



# Sparse generalized Yule–Walker estimation for large spatio-temporal autoregressions with an application to NO<sub>2</sub> satellite data

Hanno Reuvers <sup>a,b</sup>, Etienne Wijler <sup>c,\*</sup>

<sup>a</sup> Department of Econometrics, Erasmus University Rotterdam, 3062 PA Rotterdam, The Netherlands

<sup>b</sup> Iconsulting S.p.A., 00185 Rome, Italy

<sup>c</sup> Department of Econometrics and Data Science, Vrije Universiteit Amsterdam, 1081 HV Amsterdam, The Netherlands

## ARTICLE INFO

### JEL classification:

C33  
C53  
C55

### Keywords:

Spatio-temporal models  
SPLASH  
Satellite data  
Yule–Walker  
High-dimensional

## ABSTRACT

We consider a high-dimensional model in which variables are observed over time and space. The model consists of a spatio-temporal regression containing a time lag and a spatial lag of the dependent variable. Unlike classical spatial autoregressive models, we do not rely on a predetermined spatial interaction matrix, but infer all spatial interactions from the data. Assuming sparsity, we estimate the spatial and temporal dependence fully data-driven by penalizing a set of Yule–Walker equations. This regularization can be left unstructured, but we also propose customized shrinkage procedures when observations originate from spatial grids (e.g. satellite images). Finite sample error bounds are derived and estimation consistency is established in an asymptotic framework wherein the sample size and the number of spatial units diverge jointly. Exogenous variables can be included as well. A simulation exercise shows strong finite sample performance compared to competing procedures. As an empirical application, we model satellite measured nitrogen dioxide (NO<sub>2</sub>) concentrations in London. Our approach delivers forecast improvements over a competitive benchmark and we discover evidence for strong spatial interactions.

## 1. Introduction

In this paper, we propose the SPatial LAsso-type SHrinkage (SPLASH) estimator: a fully data-driven, sparse estimator for large spatio-temporal models. A unique feature of our estimator is its capability to provide interpretable estimates of spatial interactions between many spatial units in a fully data-driven way, while simultaneously tackling important econometric challenges, such as inherent endogeneity problems and overparameterization of the model. SPLASH specializes in, without being limited to, modelling outcomes of spatial units that are ordered on a regularly spaced spatial grid. Relevant examples of such settings include the dynamic modelling of political preferences across voting districts, crime statistics across municipalities, and satellite-measured air pollutants. Illustrating the latter, we apply our estimator to predict NO<sub>2</sub> concentrations in Greater London based on daily satellite images.

Spatio-temporal models are powerful tools to explain and exploit dependencies between variables that are observed over both time and space, but they come with a number of challenges. In particular, endogeneity issues arise because the contemporaneous observations occur on both sides of the model equation. Furthermore, the inclusion of both spatial and temporal lags quickly results

\* Correspondence to: Department of Econometrics and Data Science, School of Business and Economics, Vrije Universiteit Amsterdam, 1081 HV Amsterdam, The Netherlands.

E-mail address: [e.j.wijler@vu.nl](mailto:e.j.wijler@vu.nl) (E. Wijler).

<https://doi.org/10.1016/j.jeconom.2023.105520>

Received 14 December 2021; Received in revised form 27 June 2023; Accepted 9 August 2023

Available online 21 October 2023

0304-4076/© 2023 The Author(s).

Published by Elsevier B.V. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

Published by Elsevier B.V. This is an open access article under the CC BY license

in heavily parameterized models. To circumvent these issues, a large part of the literature incorporates predetermined spatial weight matrices that govern the contemporaneous interactions between spatial units. Examples of this modelling strategy are: the spatial autoregressive model with a Gaussian quasi-maximum likelihood estimator (QMLE) (Lee, 2004); the QMLE estimation of stationary spatial panels with fixed effects detailed (Yu et al., 2008); an extension of these spatial panels to include spatially autoregressive disturbances (Lee and Yu, 2010); a further extension to a non-stationary setting in which units can be spatially cointegrated (Yu et al., 2012); and a computationally beneficial generalized method of moments (GMM) estimator (Lee and Yu, 2014). In some settings, however, the specification of a pre-determined spatial weight matrix is a non-trivial, and potentially error-prone, step in the model building process. Accordingly, recent literature focuses on either incorporating multiple weight matrices (e.g. Debarsy and LeSage, 2018; Zhang and Yu, 2018) or, at the expense of estimating many parameters, directly inferring all spatial interactions from the data (e.g. Lam and Souza, 2019; Gao et al., 2019; Ma et al., 2023). Our model is a member of the latter category.

SPLASH offers the flexibility of a data-driven spatial weight matrix, without suffering from endogeneity or overparameterization. Apart from the assumption of a generous bandwidth, SPLASH leaves the spatial weight matrix and autoregressive matrix unspecified, while employing regularization to estimate sparse solutions. We show how dependencies between neighbouring units on a spatial grid translate to diagonally structured patterns of sparsity in the spatial weight matrix. Accordingly, SPLASH incorporates a tailored regularization component to exploit such structured sparsity when present, but without imposing it from the outset. Furthermore, SPLASH does not require the use of external instruments to correct for endogeneity. Building upon previous works by Dou et al. (2016) and Gao et al. (2019), we utilize the generalized Yule–Walker equations to resolve the endogeneity problem. Further generalizing the latter contributions, we extend the system of Yule–Walker equations to allow for the inclusion of exogenous variables.

Our work is related to the literature on the estimation of (nonparametric) spatial weight matrices and banded spatio-temporal vector autoregressions. We highlight a number of key contributions along with differences and similarities with respect to the current literature. Lam and Souza (2014) consider a model specification where the spatial units depend linearly on a spatial lag and exogenous regressors and show that the adaptive lasso can consistently select the correct sparsity pattern. To solve the endogeneity issue, they require the error variance to decay to zero as the time dimension grows large. Ahrens and Bhattacharjee (2015) solve the endogeneity problem using external instruments. Their two-step lasso estimation procedure selects the relevant instruments in the first step and the relevant spatial interactions in the second step. This approach, however, requires the instruments and the idiosyncratic component to be serially independent, resulting in limited applicability to spatio-temporal models. Finally, Lam and Souza (2019) augment a spatial lag model with a set of potentially endogenous variables. They decompose the spatial weight matrix into a pre-determined component based on expert knowledge and a sparse adjustment matrix that represents specification errors. The adjustment matrix is sparsely estimated based on a penalized version of instrumental variables (IV) regression. In contrast to our approach, Lam and Souza (2019) do not regularize the interactions between the dependent variables and the variables in the augmenting set, and they assume the number of such interactions to be fixed. A fixed number of interactions is inappropriate in high-dimensional settings in which the number of spatial units diverges.

Most closely related to our work are several recent contributions on banded estimation of spatio-temporal VARs. First, Gao et al. (2019), and Ma et al. (2023) consider the same model that appears in this work, and they also rely on the generalized Yule–Walker equations for estimation. The key difference with our paper lies in the method by which the model complexity is controlled during estimation. Gao et al. (2019) assume the coefficient matrices to be banded with a bandwidth that is small compared to the number of spatial units. The bandwidth is determined from the data and all parameters within the selected bandwidth are left unregularized. Our SPLASH estimator, however, has the ability to exploit (structured) sparsity within the bandwidth and thereby improve estimation and forecasting performance. In addition, apart from a generous upper bound on the bandwidth to ensure identification, SPLASH does not require an a priori choice regarding the bandwidth. The recently developed bagging approach by Ma et al. (2023) does allow for sparsity within the bands, yet it also requires the calculation of so-called solution paths. That is, a forward addition and backward deletion stage are needed to determine the variables that enter the final model specification. In contrast, the SPLASH estimator provides this solution at once. Furthermore, their approach is not designed to detect diagonally structured forms of sparsity, while the ability to do so results in clear performance improvements of SPLASH in both the simulations and the empirical application documented below. Finally, Wang and Tsay (2023) consider constrained estimation of the Yule–Walker equations for, potentially misspecified, high-dimensional VAR models. As part of their general framework, they consider banded VARs as a special case. Their model, however, does not include a contemporaneous (spatial) lag, nor is their estimator designed to uncover the type of (structured) sparsity that is essential to our application.

Several theoretical contributions are put forth in this paper. First, since SPLASH is based on the generalized Yule–Walker equations, which require estimates of the population autocovariance matrices, we derive a novel convergence result for banded estimates of the autocovariance matrices. Second, we derive finite-sample performance bounds for the estimation and prediction error of our estimator and utilize these bounds to derive asymptotic consistency of SPLASH in a variety of settings. For example, in the case of a finitely bounded bandwidth and unstructured sparsity, it follows that the number of spatial units  $N$  may grow at any polynomial rate of the number of temporal observations  $T$ . Finally, we develop a strategy to extend the system of Yule–Walker equations to accommodate additional exogenous regressors and we derive similar performance bounds of SPLASH in this broader framework as well.

We document strong performance of SPLASH in terms of estimation and forecast accuracy, both in simulated settings and in an empirical application. The simulation results highlight that SPLASH obtains highly competitive predictive accuracy even when compared to correctly-specified lattice-based methods. Furthermore, despite the heavier parametrization, the average estimation errors attained by SPLASH are among the lowest and the estimates are demonstrated to allow for straightforward interpretation

upon visualization. In our empirical application, we collect daily NO<sub>2</sub> column densities from 1 August 2018 to 31 March 2023, recorded by the TROPOspheric Monitoring Instrument (TROPOMI) on board of the Copernicus Sentinel-5 Precursor satellite. Each spatial unit represents an aggregation of a small number of pixels on the satellite image. We again find that SPLASH is a competitive alternative to popular lattice-based methods in terms of forecasting performance, and that spatial-temporal methods in general deliver significant forecast improvements over regularized estimation of a reduced form VAR. In addition, we find evidence for spatial interactions between first-order and second-order neighbours (i.e. neighbours of neighbours). By visualizing these spatial interactions, we shed light on the dominant flow vectors of NO<sub>2</sub> which can aid the development of regional NO<sub>2</sub> containment policies.

This paper is organized as follows. In Section 2 we introduce the spatio-temporal vector autoregressive model and discuss its stability and sparsity properties. Next, we develop the sparse estimation strategy of this model in Section 3, which includes banded covariance matrix estimation (3.1) and the SPLASH estimator (3.2). The simulation results are described in Section 4, followed by the empirical application in Section 5. We conclude in Section 6.

