_0)' \boldsymbol{\eta}^0(\mathbf{s}_0) = \left[ Y_t(\mathbf{s}_0) - \hat{E}[Y_t(\mathbf{s}_0)|X_t(\mathbf{s}_0)] \right] - \left[ \mathbf{Z}_t(\mathbf{s}_0) - \hat{E}[\mathbf{Z}_t(\mathbf{s}_0)|X_t(\mathbf{s}_0)] \right]' \boldsymbol{\eta}^0(\mathbf{s}_0),$$

by adding and subtracting  $E[Y_t(\mathbf{s}_0)|X_t(\mathbf{s}_0)]$  and  $E[\mathbf{Z}_t(\mathbf{s}_0)|X_t(\mathbf{s}_0)]' \boldsymbol{\eta}^0(\mathbf{s}_0)$ , we have

$$\begin{aligned} \hat{Y}_t(\mathbf{s}_0) - \hat{\mathbf{Z}}_t(\mathbf{s}_0)' \boldsymbol{\eta}^0(\mathbf{s}_0) &= \varepsilon_t(\mathbf{s}_0) + \left[ E[Y_t(\mathbf{s}_0)|X_t(\mathbf{s}_0)] - \hat{E}[Y_t(\mathbf{s}_0)|X_t(\mathbf{s}_0)] \right] \\ &\quad - \left[ E[\mathbf{Z}_t(\mathbf{s}_0)|X_t(\mathbf{s}_0)] - \hat{E}[\mathbf{Z}_t(\mathbf{s}_0)|X_t(\mathbf{s}_0)] \right]' \boldsymbol{\eta}^0(\mathbf{s}_0), \end{aligned}$$

523 where  $\varepsilon_t(\mathbf{s}_0) = [Y_t(\mathbf{s}_0) - E[Y_t(\mathbf{s}_0)|X_t(\mathbf{s}_0)]] - [\mathbf{Z}_t(\mathbf{s}_0) - E[\mathbf{Z}_t(\mathbf{s}_0)|X_t(\mathbf{s}_0)]]' \boldsymbol{\eta}^0(\mathbf{s}_0)$ , and both

524  $\hat{E}[Y_t(\mathbf{s}_0)|X_t(\mathbf{s}_0)]$  and  $\hat{E}[\mathbf{Z}_t(\mathbf{s}_0)|X_t(\mathbf{s}_0)]$  are local linear estimators.

Define  $B_1 = (1/T_0) \sum_{t=r+1}^T \hat{\mathbf{Z}}_t(\mathbf{s}_0) \hat{\mathbf{Z}}_t(\mathbf{s}_0)'$ . By using the uniform consistency of local linear fitting under near epoch dependence in Li et al. (2012) and following the same argument of the proof of Theorem 1 of Al-Sulami et al. (2017), we have

$$B_1 \xrightarrow{P} \boldsymbol{\Sigma}(\mathbf{s}_0) = E[\mathbf{Z}_t^*(\mathbf{s}_0) \mathbf{Z}_t^{*'}(\mathbf{s}_0)], \quad (\text{a3})$$

and

$$A_1 = T_0 c_{T_0}^2(\mathbf{s}_0) \boldsymbol{\omega}(\mathbf{s}_0)' \boldsymbol{\Sigma}(\mathbf{s}_0) \boldsymbol{\omega}(\mathbf{s}_0) \{1 + o_P(1)\} = O_p(T_0 c_{T_0}^2(\mathbf{s}_0) \delta^2) \quad (\text{a4})$$

Also, let  $B_2 = (1/T_0) \sum_{t=r+1}^T \hat{\mathbf{Z}}_t(\mathbf{s}_0) \left( \hat{Y}_t(\mathbf{s}_0) - \hat{\mathbf{Z}}_t(\mathbf{s}_0)' \boldsymbol{\eta}^0(\mathbf{s}_0) \right)$ , then, owing to  $T_0 b^4 \rightarrow 0$  in Condition (C7)(i),  $B_2 = O_p(T_0^{-1/2})$ ,

and

$$A_2 = -2T_0 c_{T_0}(\mathbf{s}_0) \boldsymbol{\omega}(\mathbf{s}_0)' O_p(T_0^{-1/2}) = O_p(T_0^{1/2} c_{T_0}(\mathbf{s}_0) \delta) = O_p(T_0 c_{T_0}^2(\mathbf{s}_0) \delta). \quad (\text{a5})$$

Moreover, we have

$$A_3 = T_0 c_{T_0}^2(\mathbf{s}_0) \left( \sum_{i=1}^{p_0} n_i + q_0 \right) \delta = O_p \left( T_0 c_{T_0}^2(\mathbf{s}_0) \delta \right) \quad (\text{a6})$$

From (a4), (a5) and (a6),  $A_1$  is the largest term. Then, for any given  $\epsilon > 0$ , there exists a large constant  $\delta$  such that

$$P \left\{ \inf_{\|\boldsymbol{\omega}(\mathbf{s}_0)\|=\delta} Q \left( \boldsymbol{\eta}^0(\mathbf{s}_0) + c_{T_0}(\mathbf{s}_0) \boldsymbol{\omega}(\mathbf{s}_0) \right) > Q \left( \boldsymbol{\eta}^0(\mathbf{s}_0) \right) \right\} \geq 1 - \epsilon,$$

525 which means that with probability tending to one, there exists a minimizer at the ball  $\mathfrak{B}(\mathbf{s}_0) =$   
 526  $\{\boldsymbol{\eta}^0(\mathbf{s}_0) + c_{T_0}(\mathbf{s}_0) \boldsymbol{\omega}(\mathbf{s}_0) : \|\boldsymbol{\omega}(\mathbf{s}_0)\| \leq \delta\}$ , which is a global minimizer due to convexity of the  
 527 objective function  $Q(\boldsymbol{\eta}(\mathbf{s}_0))$ . Then  $\|\hat{\boldsymbol{\eta}}(\mathbf{s}_0) - \boldsymbol{\eta}^0(\mathbf{s}_0)\| = O_p \left( T_0^{-1/2} + a_{T_0}^*(\mathbf{s}_0) \right)$ .