**Notation**

The indicator function  $\mathbb{1}_{\{A\}}$  equals 1 if  $A$  is true and zero otherwise. For a vector  $\mathbf{x} \in \mathbb{R}^N$ , the  $L_p$ -norm of  $\mathbf{x}$  is denoted  $\|\mathbf{x}\|_p = (\sum_{i=1}^N |x_i|^p)^{1/p}$ , with  $\|\mathbf{x}\|_\infty = \max_i |x_i|$  as an important special case. The total number elements in  $\mathbf{x}$  is denoted by  $|\mathbf{x}|$  and the number of non-zero elements in  $\mathbf{x}$  is denoted by  $\mathcal{M}(\mathbf{x}) = \sum_{i=1}^N \mathbb{1}\{x_i \neq 0\}$ . The Orlicz norm is defined as  $\|\cdot\|_\psi = \inf \{c > 0 : \mathbb{E}[\psi(|\cdot|/c)] \leq 1\}$  for any  $\psi(\cdot) : \mathbb{R}^+ \rightarrow \mathbb{R}^+$  being a convex, increasing function with  $\psi(0) = 0$  and  $\psi(x) \rightarrow \infty$  as  $x \rightarrow \infty$ . In addition, we rely on several matrix norms. For a matrix  $\mathbf{A} \in \mathbb{R}^{M \times N}$ , the matrix norms induced by the vector  $L_p$ -norms are given by  $\|\mathbf{A}\|_p = \sup_{\mathbf{x} \in \mathbb{R}^M} (\|\mathbf{A}\mathbf{x}\|_p / \|\mathbf{x}\|_p)$ . Noteworthy examples are:  $\|\mathbf{A}\|_1 = \max_{1 \leq j \leq N} \sum_{i=1}^M |a_{ij}|$ , the spectral norm  $\|\mathbf{A}\|_2 = [\lambda_{\max}(\mathbf{A}'\mathbf{A})]^{1/2}$  where  $\lambda_{\max}(\cdot)$  stands for the maximum eigenvalue, and  $\|\mathbf{A}\|_\infty = \max_{1 \leq i \leq M} \sum_{j=1}^N |a_{ij}|$ . Finally, we define  $\|\mathbf{A}\|_{\cdot} = \max\{\|\mathbf{A}\|_1, \|\mathbf{A}\|_\infty\}$ . Let  $S \subseteq \{1, \dots, N\}$  denote an index set with cardinality  $|S|$ . Then,  $\mathbf{x}_S$  denotes the  $|S|$ -dimensional vector with the elements of  $\mathbf{x}$  indexed by  $S$ , whereas  $\mathbf{A}_S$  denotes the  $(M \times |S|)$ -dimensional matrix containing the columns of  $\mathbf{A}$  indexed by  $S$ . In addition, we define  $\mathcal{D}_A(k) = \{a_{ij} \mid |i - j| = k\}$  as the collections of elements lying on (pairs of) the diagonals in the matrix  $\mathbf{A}$ , and  $\mathcal{B}_h(\mathbf{A}) = (a_{ij} \mathbb{1}_{\{|i-j| \leq h\}})$  as the banded counterpart of  $\mathbf{A}$  with bandwidth  $h$ .

**2. The spatio-temporal vector autoregression**

As in the recent paper by Gao et al. (2019), we consider the spatio-temporal vector autoregression

$$\mathbf{y}_t = \mathbf{A}\mathbf{y}_t + \mathbf{B}\mathbf{y}_{t-1} + \boldsymbol{\epsilon}_t, \quad t = 1, \dots, T, \tag{1}$$

where  $\mathbf{y}_t = (y_{1t}, \dots, y_{Nt})'$  stacks the observations at time  $t$  over a collection of  $N$  spatial units. The contemporaneous spatial dependence between these spatial units is governed by the matrix  $\mathbf{A} = (a_{ij})_{i,j=1}^N$  with  $a_{ii} = 0$  for  $i = 1, \dots, N$ . The matrix  $\mathbf{B} = (b_{ij})_{i,j=1}^N$  incorporates dependence on past realizations. Finally,  $\boldsymbol{\epsilon}_t$  is the innovation vector driving the variation in  $\mathbf{y}_t$ .

A possible application of model (1) is the use of satellite data to study the concentration of an air pollutant over time (see the empirical application, Section 5, for more details). Overlaying each satellite image with a spatial grid of  $N$  grid cells, we record the pixel value of each grid cell at time  $t$  in the vector  $\mathbf{y}_t = (y_{1t}, \dots, y_{Nt})'$ . From a physics viewpoint, the spatial lag  $\mathbf{A}\mathbf{y}_t$  and the temporal lag  $\mathbf{B}\mathbf{y}_{t-1}$  can capture diffusion processes in the atmosphere. That is, the concentration of the pollutant in a specific grid cell is determined by its past concentration and the inflow(outflow) of the pollutant from(to) neighbouring cells. These interactions may reasonably be expected to be localized in the sense that grid cells should be mainly influenced by the concentrations nearby, an observation that we return to in Section 2.2. The main notation and theoretical results are developed on the basis of model (1). In Section 3.3, we extend the model by including exogenous variables. In the context of air pollutants, this extension enables the explicit modelling of mechanisms that influence the dynamic evolution of pollutant concentrations beyond dispersion, such as atmospheric transport via air currents or new emissions.

**2.1. Stability and dependence**

To ensure that the spatio-temporal VAR defines a stable dynamic system, we restrict the parameter space in the following way.

**Assumption 1 (Stability).** We require: (a)  $\|\mathbf{A}\|_{\cdot} = \max\{\|\mathbf{A}\|_1, \|\mathbf{A}\|_\infty\} \leq \delta_A < 1$ , and (b)  $\|\mathbf{B}\|_{\cdot} \leq \delta_B$  and  $\delta_C := \frac{\delta_B}{1-\delta_A} < 1$ .

Assumption 1 ensures that (1) has a stable reduced form VAR(1) specification. This follows from the following two observations. First, Assumption 1(a) bounds the maximum row and column sums of  $\mathbf{A}$  and thereby constraints the contemporaneous dependence between the time series. Invertibility of  $\mathbf{I}_N - \mathbf{A}$  is guaranteed since  $\|\mathbf{A}\|_2 \leq \sqrt{\|\mathbf{A}\|_1 \|\mathbf{A}\|_\infty} \leq \delta_A < 1$ . Hence, we can define the reduced-form representation as

$$\mathbf{y}_t = \mathbf{C}\mathbf{y}_{t-1} + \mathbf{D}\boldsymbol{\epsilon}_t, \tag{2}$$

with  $\mathbf{C} = (\mathbf{I}_N - \mathbf{A})^{-1}\mathbf{B}$  and  $\mathbf{D} = (\mathbf{I}_N - \mathbf{A})^{-1}$ . It follows from  $\|\mathbf{D}\|_{\cdot} \leq \sum_{j=0}^{\infty} \|\mathbf{A}\|_{\cdot}^j = \frac{1}{1-\delta_A}$  that the absolute row and column sums of  $(\mathbf{I}_N - \mathbf{A})^{-1}$  are bounded, which is the counterpart of assumption B2 in Dou et al. (2016). Second, Assumption 1(b) controls the serial dependence of  $\mathbf{y}_t$ . Indeed, we conclude from  $\|\mathbf{C}\|_2 \leq \|\mathbf{C}\|_{\cdot} \leq \delta_C < 1$  that both unit roots and explosive behaviour of the reduced form specification are ruled out.

**Remark 1.** *Assumption 1* is defined in terms of  $\|\cdot\|_+$ . Since  $\|A\|_1 = \|A'\|_\infty \leq \|A\|_+$  for any matrix  $A$ , the norm  $\|\cdot\|_+$  is convenient when bounding products of matrices containing transposes.

**Remark 2.** *Assumption 1* echoes the spatial econometrics literature in which the spatial parameter  $\lambda$  is bounded from above and the prespecified spatial weight matrix, say  $W_N$ , is standardized (see, e.g. Lee, 2004 and Lee and Yu, 2010). Typically, the product  $\lambda W_N$  – the natural counterpart of the matrix  $A$  – is required to fulfil conditions similar to  $\|A\|_+ \leq \delta_A < 1$ . For instance, it is not uncommon to row-normalize  $W_N$  (each absolute row sum equal to 1) and restrict  $\lambda < 1$ , see pages 1903–1904 of Lee (2004). If  $W_N$  is symmetric, then  $\|\lambda W_N\|_+ < 1$  holds.

Naturally, the dependence structure of  $\{y_t\}_{t \in \mathbb{Z}}$  is also determined by the innovation process, on which we impose the following assumption.

**Assumption 2 (Innovations).**

- (a) The sequence  $\{\epsilon_t\}_{t \in \mathbb{Z}}$  is a covariance stationary, martingale difference (m.d.) process with respect to the filtration  $\mathcal{F}_{t-1} = \sigma(\epsilon_{t-1}, \epsilon_{t-2}, \dots)$ , and geometrically strong mixing ( $\alpha$ -mixing). That is, the mixing coefficients  $\{\alpha_m\}$  satisfy  $\alpha_m \leq c_2 e^{-\gamma_\alpha m}$  for all  $m$  and some constants  $c_2, \gamma_\alpha > 0$ .
- (b) The smallest eigenvalue of  $\Sigma_\epsilon = \mathbb{E}(\epsilon_t \epsilon_t') = (\sigma_{ij})_{i,j=1}^N$  has a positive lower bound, and  $\|\Sigma_\epsilon\|_+ \leq C_\epsilon < \infty$ .
- (c) Either one of the following assumptions holds:
  - (c1) For  $\psi(x) = x^d$ , we have  $\sup_{i,t} \|\epsilon_{it}\|_\psi = (\mathbb{E}|\epsilon_{it}|^d)^{1/d} \leq \mu_d < \infty$  for  $d \geq 4$ .
  - (c2) For  $\psi(x) = \exp(x) - 1$ , we have  $\sup_{i,t} \|\epsilon_{it}\|_\psi \leq \mu_\infty < \infty$ .

Similarly to Masini et al. (2022), *Assumption 2* restricts the stochastic properties of the innovation process  $\{\epsilon_t\}$ . First, *Assumption 2(a)* ensures that the memory of the innovation process fades sufficiently fast and limits the cross-sectional dependence between the elements in  $\epsilon_t$ . The m.d. assumption implies that  $\mathbb{E}(\epsilon_t y'_{t-j}) = \mathbf{0}$  while the mixing assumption controls the serial correlation in the data. *Assumption 2(b)* restricts the contemporaneous dependence in  $\epsilon_t$ . The lower bound on the minimum eigenvalue of  $\Sigma_\epsilon$  implies that each innovation contains unique information, while the upper bound on the maximum eigenvalue of  $\Sigma_\epsilon$  bounds the degree of cross-sectional dependence, for example by excluding common factor structures. Polynomial or exponential tail decay of the distribution of the innovations is imposed through either *Assumption 2(c1)* or *Assumption 2(c2)*, respectively. The type of tail decay will directly influence the growth rates we can allow for  $N$  and  $T$ . The discussions in Masini et al. (2022) demonstrate that *Assumption 2* allows for a wide range of innovation models.

While these assumptions are sufficient to estimate the reduced form VAR in (2), there is a potential issue with endogeneity. In particular, since  $y_t$  appears on both sides of (1), equation-by-equation OLS estimation is inconsistent for general  $A$  and  $B$ .<sup>1</sup> This issue can be addressed through the use of instrumental variables (e.g. Ahrens and Bhattacharjee, 2015; Lam and Souza, 2019) or via the generalized Yule–Walker equations, as considered in Gao et al. (2019) and this paper. Under *Assumptions 1* and *2*, we can post-multiply (1) by  $y'_{t-1}$  and take expectations to obtain

$$\Sigma_1 = A \Sigma_1 + B \Sigma_0. \tag{3}$$

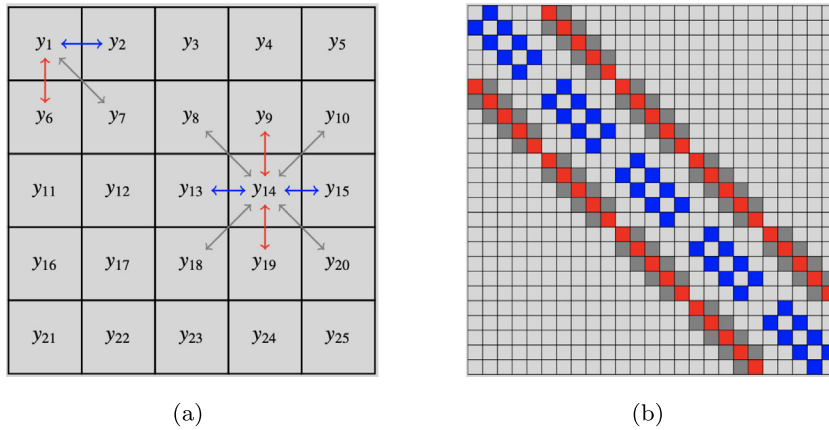
By plugging in estimates of the population covariance matrices  $\Sigma_1$  and  $\Sigma_0$  in (3), we obtain a system of equations from which the parameters in  $A$  and  $B$  can be estimated. However, with each row in (3) defining  $N$  equations with  $2N - 1$  unknowns, additional structure is required to identify and estimate the parameters.

**2.2. Structured sparsity**

There are several possibilities to introduce structure into  $A$  and  $B$ . Early spatial econometrics models, e.g. the spatial autoregressive (SAR) model or spatial Durbin model (SDM), incorporate spatial effects through the product  $\lambda W_N$  (with  $W_N$  pre-specified). The specification  $A = \lambda W_N$  imposes substantial structure on  $A$  and leaves only the single parameter  $\lambda$  to estimate. Dou et al. (2016) consider a more general setting in which each row of  $W_N$  receives its own spatial autoregressive parameter. Specifically, they set  $A = \text{diag}(\lambda_0) W_N$  and  $B = \text{diag}(\lambda_1) + \text{diag}(\lambda_2) W_N$ , and estimate the  $3N$  coefficients in  $(\lambda'_0, \lambda'_1, \lambda'_2)'$ . Gao et al. (2019) require  $A$  and  $B$  to be banded matrices. We also employ a bandedness assumption, although our bandwidth may be substantially wider than in Gao et al. (2019).

**Assumption 3 (Banded Matrices).** Recall  $A = (a_{ij})_{i,j=1}^N$ ,  $B = (b_{ij})_{i,j=1}^N$ , and  $\Sigma_\epsilon = (\sigma_{ij})_{i,j=1}^N$ . We have: (a)  $a_{ij} = b_{ij} = 0$  for all  $|i - j| > k_0$  with  $k_0 \leq \lfloor (N - 1)/4 \rfloor$ , and (b)  $\sigma_{ij} = 0$  for all  $|i - j| > l_0$ .

<sup>1</sup> As pointed out by an anonymous referee, notable exceptions occur when  $A$  is diagonal, which is ruled out in the current setting, or when  $A$  is triangular. The latter setting may be justifiable for certain structural VARs, but seems unrealistic in most spatial settings.



**Fig. 1.** Sparsity patterns in the matrices  $A$  and  $B$  due to localized interactions. (a) A  $(5 \times 5)$  grid of spatial units with arrows depicting the nearest horizontal (blue), vertical (red), and diagonal (grey) interactions for  $y_1$  and  $y_{14}$ . (b) The sparsity pattern resulting from nearest horizontal, nearest vertical and nearest diagonal interactions. The first row in the  $(25 \times 25)$  matrix has three non-zero elements corresponding to the horizontal interaction  $y_1 \rightarrow y_2$  causing a non-zero  $(1,2)$  element (blue), the vertical interaction  $y_1 \rightarrow y_6$  causing a non-zero  $(1,6)$  element (red), and  $y_1 \rightarrow y_7$  causing a non-zero  $(1,7)$  element (grey). The full colour pattern is recovered by applying this reasoning to all remaining cells in the grid. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

**Assumption 3** serves two purposes. First, for each spatial unit  $i = 1, \dots, N$ , the matrices  $A$  and  $B$  are banded to have no more than  $N$  unknown parameters per equation. With  $N$  moment conditions for each  $i$ , **Assumption 3(a)** is the widest bandwidth that enables identification of the parameters. In many applications, including the one considered in this paper, the actual bandwidth may be smaller than the upper bound imposed in **Assumption 3**. While a smaller bandwidth is neither required nor imposed, it will result in additional sparsity that our estimator is designed to exploit. Second, we demonstrate in **Theorem 1** that the combination of **Assumptions 3(a)–(b)** implies that the  $(N \times N)$  autocovariance matrices  $\Sigma_j = \mathbb{E}(y_t y'_{t-j})$  are *approximately* banded. This in turn is exploited in the Yule–Walker estimation approach, in which we incorporate banded estimation of the autocovariance matrices to improve the convergence rates of our estimator.

Even under **Assumption 3**, the number of unknown parameters in  $A$  and  $B$  may grow quadratically in  $N$ . Accurate estimation of all these parameters becomes infeasible even for moderately sized systems. To alleviate this curse of dimensionality, we rely on sparsity, which in the present context occurs when two spatial units do not interact with each other (given their interactions with other units within the system). We demonstrate how a particular, and exploitable, sparsity pattern arises when the spatial units are ordered in a structured way.

As an illustrative example, let us consider repeated measurements on the  $(5 \times 5)$  spatial grid shown in **Fig. 1(a)**. The  $N = 25$  spatial units are labelled  $y_1$  up to  $y_{25}$  and enumerated row-wise. This ordering of the spatial entities creates an implicit notion of proximity and we intuitively expect economic/physical interactions to be most pronounced at short length scales. For instance, returning to our running example of air pollution, we expect elevated concentrations of pollutants (e.g. due to high traffic or polluting industries) to diffuse through space and mostly affect neighbouring locations. **Fig. 1(a)** shows the short-range interactions of  $y_1$  and  $y_{14}$  towards their horizontal (blue), vertical (red) and diagonal (grey) nearest neighbour. If these are the only possible interactions, then the resulting sparsity pattern in the  $(25 \times 25)$  matrices  $A$  and  $B$  would correspond to the coloured elements in **Fig. 1(b)** (see its caption for further explanations). The nonzero elements in  $A$  and  $B$  are thus seen to cluster in specific, dense diagonals with an occasional zero whenever neighbours are absent at the grid boundary. It is easy to verify that more distant neighbours (in either direction) are still found on diagonals but located further away from the main diagonal. Not knowing a priori in which direction and at which distance these interactions become unimportant, we develop a penalized estimator that is able to detect and efficiently estimate such diagonal sparsity patterns *without imposing* this sparsity structure from the outset.

### 3. Sparse estimation

In this section, we develop the estimation and statistical theory for a sparse estimator of the spatio-temporal VAR in (1). As mentioned in Section 2.2, we exploit the banded structure in  $A$  and  $B$  via banded estimation of the sample autocovariance matrices. Accordingly, we start off by deriving finite-sample error bounds attained by banded sample auto-covariance matrices as estimators of their population counterparts. Afterwards, we fully develop our estimator and investigate its theoretical properties.

#### 3.1. Banded autocovariance matrix estimation

**Assumption 3** plays a crucial role in identifying the parameters in  $A$  and  $B$ , while simultaneously opening up the door for efficiency gains in the estimation of the autocovariance matrices in (3). Since the inverse of a banded matrix is in general not



banded itself, the matrix  $C = (I_N - A)^{-1}B$  containing the autoregressive coefficients of the reduced-form representation of the spatio-temporal VAR in (2) cannot be expected to be banded. However, given its structure, it may be expected to be well-approximated by a banded matrix. Consequently, combining the *approximate* bandedness of  $C$  with the bandedness of  $\Sigma_\epsilon$ , a reasonable conjecture would be that the autocovariance matrices of  $y_t$  are approximately banded as well. Indeed, we formalize this conjecture in the following theorem, along with the expected error bounds one may expect from utilizing a banded estimation approach.