### 528 3.2. Proof of Theorem 2

First, we show that  $\hat{\boldsymbol{\lambda}}_2(\mathbf{s}_0) \rightarrow \mathbf{0}$  with probability tending to one, where

$$\boldsymbol{\lambda}_2(\mathbf{s}_0) = ((\lambda_{0k^*,p_0+1}, \lambda_{0k^*,p_0+2}, \dots, \lambda_{0k^*,p}), k^* = 1, \dots, N; (\lambda_{0,n_i+1,i}, \lambda_{0,n_i+2,i}, \lambda_{0N,i}), i = 1, \dots, p_0)'$$

with regularization parameter

$$\begin{aligned} \boldsymbol{\gamma}_2(\mathbf{s}_0) = & \text{diag} \left( (\gamma_{p_0+1}^{k^*}(\mathbf{s}_0), \gamma_{p_0+2}^{k^*}(\mathbf{s}_0), \dots, \gamma_p^{k^*}(\mathbf{s}_0)), k^* = 1, \dots, N; \right. \\ & \left. (\gamma_i^{n_i+1}(\mathbf{s}_0), \gamma_i^{n_i+2}(\mathbf{s}_0), \dots, \gamma_i^N(\mathbf{s}_0)), i = 1, \dots, p_0 \right). \end{aligned}$$

It holds that

$$\frac{\partial Q(\hat{\boldsymbol{\eta}}(\mathbf{s}_0))}{\partial \boldsymbol{\lambda}_2(\mathbf{s}_0)} = \frac{\partial L(\hat{\boldsymbol{\eta}}(\mathbf{s}_0))}{\partial \boldsymbol{\lambda}_2(\mathbf{s}_0)} + T_0 \boldsymbol{\gamma}_2(\mathbf{s}_0) \text{sgn}(\boldsymbol{\lambda}_2(\mathbf{s}_0))$$

By a Taylor's expansion

$$\frac{\partial L(\hat{\boldsymbol{\eta}}(\mathbf{s}_0))}{\partial \boldsymbol{\lambda}_2(\mathbf{s}_0)} = \frac{\partial L(\boldsymbol{\eta}^0(\mathbf{s}_0))}{\partial \boldsymbol{\lambda}_2(\mathbf{s}_0)} + \frac{\partial^2 L(\boldsymbol{\eta}^0(\mathbf{s}_0))}{\partial \boldsymbol{\lambda}_2(\mathbf{s}_0) \partial \boldsymbol{\lambda}_2(\mathbf{s}_0)'} (\hat{\boldsymbol{\eta}}(\mathbf{s}_0) - \boldsymbol{\eta}^0(\mathbf{s}_0)) \{1 + o_p(1)\}$$

From the law of large number and a central limit theorem, it can be shown that the first term is of order  $O_p(T_0^{1/2})$  and  $T_0^{-1} \frac{\partial^2 L(\boldsymbol{\eta}^0(\mathbf{s}_0))}{\partial \boldsymbol{\lambda}_2(\mathbf{s}_0) \partial \boldsymbol{\lambda}_2(\mathbf{s}_0)'} \xrightarrow{P} \boldsymbol{\Sigma}_{\boldsymbol{\lambda}_2}(\mathbf{s}_0)$ .

Thus

$$\frac{\partial Q(\hat{\boldsymbol{\eta}}(\mathbf{s}_0))}{\partial \boldsymbol{\lambda}_2(\mathbf{s}_0)} = \frac{\partial L(\boldsymbol{\eta}^0(\mathbf{s}_0))}{\partial \boldsymbol{\lambda}_2(\mathbf{s}_0)} + T_0 \boldsymbol{\Sigma}_{\boldsymbol{\lambda}_2}(\mathbf{s}_0) (\hat{\boldsymbol{\eta}}(\mathbf{s}_0) - \boldsymbol{\eta}^0(\mathbf{s}_0)) \{1 + o_p(1)\} + T_0 \boldsymbol{\gamma}_2(\mathbf{s}_0) \text{sgn}(\boldsymbol{\lambda}_2(\mathbf{s}_0))$$

529 From Theorem 1, the second term is of order  $O_p(T_0^{1/2})$ . Moreover, the last term is of order  
 530  $T_0^{1/2} O_p(T_0^{1/2} d_{T_0}^*(\mathbf{s}_0))$ . Since  $\sqrt{T_0} d_{T_0}^*(\mathbf{s}_0) \rightarrow \infty$ , the last term is the largest term and hence the sign  
 531 of  $\frac{\partial Q(\hat{\boldsymbol{\eta}}(\mathbf{s}_0))}{\partial \boldsymbol{\lambda}_2(\mathbf{s}_0)}$  is determined by the sign of  $\boldsymbol{\lambda}_2(\mathbf{s}_0)$ . Thus, owing to  $\frac{\partial Q(\hat{\boldsymbol{\eta}}(\mathbf{s}_0))}{\partial \boldsymbol{\lambda}_2(\mathbf{s}_0)} = 0$ , we have  $\boldsymbol{\lambda}_2(\mathbf{s}_0) = \mathbf{0}$   
 532 with probability tending to 1.

533 By a similar argument, we can show that  $\alpha_{0l} = 0$  for  $l = q_0 + 1, q_0 + 2, \dots, q$ . Therefore,

534  $P(\hat{\boldsymbol{\lambda}}_2(\mathbf{s}_0) = \mathbf{0}) \rightarrow 1$  and  $P(\hat{\boldsymbol{\alpha}}_2(\mathbf{s}_0) = \mathbf{0}) \rightarrow 1$ .