**Theorem 1 (Convergence Rates for Banded Sample Autocovariance Matrices).** For any matrix  $M = (m_{ij})$ , its  $h$ -banded counterpart is defined as  $\mathcal{B}_h(M) = (m_{ij} \mathbb{1}_{\{|i-j| \leq h\}})$ . Define the  $(N \times 2N)$  matrix  $\hat{V}_h = [\mathcal{B}_h(\hat{\Sigma}_1)' \ \mathcal{B}_h(\hat{\Sigma}_0)]$  with  $\hat{\Sigma}_1 = \frac{1}{T} \sum_{t=2}^T y_t y_{t-1}'$  and  $\hat{\Sigma}_0 = \frac{1}{T} \sum_{t=2}^T y_t y_t'$ , and choose

$$h := h(\epsilon) = \left( \frac{\max \left\{ \log \left( \frac{K - \delta_\epsilon}{\delta_\epsilon} \right), \log \left( \frac{C_1}{\epsilon} \right) \right\}}{|\log(\delta_A)|} + 1 \right) \left( \frac{\log \left( \frac{C_2}{\epsilon} \right)}{|\log(\delta_C)|} + 2 \right) (k_0 - 1) + 2l_0 + 1, \tag{4}$$

for some  $K \in (\delta_c, 1)$ , with  $C_1 = \frac{34C_\epsilon}{(1-\delta_A)^2(1-K^2)^2}$  and  $C_2 = \frac{4C_\epsilon}{(1-\delta_A)^2(1-\delta_c^2)}$ . Then,  $\|\hat{V}_h - V\|_F \leq 4\epsilon$  with a probability of at least

- (a)  $1 - 2\mathcal{P}_1(\epsilon, N, T)$  under Assumptions 1–3 using Assumption 2(b1) (polynomial tails),
- (b)  $1 - 2\mathcal{P}_2(\epsilon, N, T)$  under Assumptions 1–3 using Assumption 2(b2) (exponential tails),

where

$$\mathcal{P}_1(\epsilon, N, T) = (2h + 1)N \left[ \left( b_1 T^{(1-\delta)/3} + \frac{[2h + 1]b_3}{\epsilon} \right) \exp \left( -\frac{T^{(1-\delta)/3}}{2b_1^2} \right) + \frac{b_2[2h + 1]^d}{\epsilon^d T^{\frac{\delta}{2}(d-1)}} \right],$$

for some  $0 < \delta < 1$ , and

$$\mathcal{P}_2(\epsilon, N, T) = (2h + 1)N \left[ \frac{\kappa_1[2h + 1]}{\epsilon} + \frac{2}{\kappa_2} \left( \frac{T\epsilon^2}{[2h + 1]^2} \right)^{\frac{1}{7}} \right] \exp \left( -\frac{1}{\kappa_3} \left( \frac{T\epsilon^2}{[2h + 1]^2} \right)^{\frac{1}{7}} \right).$$

The constants  $b_i, \kappa_i$ , for  $i = 1, 2, 3$ , are positive and independent of  $N$  and  $T$ .

Theorem 1 shows that banded estimators for  $\Sigma_0$  and  $\Sigma_1$  provide an accurate approximation to  $V = [\Sigma_1' \ \Sigma_0]$ . Each of these banded matrices has at most  $2h(\epsilon) + 1$  nonzero elements in their columns/rows. In other words, given  $\epsilon, l_0$  and  $k_0$ , Assumptions 1–3 guarantee that  $\Sigma_0$  and  $\Sigma_1$  can be well-approximated by matrices with bandwidths smaller than  $N$ . In the following section, we develop our sparse estimator for  $A$  and  $B$  and show that the established approximability of the population autocovariance matrices improves the convergence rate of the estimator.

### 3.2. The SPLASH $(\alpha, \lambda)$ estimator

In this section, we introduce a sparse estimator for large spatio-temporal models generated by (1). Our estimator essentially combines generalized Yule–Walker estimation with a sparse group penalty (e.g. Simon et al., 2013). Estimation via the Yule–Walker equations is adopted to control for endogeneity, while the added penalization shrinks the estimates towards the diagonally structured sparsity displayed in Fig. 1. Compared to Gao et al. (2019), we hereby gain the ability exploit sparsity within the imposed bandwidths of  $A$  and  $B$ . In practice, this means that we can assume a much larger bandwidth without substantially sacrificing estimation accuracy.

#### 3.2.1. Definition

The formal definition of our estimator requires further notation. Part of this notation comes naturally if we briefly review the generalized Yule–Walker estimator. First, we rewrite the generalized Yule–Walker conditions in (3) more compactly as

$$\Sigma_1' = [\Sigma_1' \ \Sigma_0] [A \ B]' =: VC'. \tag{5}$$

The  $i$ th column of  $C'$  contains all coefficients that belong to the  $i$ th equation in (1). Assumption 3 requires several of these coefficients to be zero, so we exclude these from the outset. We collect all remaining (possibly) nonzero coefficients of the  $i$ th equation in the vector  $c_i$ , and define  $V_i$  as the matrix containing the corresponding columns from  $V$ . In the population, we have  $V_i c_i = \Sigma_1' e_i =: \sigma_i$  for  $i = 1, \dots, N$ , with  $e_i$  being the  $i$ th column of the  $(N \times N)$  identity matrix. Sample counterparts of  $V_i$  and  $\sigma_i$  are readily available from the sample autocovariance matrices. More explicitly, Gao et al. (2019) set  $\hat{\sigma}_i = \frac{1}{T} \sum_{t=2}^T y_{t-1} y_{it}'$  and construct  $\hat{V}_i$  from the appropriate columns of  $\hat{V} = [\hat{\Sigma}_1' \ \hat{\Sigma}_0]$ . Motivated by  $\sigma_i - V_i c_i = \mathbf{0}$ , they define their estimator  $\hat{c}_i^{GMWY}$  as the following minimizer:

$$\hat{c}_i^{GMWY} = \arg \min_c \|\hat{\sigma}_i - \hat{V}_i c\|_2^2. \tag{6}$$

We will adjust this objective function in three ways. First, we define our estimator in terms of banded estimated covariance matrices. Second, we add a group penalty to sparsely estimate the parameters in  $A$  and  $B$ . To exploit the diagonally structured sparsity in Fig. 1, if it is present, we allow for a penalization of complete diagonals. This form of penalization, however, renders

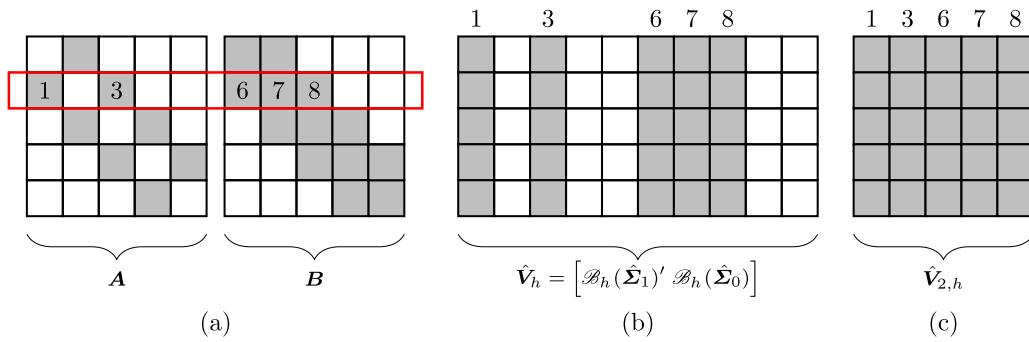


Fig. 2. A visualization on the construction of  $\hat{V}_{2,h}$  for  $N = 5$ . (a) If  $h = 1$ , then grey elements in  $A$  and  $B$  are (potentially) nonzero whereas white elements are zero by construction. Enumerating along the second row, the active elements are in the set  $\{1, 3, 6, 7, 8\}$ . (b) We select the columns from  $\hat{V}_h$  corresponding to the active set. (c) The matrix  $\hat{V}_{2,h}$  is the submatrix of  $\hat{V}_h$  with only active columns.

equation-by-equation estimation of the parameters infeasible. Therefore, recalling the definitions of  $\mathcal{B}_h(\hat{\Sigma}_1)$  and  $\hat{V}_h$  in Theorem 1, we define  $\hat{\sigma}_h = \text{vec}(\mathcal{B}_h(\hat{\Sigma}_1)')$  and  $\hat{V}_h^{(d)} = \text{diag}(\hat{V}_{1,h}, \dots, \hat{V}_{N,h})$ , where  $\hat{V}_{i,h}$  contains the columns of  $\hat{V}_h$  that multiply the elements in  $c_i$  (see Fig. 2 for an illustration). Finally, we construct the penalty function. We define an index set that partitions the vector  $c$  into sub-vectors, denoted  $\{c_g\}$ , that contain the non-zero diagonals of  $A$  and  $B$  that are admissible under Assumption 3 as

$$\begin{aligned} \mathcal{G}_A &:= \{g \subset \mathbb{N} : c_g = \mathcal{D}_A(k), k \in \{1, \dots, \lfloor (N-1)/4 \rfloor\}\}, \\ \mathcal{G}_B &:= \{g \subset \mathbb{N} : c_g = \mathcal{D}_B(k), k \in \{0, \dots, \lfloor (N-1)/4 \rfloor\}\}, \end{aligned} \tag{7}$$

respectively, and let  $\mathcal{G} = \mathcal{G}_A \cup \mathcal{G}_B$ . Based on this notation, we define our objective function as

$$\mathcal{L}_\alpha(c; \lambda) = \|\hat{\sigma}_h - \hat{V}_h^{(d)} c\|_2^2 + \lambda \underbrace{\left( (1 - \alpha) \sum_{g \in \mathcal{G}} \sqrt{|g|} \|c_g\|_2 + \alpha \|c\|_1 \right)}_{=: P_\alpha(c)} = \|\hat{\sigma}_h - \hat{V}_h^{(d)} c\|_2^2 + \lambda P_\alpha(c). \tag{8}$$

The spatial lasso-type shrinkage estimator, abbreviated SPLASH( $\alpha, \lambda$ ) or SPLASH for short, is defined as the minimizer of (8), i.e.  $\hat{c} = \arg \min_c \mathcal{L}_\alpha(c; \lambda)$ . The influence of the penalty function  $P_\alpha(c)$  is governed by the penalty parameter  $\lambda$  and the second hyperparameter  $\alpha$  balances group-structured sparsity versus individual sparsity. At the extremities of  $\alpha \in [0, 1]$  we find the group lasso ( $\alpha = 0$ ) and the lasso ( $\alpha = 1$ ). Intermediate values of  $\alpha$  will shrink groups of diagonal coefficients in  $A$  and  $B$ , and individual parameters. The SPLASH solution promotes completely sparse diagonals and sparse elements within nonzero diagonals, and thus shrinks towards the sparsity patterns displayed in Fig. 1(b).

### 3.2.2. Implementation

Algorithm 1 describes how to apply the SPLASH estimator. There are three key parts: (1) the selection of the bandwidth parameter  $h$ , (2) the determination of the hyperparameters  $\alpha$  and  $\lambda$ , and (3) the minimization of the objective function  $\mathcal{L}_\alpha(c; \lambda)$ . Further details on parts (1) and (3) can be found in the articles that are referenced in Algorithm 1. The hyperparameters are determined using time series cross-validation. The grid values being outlined in the algorithm are mere recommendations. Finer grids can yield better solutions at the cost of higher computational demands. If computational resources are scarce, then hyperparameters might also be pre-specified (e.g. fixing  $\alpha$  to say 0.5 in order to obtain a one-dimensional grid search). Finally, we mention that there is no need to manually implement the steps in Algorithm 1. An efficient implementation of this procedure is readily available as an R package on the website of one of the authors.<sup>2</sup>

### 3.2.3. Theoretical properties

We require an additional assumption on the DGP to ensure that  $A$  and  $B$  in (1) are uniquely identified. To this end, we leverage on Assumption 3, which enables unique identification of  $A$  and  $B$  via a straightforward full-rank condition on sub-matrices of the autocovariance matrices that appear in the generalized Yule–Walker equations.

**Assumption 4 (Restricted Minimum Eigenvalue).** Assume that

$$\phi_{\min}(\mathbf{x}) := \min_{\mathbf{x} \in \mathbb{R}^{2N} : \|\mathbf{x}\|_2 \leq N} \frac{\|\mathbf{V}\mathbf{x}\|_2}{\|\mathbf{x}\|_2} \geq \phi_0 > 0.$$

<sup>2</sup> <https://sites.google.com/view/etiennewijler/code>

**Algorithm 1** SPLASH( $\alpha, \lambda$ ) implementation

- 1: Compute the  $(N \times N)$  matrices  $\hat{\Sigma}_0 = \frac{1}{T} \sum_{t=2}^T y_t y_t'$  and  $\hat{\Sigma}_1 = \frac{1}{T} \sum_{t=2}^T y_t y_{t-1}'$ .
- 2: Determine the bandwidth  $h$  using the bootstrap procedure in Guo et al. (2016, p. 7).
- 3: Using the selected bandwidth of step 2, calculate  $\mathcal{B}_h(\hat{\Sigma}_0)$ ,  $\mathcal{B}_h(\hat{\Sigma}_1)$ , and construct  $\hat{\sigma}_h$ ,  $\hat{V}^{(d)}$  and  $\mathcal{G}$  as in Section 3.2.1.
- 4: Define a grid of hyperparameter values. Depending on computational power, we suggest  $\alpha_{\text{grid}} = \{0, 0.25, 0.5, 0.75, 1\}$  and a vector of  $K$  logarithmically spaced points for  $\lambda$ . That is, with  $\lambda_{\text{max}} = \max\left(\max_{g \in \mathcal{G}} \frac{1}{T} \left\| \hat{V}_{h,g}^{(d)'} \hat{\sigma}_h \right\|_2 / \sqrt{|g|}, \max_{1 \leq i \leq N_c} \left| \hat{V}_{h,i}^{(d)'} \hat{\sigma}_h \right| \right)$  and  $\lambda_{\text{min}} = \eta \lambda_{\text{max}}$  for some small  $\eta > 0$ , take  $\lambda_{\text{grid}} = \{\lambda_1, \dots, \lambda_K\}$  where

$$\lambda_i = \exp\left(\ln(\lambda_{\text{max}}) - \frac{i-1}{K-1} \left[\ln(\lambda_{\text{max}}) - \ln(\lambda_{\text{min}})\right]\right) \quad i = 1, \dots, K.$$

- 5: For each combination  $(\alpha, \lambda)$ , use the SGL algorithm as in Simon et al. (2013) to compute the SPLASH solution  $\hat{c}(\alpha, \lambda)$ . Select the solution with the hyperparameters that perform best via time series cross-validation (see, e.g. Hyndman and Athanasopoulos, 2018).

Assumption 4 states that every sub-matrix containing at most  $N$  columns from  $V$  has full column-rank and a minimum singular value bounded away from zero. Related assumptions appear in Bickel et al. (2009, Section 4), who refer to  $\phi_{\min}(x)$  as a *restricted eigenvalue* and use this quantity to construct sufficient conditions for their restricted eigenvalue condition. Assumption 4 fits our framework particularly well, as the assumed maximum bandwidth of the matrices  $A$  and  $B$  in Assumption 3 imply that the diagonal blocks of the matrix  $V^{(d)}$  never contain more than  $N$  unique columns of  $V$ . Using this property, we show in Lemma 1 of Appendix A that a Sparse Group Lasso compatibility condition is implied by Assumption 4.

Equipped with Assumption 4, we find the following finite-sample bounds on the prediction and estimation error of SPLASH.

**Theorem 2.** Define

$$\bar{\omega}_\alpha = \max\left\{(1-\alpha) \sum_{g \in \mathcal{G}_S} \sqrt{|g|}, \alpha \sqrt{|S|}\right\},$$

where  $\mathcal{G}_s = \{g \in \mathcal{G} : c_g \neq 0\}$  and  $S = \{j : c_j \neq 0\}$ . Under Assumptions 3–4 and  $\|V\|_- \leq C_V$ , it holds that

$$\left\| \hat{V}_h^{(d)}(\hat{c} - c) \right\|_2^2 + \lambda \left( (1-\alpha) \sum_{g \in \mathcal{G}} \sqrt{|g|} \|\hat{c}_g - c_g\|_2 + \alpha \|\hat{c} - c\|_1 \right) \leq \frac{64 \bar{\omega}_\alpha^2 \lambda^2}{\phi_0^2} \tag{9}$$

with a probability of at least

- (a)  $1 - 10 \mathcal{P}_1(f(\lambda, \phi_0), N, T)$  when Assumption 2(b1) (polynomial tail decay) is valid, or
- (b)  $1 - 10 \mathcal{P}_2(f(\lambda, \phi_0), N, T)$  when Assumption 2(b2) (exponential tail decay) is valid,

where  $\mathcal{P}_1(x, N, T)$  and  $\mathcal{P}_2(x, N, T)$  are defined in Theorem 1 and  $f(\lambda, \phi_0) = \min\left(\frac{\lambda^{1/2}}{24}, \frac{\lambda}{96C_V}, \frac{\phi_0}{12}\right)$ .

Theorem 2 contains a finite-sample bound on the prediction and estimation error for the SPLASH( $\alpha, \lambda$ ) estimator. It offers some interesting insights. First, we focus on the probability with which inequality (9) holds. For VAR estimation with a penalized least-squares objective function, such probabilities are governed by tail probabilities of the process  $\{\frac{1}{T} \sum_{t=1}^T y_{it} \epsilon_{jt}\}$  (see, e.g. Lemma 4 in Kock and Callot, 2015, or Lemmas 5–6 in Medeiros and Mendes, 2016). Because Yule–Walker estimation relies primarily on autocovariance matrix estimation, our probability depends on the tail decay of the distribution of  $\{\|\hat{V}_h - V\|_-\}$ . Overall, the probability of (9) improves through faster tail decay of the innovation distribution (compare cases (a) and (b)) and banded autocovariance matrix estimation (Theorem 1). Second, we look closer at the performance upper bound itself. The right-hand side of (9) demonstrates that the upper bound on the prediction and estimation error is increasing in  $\bar{\omega}_\alpha$ , which in turn is increasing in the bandwidths  $k_0$  and  $l_0$ , increasing in the group sizes ( $\alpha < 1$ ), and increasing in the number of relevant interactions  $|S|$  ( $\alpha > 0$ ). Furthermore, the prediction and estimation error increases in the degree of penalization. Whereas this seemingly suggests to minimize  $\lambda$  as to improve performance bounds, we emphasize that the effect of regularization in Theorem 2 is two-fold: increasing regularization deteriorates the performance bound, but increases the probability of the set on which the performance bound holds. Intuitively, shrinkage induces finite-sample bias which worsens accuracy, but simultaneously reduces sensitivity to noise, thereby enabling performance guarantees at higher degrees of certainty.

The aforementioned effects can also be demonstrated by means of an asymptotic analysis. Based on Theorem 2, we derive the conditions for convergence of the prediction and estimation errors in the following corollary. The exact convergence rates are also provided.

**Corollary 1.** Let  $\lambda \in O(T^{-q_\lambda})$ ,  $N \in O(T^{q_N})$ ,  $|\mathcal{G}_S| \in O(T^{q_g})$ ,  $|S| \in O(T^{q_s})$ ,  $k_0, l_0 \in O(T^{q_k})$ , where  $q_\lambda, q_N, q_s$ , and  $q_k$  are fixed and positive constants. Maintain Assumptions 1–4 and assume that either (i)  $q_\lambda < \frac{\delta(d-1)}{2d} - \frac{(d+1)q_k}{d} - \frac{q_N}{d}$  for some  $0 < \delta < 1$  and Assumption 2(b1) holds, or (ii)  $q_\lambda < \frac{1}{2} - q_k$  and Assumption 2(b2) holds. Then,



$$(a) \left\| \hat{V}_h^{(d)} (\hat{c} - c) \right\|_2^2 = O_p \left( (1 - \alpha) T^{2q_g + q_N - 2q_\lambda} + \alpha T^{q_s - 2q_\lambda} \right),$$

$$(b) P_\alpha (\hat{c} - c) = (1 - \alpha) \sum_{g \in \mathcal{G}} \sqrt{|g|} \left\| \hat{c}_g - c_g \right\|_2 + \alpha \|\hat{c} - c\|_1 = O_p \left( (1 - \alpha) T^{2q_g + q_N - q_\lambda} + \alpha T^{q_s - q_\lambda} \right).$$

**Corollary 1** provides insights into the determinants of the convergence rate. In particular, the result confirms that the convergence rate decreases in the bandwidths  $k_0$  and  $l_0$ , the number of spatial units  $N$ , the number of interactions  $|S|$  and the degree of penalization  $\lambda$ . To ensure that the set on which the performance bound in **Theorem 2** holds occurs with probability converging to one, conditions (i) and (ii) impose that the degree of penalization does not decay too fast in  $T$ . The optimal convergence rate is obtained by choosing  $q_\lambda$  as large as possible without violating these conditions. Some concrete examples are provided in **Remark 3**.

**Remark 3.** Insightful special cases can be examined based on **Corollary 1**. For the sake of brevity, we consider two cases while focusing on the estimation error  $P_\alpha (\hat{c} - c)$  and assuming errors with at least  $d$  finite moments (**Assumption 2(b1)**). In the absence of within-group shrinkage ( $\alpha = 0$ ), **Corollary 1** demonstrates that  $P_0 (\hat{c} - c) = O_p \left( T^{2q_g + q_N - q_\lambda} \right)$ , with  $q_\lambda < \frac{1}{2} - \frac{1}{2d} - \frac{(d+1)q_k}{d} - \frac{q_N}{d}$ . The estimator now converges almost at rate  $\frac{T^{1/2-1/2d}}{N^{1+1/d} k_0^{(d+1)/d} |\mathcal{G}_S|}$ . For fixed  $N$  and large  $d$ , this is close to the common  $\sqrt{T}$ -rate of fixed-dimensional settings without regularization. If shrinkage is imposed at the individual interaction level only ( $\alpha = 1$ ), then  $P_1 (\hat{c} - c) = O_p (T^{q_s - q_\lambda})$  and the estimation error converges almost at the rate  $\frac{T^{1/2-1/2d}}{|S|N^{1/d} k_0^{(d+1)/d}}$ . Noting that  $N^{1+1/d} |\mathcal{G}_S| > |S| N^{1/d}$ , we see that **SPLASH**(1,  $\lambda$ ) attains a convergence rate at least as fast **SPLASH**(0,  $\lambda$ ), and possibly faster when the sparsity is unstructured or the diagonals are highly sparse.