### 535 3.3. Proof of Theorem 3

From Theorem 2, we have  $P(\hat{\boldsymbol{\lambda}}_2(\mathbf{s}_0) = \mathbf{0}) \rightarrow 1$  and  $P(\hat{\boldsymbol{\alpha}}_2(\mathbf{s}_0) = \mathbf{0}) \rightarrow 1$ , and hence,  $P(\hat{\boldsymbol{\eta}}_2(\mathbf{s}_0) = \mathbf{0}) \rightarrow 1$ . Thus, with probability tending to one, the minimizer of the objective function  $Q(\boldsymbol{\eta}(\mathbf{s}_0))$  is the same as the minimizer of  $Q(\boldsymbol{\eta}_1(\mathbf{s}_0))$ . Therefore,

$$\frac{\partial Q(\boldsymbol{\eta}_1(\mathbf{s}_0))}{\partial \boldsymbol{\eta}_1(\mathbf{s}_0)} \Big|_{\boldsymbol{\eta}_1(\mathbf{s}_0) = \hat{\boldsymbol{\eta}}_1(\mathbf{s}_0) = \mathbf{0}} = \mathbf{0}.$$

Recall  $\boldsymbol{\eta}_1(\mathbf{s}_0) = (\boldsymbol{\lambda}_1(\mathbf{s}_0)', \boldsymbol{\alpha}_1(\mathbf{s}_0)')'$  with  $\boldsymbol{\lambda}_1(\mathbf{s}_0) = ((\lambda_{01,i}, \lambda_{02,i}, \dots, \lambda_{n_i,i}), i = 1, \dots, p_0)'$  and  $\boldsymbol{\alpha}_1(\mathbf{s}_0) = (\alpha_{01}, \alpha_{02}, \dots, \alpha_{q_0})'$ . Let  $\boldsymbol{\zeta}_1(\mathbf{s}_0) = \text{diag}(\boldsymbol{\gamma}_1(\mathbf{s}_0)', \boldsymbol{\beta}_1(\mathbf{s}_0)')'$  be the regularization parameter matrix with  $\boldsymbol{\gamma}_1(\mathbf{s}_0) = \text{diag}((\gamma_i^1(\mathbf{s}_0), \gamma_i^2(\mathbf{s}_0), \dots, \gamma_i^{n_i}(\mathbf{s}_0)), i = 1, \dots, p_0)$  and  $\boldsymbol{\beta}_1(\mathbf{s}_0) = (\beta_1(\mathbf{s}_0), \beta_2(\mathbf{s}_0), \dots, \beta_{q_0}(\mathbf{s}_0))'$ . Then, we have

$$\mathbf{0} = \frac{\partial Q(\hat{\boldsymbol{\eta}}_1(\mathbf{s}_0))}{\partial \boldsymbol{\eta}_1(\mathbf{s}_0)} = \frac{\partial L(\hat{\boldsymbol{\eta}}_1(\mathbf{s}_0))}{\partial \boldsymbol{\eta}_1(\mathbf{s}_0)} + T_0 \boldsymbol{\zeta}_1(\mathbf{s}_0) \text{sgn}(\hat{\boldsymbol{\eta}}_1(\mathbf{s}_0)),$$

where  $\text{sgn}(\hat{\boldsymbol{\eta}}_1)$  is a point-wise sign function of vector  $\hat{\boldsymbol{\eta}}_1$ . Thus, by a Taylor's expansion,

$$\mathbf{0} = \frac{\partial L(\boldsymbol{\eta}_1^0(\mathbf{s}_0))}{\partial \boldsymbol{\eta}_1(\mathbf{s}_0)} + \frac{\partial^2 L(\boldsymbol{\eta}_1^0(\mathbf{s}_0))}{\partial \boldsymbol{\eta}_1(\mathbf{s}_0) \partial \boldsymbol{\eta}_1(\mathbf{s}_0)'} (\hat{\boldsymbol{\eta}}_1(\mathbf{s}_0) - \boldsymbol{\eta}_1^0(\mathbf{s}_0)) \{1 + o_p(1)\} + T_0 \boldsymbol{\zeta}_1(\mathbf{s}_0) \text{sgn}(\hat{\boldsymbol{\eta}}_1(\mathbf{s}_0)).$$

Further, it can be shown that

$$\frac{1}{\sqrt{T_0}} \frac{\partial L(\boldsymbol{\eta}_1^0(\mathbf{s}_0))}{\partial \boldsymbol{\eta}_1(\mathbf{s}_0)} \xrightarrow{D} N(\mathbf{0}, \boldsymbol{\Gamma}_{\boldsymbol{\eta}_1}(\mathbf{s}_0)) \text{ and } T_0^{-1} \frac{\partial^2 L(\boldsymbol{\eta}_1^0(\mathbf{s}_0))}{\partial \boldsymbol{\eta}_1(\mathbf{s}_0) \partial \boldsymbol{\eta}_1(\mathbf{s}_0)'} \xrightarrow{P} \boldsymbol{\Sigma}_{\boldsymbol{\eta}_1}(\mathbf{s}_0).$$

Therefore,

$$\mathbf{0} = \frac{1}{\sqrt{T_0}} \frac{\partial L(\boldsymbol{\eta}_1^0(\mathbf{s}_0))}{\partial \boldsymbol{\eta}_1(\mathbf{s}_0)} + \sqrt{T_0} \boldsymbol{\Sigma}_{\boldsymbol{\eta}_1}(\mathbf{s}_0) (\hat{\boldsymbol{\eta}}_1(\mathbf{s}_0) - \boldsymbol{\eta}_1^0(\mathbf{s}_0)) \{1 + o_p(1)\} + \sqrt{T_0} \boldsymbol{\zeta}_1(\mathbf{s}_0) \text{sgn}(\hat{\boldsymbol{\eta}}_1(\mathbf{s}_0)),$$