### 3.2.4. Future work: Selection and inference

In this section, we propose two prospective avenues of future research opened up by the development of **SPLASH**. In particular, we discuss how the estimator may be extended to provide consistent variable selection or to enable uniformly valid statistical inference.

We conjecture that a weighted penalty scheme, similar in spirit to the adaptive lasso in **Zou (2006)**, will be able to estimate the zero coefficients as exactly zero with probability converging to one, while maintaining the ability to consistently estimate the non-zero parameters in (1). Such an estimator could be defined as the minimizer of

$$\mathcal{L}_\alpha(c; \lambda, \omega) = \left\| \hat{\sigma}_h - \hat{V}_h^{(d)} c \right\|_2^2 + \lambda \left( (1 - \alpha) \sum_{g \in \mathcal{G}} \sqrt{|g|} \left\| c_g \right\|_2 + \alpha \sum_{j=1}^M \frac{|c_j|}{\omega_j} \right),$$

where  $M$  denotes number of coefficients in the model. Compared to the objective function in (8), the vector  $\omega = (\omega_1, \dots, \omega_M)'$  is a vector of positive weights that should be large for parameters that require little shrinkage and small for parameters that are small or zero. A natural choice that satisfies this requirement is  $\omega_i = |\hat{c}_i|^\gamma$  for some  $\gamma \geq 1$ , where  $\hat{c}_i$  is the  $i$ th element of the **SPLASH** estimator  $\hat{c}$  defined in (8). We leave the investigation of the theoretical properties and empirical performance of this “*adaptive SPLASH estimator*” for future research.

Another prospective extension of our method concerns the ability to conduct honest statistical inference. The main challenges for inference based on sparse estimators are: (1) the intractability of the limiting distribution and/or its dependence on unknown parameters (e.g. **Knight and Fu, 2000**), and (2) convergence results not holding uniformly over the parameter space (e.g. **Leeb and Pötscher, 2005, 2008**). To circumvent these challenges, a “*debiased SPLASH estimator*” for uniformly valid inference might be derived by inverting the Karush–Kuhn–Tucker (KKT) conditions, analogous to the desparsified lasso in **Van de Geer et al. (2014)**. Compared to **Van de Geer et al. (2014)**, however, additional care should be taken to account for the banded autocovariance matrix and the group penalty. A full derivation is outside the scope of this paper and subject to further investigation.

### 3.3. Exogenous variables

We generalize model specification (1) by accommodating  $K$  exogenous variables, i.e.

$$y_t = A y_t + B y_{t-1} + \sum_{k=1}^K \text{diag}(\beta_k) x_{t,k} + e_t, \quad t = 1, \dots, T. \tag{10}$$

Each vector  $x_{t,k} = (x_{1t,k}, \dots, x_{Nt,k})'$  augments the spatio-temporal vector autoregression with an extra regressor. This regressor may vary over time and is assumed exogenous, i.e. we have  $\mathbb{E}(x_{t,k} e_t') = \mathbf{0}$  for  $k = 1, \dots, K$ . For notational brevity, we consider the situation in which the exogenous regressors  $x_{it,1}, \dots, x_{it,K}$  can only directly influence spatial unit  $i$ . This explains the diagonal structure in  $\text{diag}(\beta_k)$ . In **Remark 5** we argue that this simplification does not greatly hinder generality. In contrast to **Ma et al. (2023)**, we allow  $\beta_k = (\beta_{1k}, \dots, \beta_{Nk})'$  to vary with location. We keep  $K$  fixed.

To account for the exogenous variables, we modify the generalized Yule–Walker estimator of Section 3.2. We recall  $\Sigma_j = \mathbb{E}(y_t y_{t-j}')$ , and define the matrices  $\Sigma_j^{x_k y} = \mathbb{E}(x_{t,k} y_{t-j}')$  and  $\Sigma_j^{x_k x_\ell} = \mathbb{E}(x_{t,k} x_{t-j,\ell}')$ . Two sets of Yule–Walker equations, namely

$$\Sigma_1 = A \Sigma_1 + B \Sigma_0 + \sum_{k=1}^K \text{diag}(\beta_k) \Sigma_1^{x_k y} \tag{11a}$$

and

$$(\Sigma_0^{x_j y})' = \mathbf{A}(\Sigma_0^{x_j y})' + \mathbf{B}(\Sigma_1^{x_j y})' + \sum_{k=1}^K \text{diag}(\beta_k) \Sigma_0^{x_k x_j}, \quad \text{for } j = 1, \dots, K, \tag{11b}$$

are derived by post-multiplying the model by respectively  $\mathbf{y}'_{t-1}$  and  $\mathbf{x}'_{t,k}$ , and taking expectations. Compared to (5), the Yule–Walker equations in (11a) contain the additional term  $\sum_{k=1}^K \text{diag}(\beta_k) \Sigma_1^{x_k y}$  to provide information on  $\beta_1, \dots, \beta_K$ . However, if  $\Sigma_1^{x_k y} = \mathbf{O}$  (e.g. when  $\{\mathbf{x}_{t,k}\}$  and  $\{\mathbf{y}_t\}$  are independent and  $\beta_k = \mathbf{0}$ ), then (11a) alone will not identify  $\beta_k$ . We therefore add the additional Yule–Walker equations in (11b). To develop the estimator, we combine (11a) and (11b) into

$$\begin{bmatrix} \Sigma'_1 \\ \Sigma_0^{x_1 y} \\ \vdots \\ \Sigma_0^{x_K y} \end{bmatrix} = \begin{bmatrix} \Sigma'_1 & \Sigma_0 \\ \Sigma_0^{x_1 y} & \Sigma_1^{x_1 y} \\ \vdots & \vdots \\ \Sigma_0^{x_K y} & \Sigma_1^{x_K y} \end{bmatrix} [\mathbf{A} \quad \mathbf{B}]' + \sum_{k=1}^K \begin{bmatrix} (\Sigma_1^{x_k y})' \\ \Sigma_0^{x_1 x_k} \\ \vdots \\ \Sigma_0^{x_K x_k} \end{bmatrix} \text{diag}(\beta_k) := \mathbf{V}^* \mathbf{C}' + \sum_{k=1}^K \mathbf{W}_k^* \text{diag}(\beta_k). \tag{12}$$

From this point onward, the development of the SPLASHX( $\alpha, \lambda$ ) estimator closely mimics the reasoning from Section 3.2.1. First, we focus on the  $i$ th spatial unit and collect all the nonzero coefficients of  $\mathbf{A}$  and  $\mathbf{B}$  (as stipulated by Assumption 3) in  $c_i$ . Letting  $\mathbf{V}_i^*$  denote the columns in  $\mathbf{V}^*$  related to  $c_i$  and defining both  $\sigma_i^* = [\Sigma_1 \quad (\Sigma_0^{x_1 y})' \quad \dots \quad (\Sigma_0^{x_K y})']' e_i$  and  $\mathbf{w}_{ik}^* = \mathbf{W}_k^* e_i$ , result (12) implies  $\mathbf{V}_i^* c_i + \sum_{k=1}^K \mathbf{w}_{ik}^* \beta_{ik} = \sigma_i^*$ . Second, we define (a) the sample counterparts of  $\Sigma_j, \Sigma_j^{x_k y}$  and  $\Sigma_j^{x_k x_\ell}$  as respectively  $\hat{\Sigma}_j = \frac{1}{T} \sum_{t=j+1}^T \mathbf{y}_t \mathbf{y}'_{t-j}$ ,  $\hat{\Sigma}_j^{x_k y} = \frac{1}{T} \sum_{t=j+1}^T \mathbf{x}_{t,k} \mathbf{y}'_{t-j}$  and  $\hat{\Sigma}_j^{x_k x_\ell} = \frac{1}{T} \sum_{t=j+1}^T \mathbf{x}_{t,k} \mathbf{x}'_{t-j, \ell}$ , and (b) define the quantities  $\hat{\sigma}_i^*, \hat{\mathbf{w}}_{ik}^*$  and  $\hat{\mathbf{V}}_i^*$  based on their underlying sample covariance matrix estimators. Finally, set  $\hat{\sigma}^* = (\hat{\sigma}_1^{*t}, \dots, \hat{\sigma}_N^{*t})', \hat{\mathbf{V}}^{*(d)} = \text{diag}(\hat{\mathbf{V}}_1^*, \dots, \hat{\mathbf{V}}_N^*)$ , and  $\hat{\mathbf{W}}_k^{*(d)} = \text{diag}(\hat{\mathbf{w}}_{1k}^*, \dots, \hat{\mathbf{w}}_{Nk}^*)$ . The SPLASHX( $\alpha, \lambda$ ) objective function is

$$\begin{aligned} \mathcal{L}_\alpha^*(\beta_1, \dots, \beta_K, c; \lambda) &= \left\| \hat{\sigma}^* - \hat{\mathbf{V}}^{*(d)} c - \sum_{k=1}^K \hat{\mathbf{W}}_k^{*(d)} \beta_k \right\|_2^2 \\ &+ \lambda \left( P_\alpha(c) + \sum_{k=1}^K (1 - \alpha) \sqrt{N} \|\beta_k\|_2 + \alpha \|\beta_k\|_1 \right). \end{aligned} \tag{13}$$

This objective function allows for the estimation of  $\beta_1, \dots, \beta_K$ , which may contain sparse coefficients or completely sparse vectors  $\beta_k$ , as well as sparse diagonals in the coefficient matrices  $\mathbf{A}$  and  $\mathbf{B}$ . There is a clear mathematical resemblance between the SPLASH and SPLASHX estimators. Accordingly, under appropriate modifications to Assumptions 1–4, a finding similar to Theorem 2 is attainable. For the reader’s convenience, we make the correspondence between these assumptions explicit by adhering to the original numbering while adding a “\*”. A short discussion of these assumption is found at the end of this section.

**Assumption 1\*.** We require: (a)  $\|\mathbf{A}\|_- = \max\{\|\mathbf{A}\|_1, \|\mathbf{A}\|_\infty\} \leq \delta_A < 1$ , (b)  $\|\mathbf{B}\|_- \leq C_B$  and  $\frac{C_B}{1 - \delta_A} < 1$ , and (c)  $\|\beta^*\|_1 = \max\{\|\beta_1\|_1, \dots, \|\beta_K\|_1\} \leq C_\beta$ .

**Assumption 2\*.**

- (a) The regressor  $x_{i_1 t_1, k}$  and innovation  $\epsilon_{i_2, t_2}$  are independent of each other for all  $1 \leq i_1 < i_2 \leq N$ , all  $1 \leq t_1, t_2 \leq T$ , and all  $k = 1, \dots, K$ .
- (b) The sequence  $\{\epsilon_t\}$  is a covariance stationary, martingale difference process with respect to the filtration  $\mathcal{F}_{t-1} = \sigma(\epsilon_{t-1}, \epsilon_{t-2}, \dots)$ , and geometrically strong mixing ( $\alpha$ -mixing). That is, its mixing coefficients  $\{\alpha_m\}$  satisfy  $\alpha_m \leq c_2 e^{-\gamma_\alpha m}$  for all  $m$  and some constants  $c_2, \gamma_\alpha > 0$ . The largest and smallest eigenvalues of  $\Sigma_\epsilon = \mathbb{E}(\epsilon_1 \epsilon_1')$  are bounded away from 0 and  $\infty$ .
- (c) For each  $k = 1, \dots, K$ , the sequence  $\{\mathbf{x}_{t,k}\}$  is covariance stationary and geometrically strong mixing ( $\alpha$ -mixing). That is, its mixing coefficients  $\{\alpha_m^*\}$  satisfy  $\alpha_m^* \leq c_2^* e^{-\gamma_\alpha^* m}$  for all  $m$  and some constants  $c_2^*, \gamma_\alpha^* > 0$ .
- (d) Either one of the following assumptions holds:

- (d1) For  $\psi(x) = x^d$  and  $d \geq 4$ , we require  $\sup_{i,t} \|\epsilon_{it}\|_\psi = (\mathbb{E}|\epsilon_{it}|^d)^{1/d} \leq \mu_d < \infty$  and  $\sup_{i,t} \|\mathbf{x}_{it,k}\|_\psi = (\mathbb{E}|\mathbf{x}_{it,k}|^d)^{1/d} \leq \mu_d^* < \infty$  ( $k = 1, \dots, K$ ).
- (d2) For  $\psi(x) = \exp(x) - 1$ , we have  $\sup_{i,t} \|\epsilon_{it}\|_\psi \leq \mu_\infty < \infty$  and  $\sup_{i,t} \|\mathbf{x}_{it,k}\|_\psi \leq \mu_\infty < \infty$  ( $k = 1, \dots, K$ ).

**Assumption 3\*.** Recall  $\mathbf{A} = (a_{ij})_{i,j=1}^N, \mathbf{B} = (b_{ij})_{i,j=1}^N$ , and  $\Sigma_\epsilon = (\sigma_{ij})_{i,j=1}^N$ . We have: (a)  $a_{ij} = b_{ij} = 0$  for all  $|i - j| > k_0$  with  $k_0 \leq \lfloor (N - 1)/4 \rfloor$ , and (b)  $\sigma_{ij} = 0$  for all  $|i - j| > l_0$ .

**Assumption 4\* (Restricted Minimum Eigenvalue).** Define the  $((K + 1)N \times (K + 2)N)$  matrix

$$\mathbf{Q} = \begin{bmatrix} \Sigma'_1 & \Sigma_0 & (\Sigma_1^{x_1 y})' & \dots & (\Sigma_1^{x_K y})' \\ \Sigma_0^{x_1 y} & \Sigma_1^{x_1 y} & \Sigma_0^{x_1 x_1} & \dots & \Sigma_0^{x_1 x_K} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \Sigma_0^{x_K y} & \Sigma_1^{x_K y} & \Sigma_0^{x_K x_1} & \dots & \Sigma_0^{x_K x_K} \end{bmatrix}.$$

We assume that

$$\phi_{\min}^*(\mathbf{x}) := \min_{\mathbf{x} \in \mathbb{R}^{(K+2)N} : \mathcal{M}(\mathbf{x}) \leq N+K} \frac{\|\mathbf{Q}\mathbf{x}\|_2}{\|\mathbf{x}\|_2} \geq \phi_0^*.$$

Assumption 1\*(a)–(b) are as before. Assumption 1\*(c) merely assumes an upper bound on the magnitude of the exogenous regressor coefficients. Assumption 2\* controls dependencies over time, in the cross-section, and with the exogenous regressors. Similar to assumption A8(i) in Ma et al. (2023), we enforce exogeneity through Assumption 2\*(a). Assumption 2\*(b) is a simple repetition of Assumption 2(b) and its counterpart for the  $\{\mathbf{x}_{t,k}\}$ 's is encountered as Assumption 2\*(c). All original moment conditions are also transferred to the exogenous regressors (Assumption 3\*). Clearly, Assumptions 1\*–4\* allow for an easy analogy with the earlier assumptions in this paper (at the cost of possibly being more restrictive than strictly necessary). That is, define  $\epsilon_t^* = \epsilon_t + \sum_{k=1}^K \text{diag}(\beta_k)\mathbf{x}_{t,k}$  and note the linearity of  $\epsilon_t^*$  in  $\epsilon_t$  and  $\{\mathbf{x}_{t,1}, \dots, \mathbf{x}_{t,K}\}$ . With  $K$  being fixed, all mixing properties and moments conditions simply carry over to  $\epsilon_t^*$ . Having modified the assumptions to accommodate the setting with exogenous variables, we are able to extend Theorem 2 and derive an applicable finite-sample error bound.

**Theorem 3.** Define  $\mathbf{q} = (c', \beta_1', \dots, \beta_K')$ ,  $S^* = \{j : q_j \neq 0\}$  and

$$\bar{\omega}_\alpha^* = \max \left\{ (1 - \alpha) \left( \sum_{g \in \mathcal{G}_S} \sqrt{|g|} + \sqrt{N} \sum_{k=1}^K \mathbb{1}_{\{\beta_k \neq 0\}} \right), \alpha \sqrt{|S|^*} \right\}.$$

Under Assumptions 1\*–4\* and  $\|\mathbf{Q}\|_{\cdot} \leq C_Q$ , it holds that

$$\begin{aligned} & \left\| \hat{\mathbf{V}}^{*(d)}(\hat{c} - c) + \sum_{k=1}^K \hat{\mathbf{W}}_k^{*(d)}(\hat{\beta}_k - \beta_k) \right\|_2^2 + \lambda \left[ (1 - \alpha) \left( \sum_{g \in \mathcal{G}} \sqrt{|g|} \|\hat{c}_g - c_g\|_2 + \sum_{k=1}^K \sqrt{N} \|\hat{\beta}_k - \beta_k\|_2 \right) \right. \\ & \left. + \alpha \left( \|\hat{c} - c\|_1 + \sum_{k=1}^K \|\hat{\beta}_k - \beta_k\|_1 \right) \right] \leq \frac{64\bar{\omega}_\alpha^{*2} \lambda^2}{\phi_0^{*2}} \end{aligned}$$

with a probability of at least

- (a)  $1 - 7(K + 1)(K + 2)\mathcal{P}_1^*(f^*(\lambda, \phi_0^*), N, T)$  under Assumption 2(b1) (polynomial tail decay), or
- (b)  $1 - 7(K + 1)(K + 2)\mathcal{P}_2^*(f^*(\lambda, \phi_0^*), N, T)$  under Assumption 2(b2) (exponential tail decay),

where  $f^*(\lambda, \phi_0^*) = \min \left\{ \frac{\lambda^{1/2}}{12\sqrt{6}}, \frac{\lambda}{144C_Q}, \frac{\lambda^{1/2}}{12\sqrt{6}C_\beta}, \frac{\lambda}{144C_Q C_\beta}, \frac{\phi_0^*}{12} \right\}$ ,

$$\mathcal{P}_1^*(\epsilon, N, T) = N^2 \left[ \left( b_1 T^{(1-\delta)/3} + \frac{(K + 2)N b_3}{\epsilon} \right) \exp \left( -\frac{T^{(1-\delta)/3}}{2b_1^2} \right) + \frac{b_2(K + 2)^d N^d}{e^d T^{\frac{\delta}{2}(d-1)}} \right]$$

for some  $0 < \delta < 1$ , and

$$\mathcal{P}_2^*(\epsilon, N, T) = N^2 \left[ \frac{\kappa_1(K + 2)N}{\epsilon} + \frac{2}{\kappa_2} \left( \frac{T\epsilon^2}{(K + 2)N} \right)^{1/7} \right] \exp \left( -\frac{1}{\kappa_3} \left( \frac{T\epsilon^2}{(K + 2)^2 N^2} \right)^{1/7} \right)$$

All constants ( $b_1, b_2, \kappa_1$ , etc.) are positive and independent of  $N$  and  $T$ , see Theorem 1.

**Remark 4.** The inclusion of exogenous variables affects the autocovariance structure of the data. For example, if  $\mathbf{B} = \mathbf{O}$ , then  $\mathbf{y}_t = (\mathbf{I}_n - \mathbf{A})^{-1} \left[ \sum_{k=1}^K \text{diag}(\beta_k)\mathbf{x}_{t,k} + \epsilon_t \right]$  and

$$\mathbb{E}(\mathbf{y}_t \mathbf{y}_t') = (\mathbf{I}_N - \mathbf{A})^{-1} \left[ \sum_{k,\kappa=1}^K \text{diag}(\beta_k)\mathbb{E}(\mathbf{x}_{t,k}\mathbf{x}_{t,\kappa}') \text{diag}(\beta_\kappa) + \Sigma_\epsilon \right] (\mathbf{I}_N - \mathbf{A}')^{-1}.$$

Clearly,  $\mathbb{E}(\mathbf{y}_t \mathbf{y}_t')$  now also depends on the various second moments of the exogenous covariates. We do not make any a priori assumptions on  $\mathbb{E}(\mathbf{x}_{t,k}\mathbf{x}_{t,\kappa}')$  and thus define SPLASHX( $\alpha, \lambda$ ) in terms of the unband ed autocovariance matrix estimators.

**Remark 5.** Defining the coefficient matrix in front of  $\mathbf{x}_{t,k}$  as diagonal is not restrictive. That is, by letting  $\mathbf{x}_{t,k+1}$  be a reordered version of  $\mathbf{x}_{t,k}$ , the former's addition to the model can accommodate for the situation in which the dependent variable is influenced by the exogenous variable  $\mathbf{x}_{t,k}$  from multiple locations.