536 When  $T_0$  is large enough,  $\boldsymbol{\zeta}_1(\mathbf{s}_0) \text{sgn}(\hat{\boldsymbol{\eta}}_1(\mathbf{s}_0)) = \boldsymbol{\zeta}_1(\mathbf{s}_0) \text{sgn}(\boldsymbol{\eta}_1^0(\mathbf{s}_0))$  with probability tend-  
 537 ing to 1 by Theorem 2, and  $\sqrt{T_0} a_{T_0}^*(\mathbf{s}_0) \rightarrow 0$  by the condition of this theorem. Therefore,  
 538  $\sqrt{T_0} \boldsymbol{\zeta}_1(\mathbf{s}_0) \text{sgn}(\boldsymbol{\eta}_1^0(\mathbf{s}_0)) = o_p(1)$ . That is

$$\begin{aligned} \sqrt{T_0} (\hat{\boldsymbol{\eta}}_1(\mathbf{s}_0) - \boldsymbol{\eta}_1^0(\mathbf{s}_0)) &= \boldsymbol{\Sigma}_{\boldsymbol{\eta}_1}^{-1}(\mathbf{s}_0) \left( \frac{1}{\sqrt{T_0}} \frac{\partial L(\boldsymbol{\eta}_1^0(\mathbf{s}_0))}{\partial \boldsymbol{\eta}_1(\mathbf{s}_0)} \right) + o_p(1) \\ &\xrightarrow{D} N \left( \mathbf{0}, \boldsymbol{\Sigma}_{\boldsymbol{\eta}_1}^{-1}(\mathbf{s}_0) \boldsymbol{\Gamma}_{\boldsymbol{\eta}_1} \left( \boldsymbol{\Sigma}_{\boldsymbol{\eta}_1}^{-1}(\mathbf{s}_0) \right)' \right). \end{aligned}$$

539 By the notations in Section 1 and together with Condition (C4), we have the result of this  
 540 theorem.

### 541 3.4. Proof of Theorem 4

Since  $\hat{g}_0(x) = \hat{g}_1(x, \mathbf{s}_0) - \hat{\mathbf{g}}_2(x, \mathbf{s}_0)' \hat{\boldsymbol{\eta}}(\mathbf{s}_0)$  is the estimator of  $g_0(x) = g_1(x, \mathbf{s}_0) - \mathbf{g}_2(x, \mathbf{s}_0)' \boldsymbol{\eta}(\mathbf{s}_0)$ , we have

$$\hat{g}_0(x) - g_0(x) = \{\hat{g}_1(x, \mathbf{s}_0) - g_1(x, \mathbf{s}_0)\} - \{\hat{\mathbf{g}}_2(x, \mathbf{s}_0) - \mathbf{g}_2(x, \mathbf{s}_0)\}' \boldsymbol{\eta}(\mathbf{s}_0) - \hat{\mathbf{g}}_2(x, \mathbf{s}_0)' \{\hat{\boldsymbol{\eta}}(\mathbf{s}_0) - \boldsymbol{\eta}(\mathbf{s}_0)\}.$$



From Theorem 1, we have  $\sqrt{T_0} \{\hat{\boldsymbol{\eta}}(\mathbf{s}_0) - \boldsymbol{\eta}(\mathbf{s}_0)\} = O_p(1)$ . Thus,

$$\sqrt{T_0 b} \hat{\mathbf{g}}_2(x, \mathbf{s}_0)' \{\hat{\boldsymbol{\eta}}(\mathbf{s}_0) - \boldsymbol{\eta}(\mathbf{s}_0)\} = O_p(\sqrt{b}) = o_p(1).$$

To establish the asymptotic normality of the estimate of  $g_0(x)$ , it suffices to establish the asymptotic normality of  $\hat{g}_1(x, \mathbf{s}_0) - g_1(x, \mathbf{s}_0)$  and  $\hat{\mathbf{g}}_2(x, \mathbf{s}_0) - \mathbf{g}_2(x, \mathbf{s}_0)$ , which follows from an argument similar to the proof of Theorem 2 of Al-Sulami et al. (2017). The detail is omitted.

#### 4. Simulation Study

We conduct a simulation study to evaluate our penalized procedure for the identification and estimation of important spatio-temporal lag interactions. We consider the spatio-temporal model

$$Y_t(\mathbf{s}_0) = g_0(X_{t-1}(\mathbf{s}_0)) + \sum_{i=1}^p \sum_{k=1}^N \lambda_{0k,i} Y_{t-i}(\mathbf{s}_k) + \sum_{l=1}^q \alpha_{0,l} Y_{t-l}(\mathbf{s}_0) + \varepsilon_t(\mathbf{s}_0), \quad (\text{a7})$$

where the exogenous explanatory variable  $X_t(\mathbf{s}_0)$  follows an AR(1) process. That is,  $X_t(\mathbf{s}_0) = 0.5X_{t-1}(\mathbf{s}_0) + e_t$ , where  $e_t$ 's follow i.i.d. standard normal distribution  $N(0, 1)$  and are independent of the innovations  $\varepsilon_t(\mathbf{s}_0)$ . The error terms in (a7)  $\varepsilon_t(\mathbf{s}_0)$ 's are i.i.d. following a normal distribution with mean 0 and variance  $\sigma^2$ . The variance  $\sigma^2$  is set to variance estimate in the housing price data example. Further, we assume the following nonlinear form for  $g_0(\cdot)$ :

$$g_0(X_{t-1}(\mathbf{s}_0)) = \log \left[ 1 + \left\{ (b(\mathbf{s}_0) + X_{t-1}(\mathbf{s}_0))^2 \right\}^{a(\mathbf{s}_0)} \right],$$

where  $a(\mathbf{s}_0) = 0.5 + 0.2 \cos(u_0 + v_0)$  and  $b(\mathbf{s}_0) = 0.6 + 0.3 \sin(u_0 v_0)$ , for  $\mathbf{s}_0 = (u_0, v_0) \in \mathbb{R}^2$ .

We also follow the set up of the housing price data example in Section 5. That is, there are  $N = 51$  sampling locations (50 states and DC),  $p = 6$  spatio-temporal lags, and  $q = 6$  temporal lags. In addition, we set the values of the spatio-temporal lag interactions  $\lambda_{0k,i}$  and the temporal lag interactions  $\alpha_{0,l}$ 's to be the estimated values in the case study with a slight modification such

that, an interaction is set to zero if the absolute value of the estimated interaction is smaller than 0.05.

We generate 100 Monte Carlo samples based on model (a7). The initial values of the response variable are set to zero and the simulated data at the first 50 time points are discarded as a warm-up step to reach stationarity over time. We take the simulated data at the remaining time points as a Monte Carlo sample and denote them as  $(X_t(\mathbf{s}_0), Y_t(\mathbf{s}_0))$  for  $t = 1, \dots, T$  and  $j = 1, \dots, N$ , where  $N = 51$  and we consider two lengths of time  $T = 350$  and  $T = 500$ .

For each Monte Carlo sample, we apply the penalized procedure to identify and estimate the lag interactions, as well as the unknown function  $g_0(\cdot)$ . In particular, for  $g_0(\cdot)$ , the estimation is done at 200 points between the 10th and 90th quantiles of the covariate  $X_{t-1}(\mathbf{s}_0)$ . Both the temporal bandwidth  $b$  and the spatial bandwidth  $h$  are selected by cross validation as in Section 5.

To assess the performance of the penalized procedure, we compute a proportion of correctly estimated zero lag interactions among the zero lag interactions and a proportion of the correctly estimated non-zero lag interactions among non-zero lag interactions. We illustrate the results by DC. With  $T = 350$ , the proportions for the 6 temporal lags are 93.69%, 92.70%, 96.56%, 95.35%, 96.71% and 96.43%, when the true temporal lag interactions are zero. With  $T = 500$ , these proportions increase to 98.35%, 97.86%, 96.90%, 97.63%, 97.84%, and 98.88%. Since the true interactions for lags 5 and 6 are zero (Figure 5(a)), we focus on the non-zero true interactions for the first four lags. With  $T = 350$ , the proportions of the correctly estimated non-zero coefficients are 66.50%, 86%, 21%, and 17.50% for lags 1–4, respectively. With  $T = 500$ , these proportions become 62.50%, 100%, 84%, and 10% for lags 1–4, respectively. The low proportion for lag 4 is not surprising, given the relatively small estimated value seen in Figure 5(a) and thus weaker signal to be identified compared to the first three lags.

## 5. Estimated $\lambda_{0k,i}$ for all States in the US

We provide here additional results on the estimated interactions  $\hat{\lambda}_{0k,i}$  for each of the  $N = 51$  spatial locations under model (14). In each of the heat maps below, the  $y$ -axis is for the  $j$ th spatial location and  $x$ -axis is for the  $k$ th spatial location interacting with the  $j$ th spatial location, whereas the color corresponds to estimated spatio-temporal lag interactions  $\hat{\lambda}_{jk,i}$  between the  $k$ th and the  $j$ th spatial locations, with  $j, k = 1, \dots, 51$ . Figure A1–A6 correspond to temporal lags  $i = 1, \dots, 6$ , respectively.

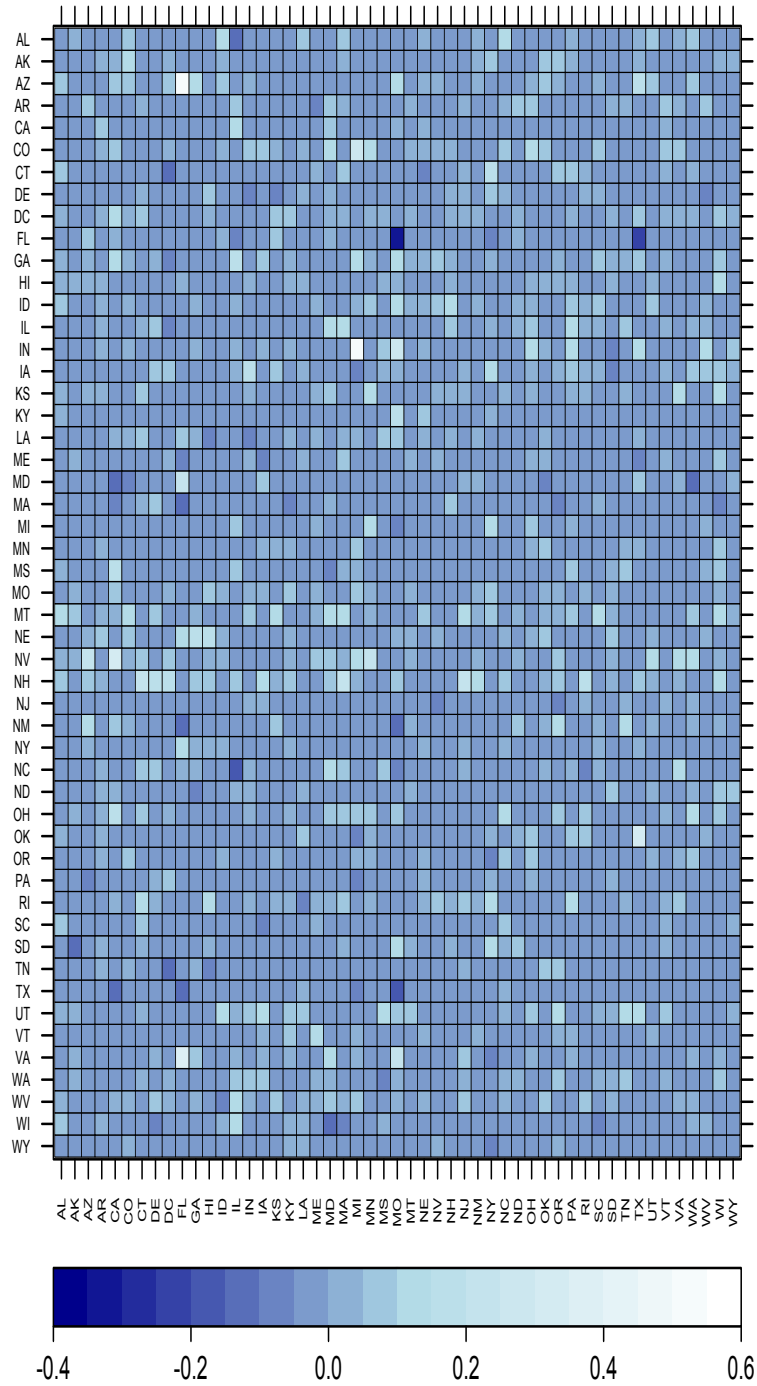


FIGURE A1: Heat map of the estimated spatio-temporal lag interactions  $\hat{\lambda}_{jk,1}$  for  $j, k = 1, \dots, 51$ , at lag  $i = 1$

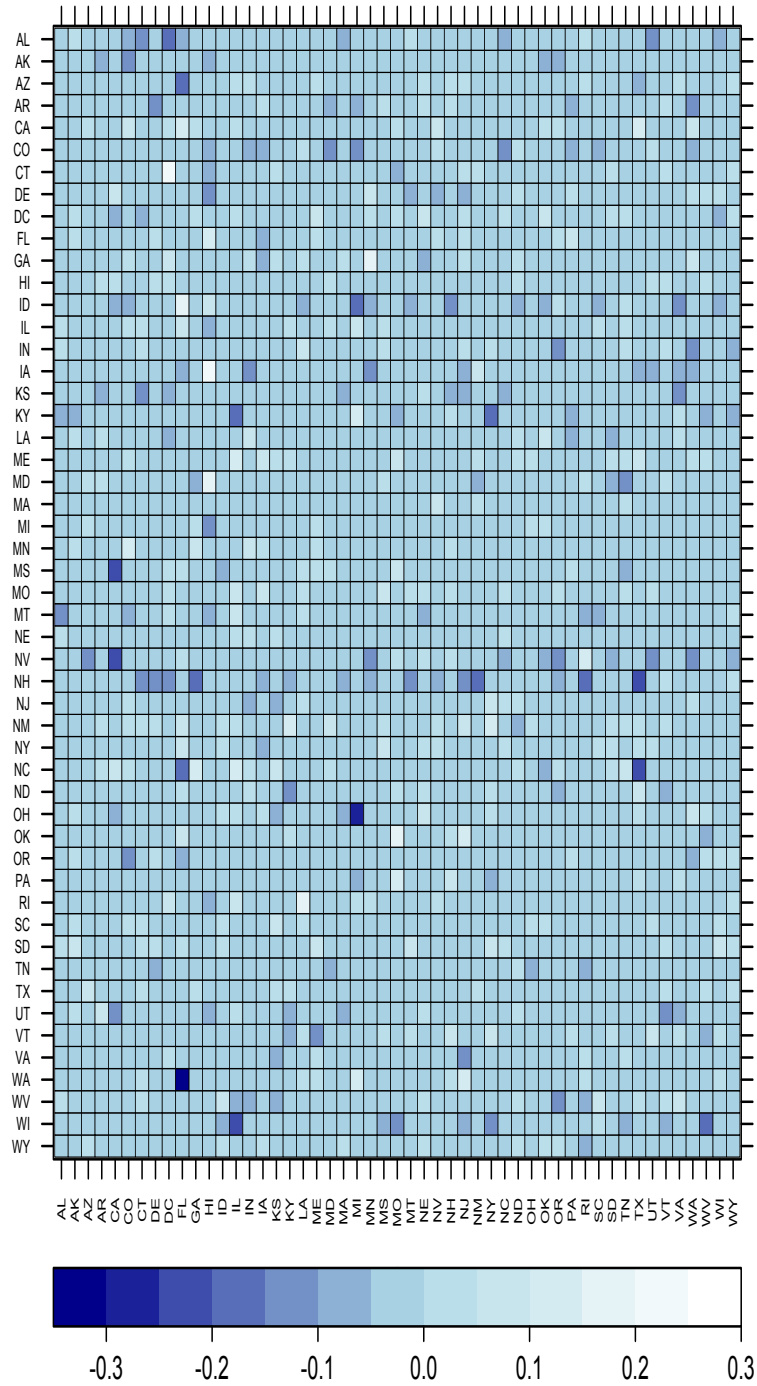


FIGURE A2: Heat map of the estimated spatio-temporal lag interactions  $\hat{\lambda}_{jk,2}$  for  $j, k = 1, \dots, 51$ , at lag  $i = 2$

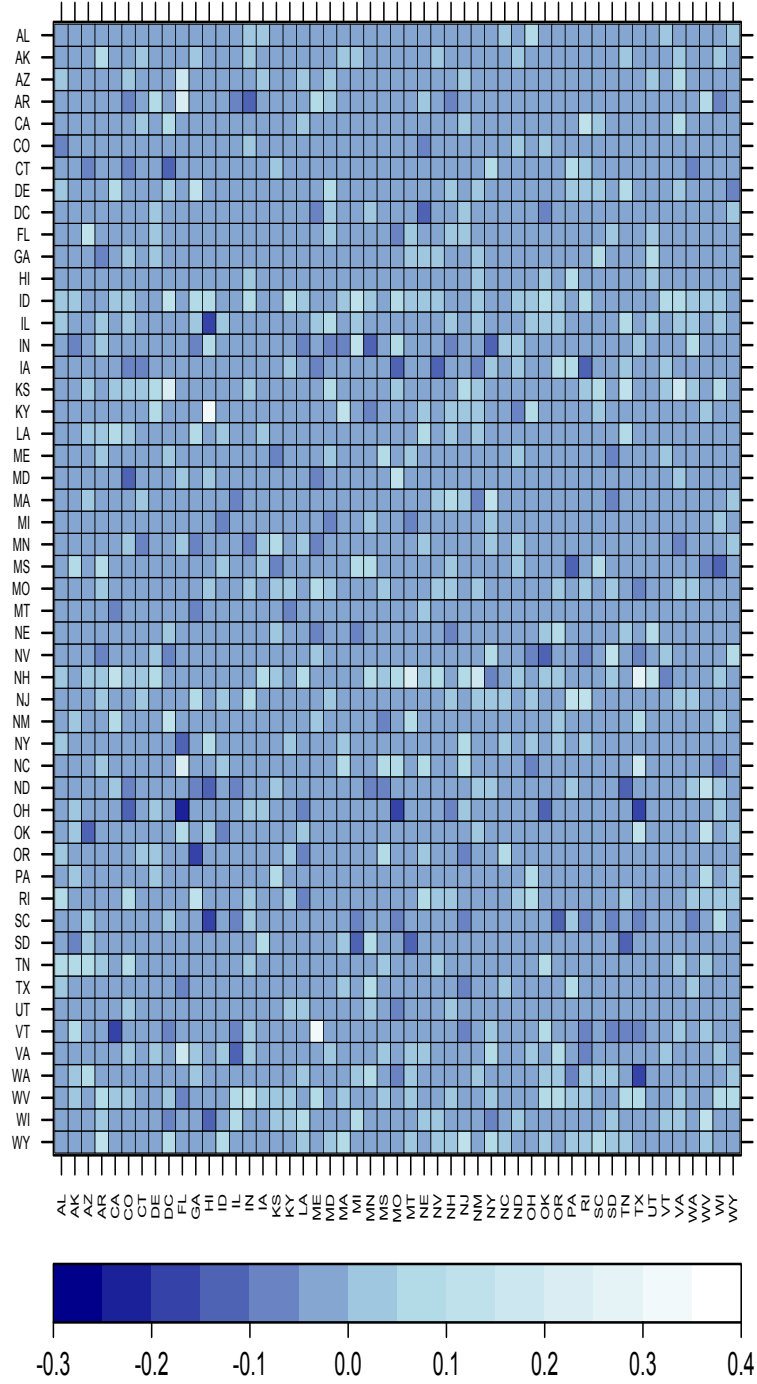


FIGURE A3: Heat map of the estimated spatio-temporal lag interactions  $\hat{\lambda}_{jk,3}$  for  $j, k = 1, \dots, 51$ , at lag  $i = 3$

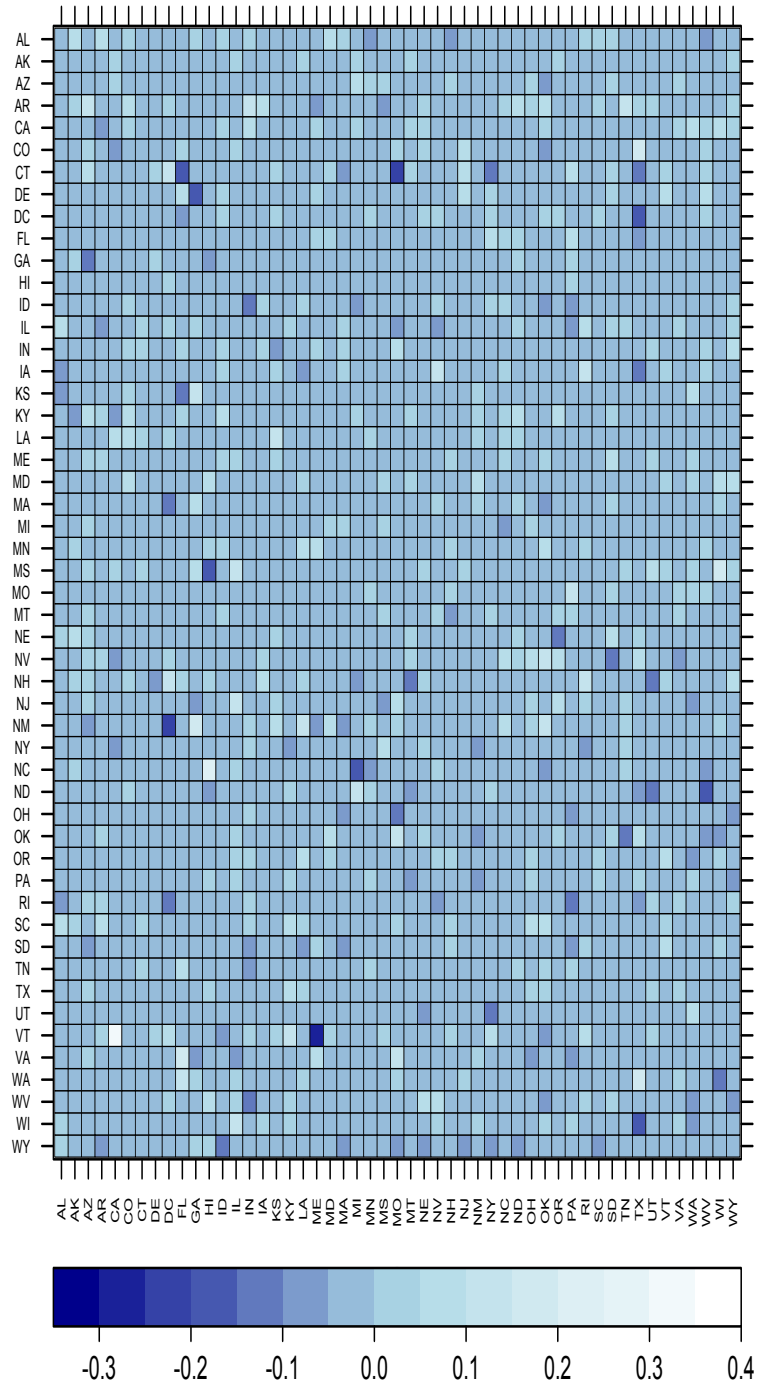


FIGURE A4: Heat map of the estimated spatio-temporal lag interactions  $\hat{\lambda}_{jk,4}$  for  $j, k = 1, \dots, 51$ , at lag  $i = 4$

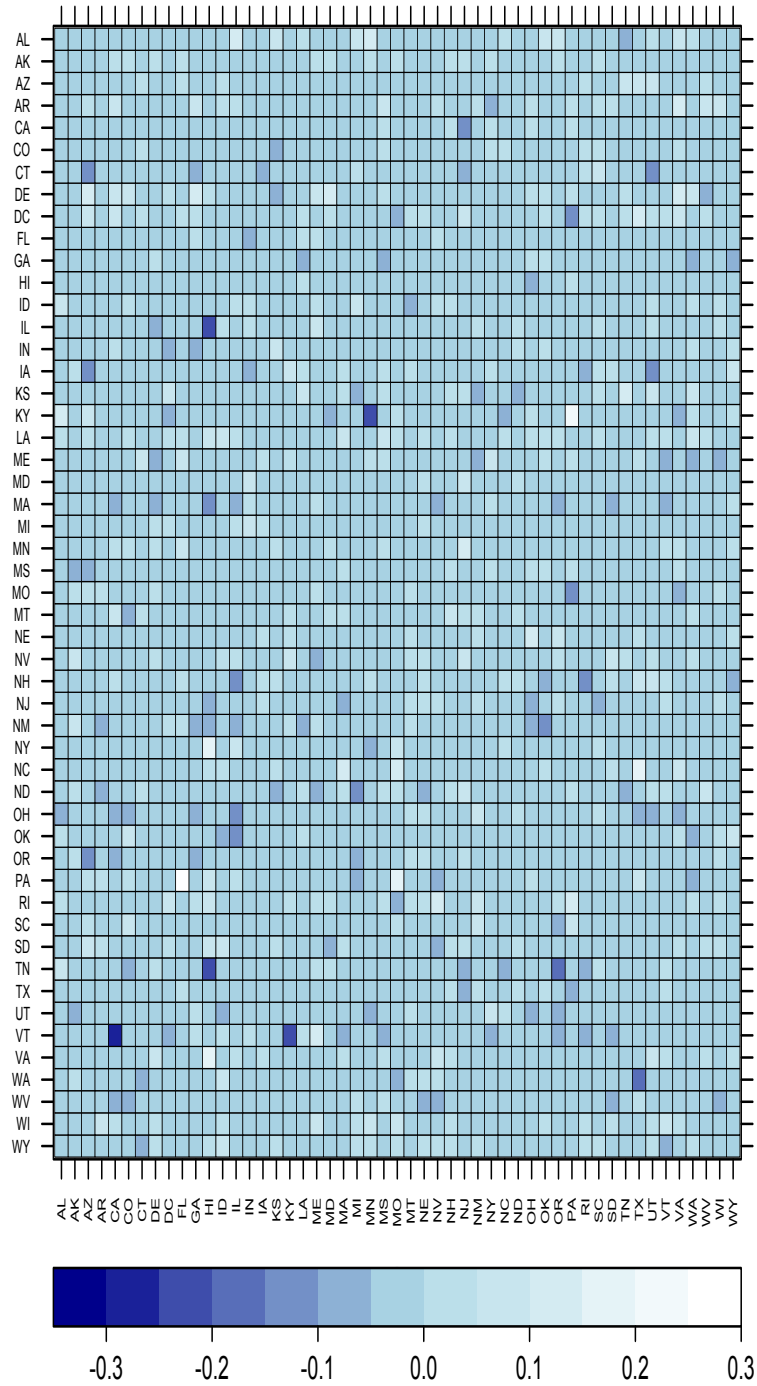


FIGURE A5: Heat map of the estimated spatio-temporal lag interactions  $\hat{\lambda}_{jk,5}$  for  $j, k = 1, \dots, 51$ , at lag  $i = 5$



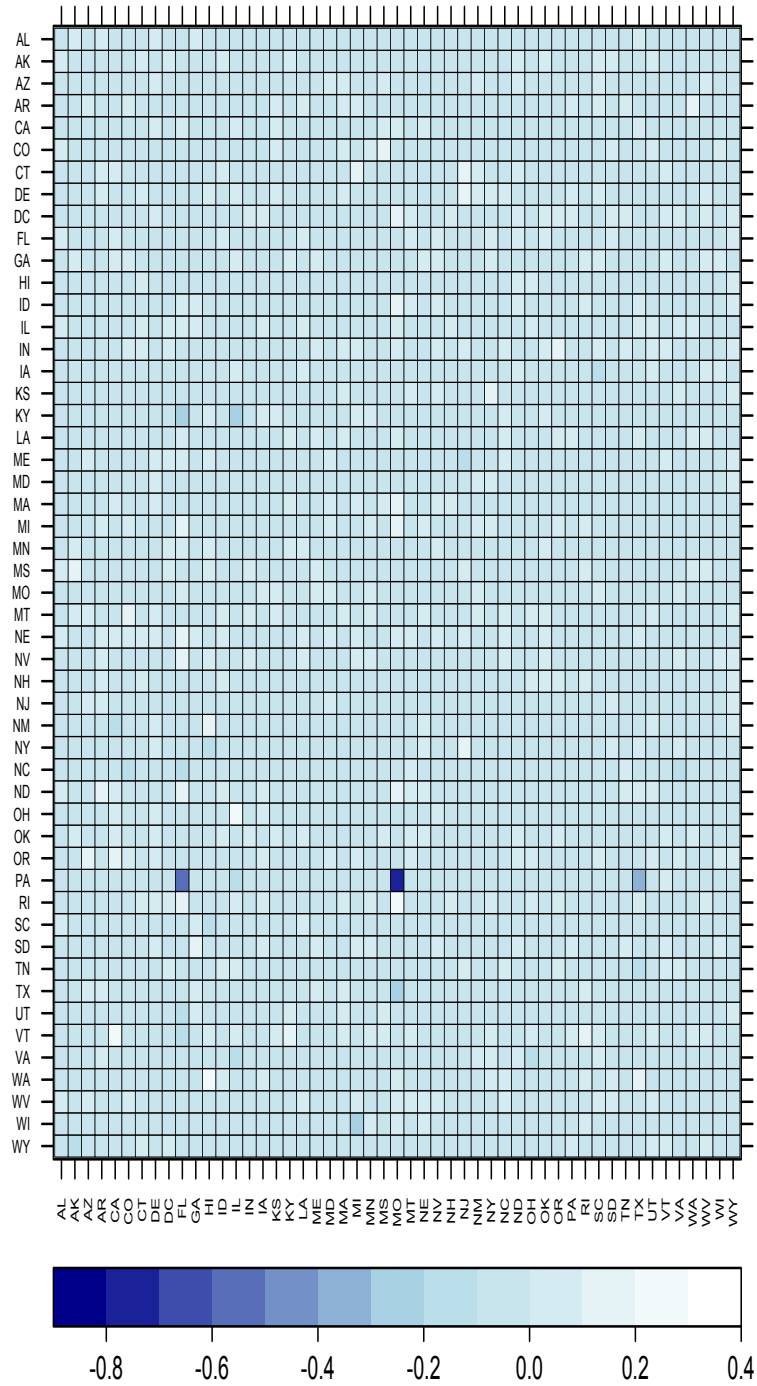


FIGURE A6: Heat map of the estimated spatio-temporal lag interactions  $\hat{\lambda}_{jk,6}$  for  $j, k = 1, \dots, 51$ , at lag  $i = 6$