## 4. Simulations

### 4.1. Simulation setting

In this section, we explore the finite sample performance of our estimator by Monte Carlo simulation. The data generating process underlying the simulations is the spatio-temporal VAR in (1). We study  $T \in \{500, 1000, 2000\}$  and draw all errors  $\epsilon_{it}$  independently

and  $N(0, 1)$  distributed. The matrices  $\mathbf{A}$  and  $\mathbf{B}$  and the cross-sectional dimension  $N$  are specified in the two designs below. All simulation results are based on  $N_{sim} = 500$  Monte Carlo replications.

**Design A (Spatial grid with neighbour interactions):** As in Fig. 1, we consider an  $(m \times m)$  grid of spatial units. For  $m = 5$  ( $m = 10$ ), this results in a cross-sectional dimension of  $N = 25$  ( $N = 100$ ). The matrix  $\mathbf{A}$  contains interactions between first horizontal and first vertical neighbours while all other coefficients are zero. The magnitude of these nonzero interactions are 0.2. For  $m = 5$  ( $m = 10$ ), the temporal matrix  $\mathbf{B}$  is a diagonal matrix with elements 0.25 (0.21) on the diagonal. The reduced form VAR matrix  $\mathbf{C} = (\mathbf{I}_N - \mathbf{A})^{-1} \mathbf{B}$  has a maximum eigenvalue of 0.814 (0.904).

**Design B (Banded specification):** We revisit simulation Case 1 in Gao et al. (2019). The matrices  $\mathbf{A}$  and  $\mathbf{B}$  are banded with a bandwidth of  $k_0 = 3$ . Specifically, the elements in the matrices  $(\mathbf{A})_{i,j=1}^N$  and  $(\mathbf{B})_{i,j}^N$  are generated according to the following three steps:

- Step 1: If  $|i - j| = k_0$ , then  $a_{ij}$  and  $b_{ij}$  are drawn independently from a uniform distribution on the two points  $\{-2, 2\}$ . All remaining elements within the bandwidth are drawn from the mixture distribution  $\omega I_{\{0\}} + (1 - \omega)N(0, 1)$  with  $\mathbb{P}(\omega = 1) = 0.4$  and  $\mathbb{P}(\omega = 0) = 0.6$ .
- Step 2: Rescale the matrices  $\mathbf{A}$  and  $\mathbf{B}$  from Step 1 to  $\eta_1 \times \mathbf{A} / \|\mathbf{A}\|_2$  and  $\eta_2 \times \mathbf{B} / \|\mathbf{B}\|_2$ , where  $\eta_1$  and  $\eta_2$  are drawn independently from  $U[0.4, 0.8]$ .
- Step 3: To avoid unstable systems, check if  $\|(\mathbf{I}_N - \mathbf{A})^{-1} \mathbf{B}\|_2 < 0.95$ . If this is not the case, then we discard the matrices  $\mathbf{A}$  and  $\mathbf{B}$  and return to Step 1.

We vary the cross-sectional dimension over  $N \in \{25, 100\}$ .

The two simulation designs are chosen to provide a balanced comparison between modelling approaches with a predetermined spatial weight matrix and modelling approaches that estimate the spatial interactions in a data-driven way.

The two methods with pre-determined spatial weight matrix are taken from Yu et al. (2008) and Lam and Souza (2019). The estimator in Yu et al. (2008) departs from a spatial-temporal autoregression (hence abbreviated as ST-AR) in which three scalar parameters and a single spatial weight matrix parametrize the spatial and temporal lag. The parameters are estimated by Gaussian quasi-maximum likelihood. As mentioned in the introduction, the method in Lam and Souza (2019), LaSo for short, allows for sparse adjustments from a predetermined spatial weight matrix. To match our simulation designs, we include  $\mathbf{y}_{t-1}$  in the matrix of covariates (see their Section 3.3.4). For both methods, we let the spatial weight matrix coincide with the “first-nearest neighbour” interactions implied by Design A. This choice clearly favours these methods in Design A while causing a misspecified model in Design B.

The second group of methods determines the spatial interactions in a data-driven way. First, we consider three implementations of our SPLASH estimator: (1) SPLASH(0,  $\lambda$ ) promotes non-sparse groups only, (2) SPLASH(0.5,  $\lambda$ ) gives equal weight to sparsity at the group and individual level, and (3) SPLASH(1,  $\lambda$ ) encourages unstructured sparsity only. In congruence with Theorems 1 and 2, we rely on banded autocovariance matrices  $\mathcal{B}_h(\hat{\Sigma}_0)$  and  $\mathcal{B}_h(\hat{\Sigma}_1)$  with bandwidths selected by the bootstrap procedure in Guo et al. (2016, p. 7). The included estimators from Gao et al. (2019) are references by  $\text{GMWY}(k)$  and  $\text{GMWY}(k_0)$ . The first of these estimators implements generalized Yule-Walker estimation for banded  $\mathbf{A}$  and  $\mathbf{B}$  using the selection rule in (2.17) of their paper. Having observed a rather poor performance of this method during some exploratory simulations for Design B, we provide the true bandwidths of  $\mathbf{A}$  and  $\mathbf{B}$  to the (infeasible) estimator denoted  $\text{GMWY}(k_0)$ . To allow for a comparison with the original simulation results, neither of these GMWY estimators applies banding to the covariance matrix estimators  $\hat{\Sigma}_0 = \frac{1}{T} \sum_{t=2}^T \mathbf{y}_t \mathbf{y}_t'$  and  $\hat{\Sigma}_1 = \frac{1}{T} \sum_{t=2}^T \mathbf{y}_t \mathbf{y}_{t-1}'$ . During additional exploratory simulations, we have seen that the performance of the GMWY estimator does not change much if banded covariance matrices would be used.

Finally, we also include an  $L_1$ -penalized reduced form VAR(1) estimator (abbreviated PVAR). In detail, we consider the reduced form VAR(1) specification  $\mathbf{y}_t = \mathbf{C} \mathbf{y}_{t-1} + \mathbf{u}_t$  and estimate  $\mathbf{C}$  by minimizing  $\mathcal{L}_{\text{par}}(\mathbf{C}) = \sum_{t=2}^T \|\mathbf{y}_t - \mathbf{C} \mathbf{y}_{t-1}\|_2^2 + \lambda \sum_{i,j=1}^N |c_{ij}|$ . This estimator is well-researched in the literature (see, e.g. Kock and Callot, 2015; Gelper et al., 2016; Masini et al., 2022), albeit in different settings. The inclusion of this estimator serves to uncover any difference in predictive performance stemming from the estimation of a spatio-temporal representation, as opposed to estimating a purely temporal specification.

The forecasting performance of each estimator will be assessed using the *Relative Mean-Squared Forecast Error (RMSFE)*. Using a superscript  $j$  to index a specific Monte Carlo replication, the RMSFE is calculated as

$$\text{RMSFE} = \frac{\sum_{j=1}^{N_{sim}} \|\mathbf{y}_{T+1}^j - \hat{\mathbf{C}}^j \mathbf{y}_T^j\|_2^2}{\sum_{j=1}^{N_{sim}} \|\mathbf{y}_{T+1}^j - \mathbf{C} \mathbf{y}_T^j\|_2^2} \tag{14}$$

The estimation accuracy of the spatio-temporal methods is analysed via the mean estimation error (MEE). Using the superscript  $j$  as before, the (MEE) is defined as

$$\text{MEE} = \frac{1}{N_{sim}} \sum_{j=1}^{N_{sim}} \left\| \left( \hat{\mathbf{A}}^j, \hat{\mathbf{B}}^j \right) - \left( \mathbf{A}^j, \mathbf{B}^j \right) \right\|_2 \tag{15}$$

To provide some insights into the variability in estimation accuracy across simulation trials, we additionally report the standard errors of the MEEs, which we compute as

$$\text{SE}(\text{MEE}) = \sqrt{\frac{1}{N_{sim} - 1} \sum_{j=1}^{N_{sim}} \left( \left\| \left( \hat{\mathbf{A}}^j, \hat{\mathbf{B}}^j \right) - \left( \mathbf{A}^j, \mathbf{B}^j \right) \right\|_2 - \text{EE} \right)^2} \tag{16}$$

Next, a word on the determination of hyperparameters. For each regularized estimator, we construct a unique grid containing  $K = 20$  penalty values that are evenly spaced on a logarithmic scale. The maximum value in each grid,  $\lambda_{\max}$ , is the smallest value that produces the zero solution (see step 4 in Algorithm 1 for its expression for SPLASH). Given  $\lambda_{\max}$ , we define the smallest penalty as  $\lambda_{\min} = \eta \lambda_{\max}$  with  $\eta = 10^{-4}$  ( $\eta = 10^{-6}$ ) for Design A (B). Estimating SPLASH solutions for each  $\lambda \in \{\lambda_1, \dots, \lambda_{20}\}$ , each  $\alpha \in \{0, 0.5, 1\}$ , and each individual simulation trial remains too computationally expensive (especially for large  $N$ ). We instead perform a small-scale preliminary analysis in which we draw a small set of simulations from Designs A and B on which we estimate all solutions for a given value of  $T$ . Then, we choose the order  $i_T \in \{1, \dots, 20\}$  that minimizes the RMSFE in this preliminary set of simulations. This process of choosing the order  $i_T$  on a log-equidistant grid for each value of  $T$ , is equivalent to setting  $\lambda = m_T \lambda_{\max}$  with  $m_T = 10^{-4(i_T-1)/20}$  or  $m_T = 10^{-6(i_T-1)/20}$  for designs A and B, respectively. For Design A (B), our selected orders for  $T = \{500, 1000, 2000\}$  are  $i_T = \{9, 10, 11\}$  ( $i_T = \{10, 11, 12\}$ ), corresponding to  $m_T \approx 0.025, 0.015, 0.01$  ( $m_T \approx 0.002, 0.001, 0.0005$ ), respectively. Having fixed the preferred order or multiplier, it suffices to compute the SPLASH solution for each  $\alpha \in \{0, 0.5, 1\}$ . The penalized VAR method is computationally less expensive. Accordingly, we choose its penalty parameter based on a time series cross-validation (TSCV) scheme (e.g. Hyndman and Athanasopoulos, 2018). Since a preliminary analysis showed that the penalty selection approach used for SPLASH resulted in sub-par performance for the LaSo method, we choose the latter's penalty via TSCV as well. In our implementation of TSCV, the first 80% of the data is used to fit multiple solutions on, which are then evaluated based on the MSFE obtained on the latter 20% of the data. The preferred penalty is chosen as the solution that attains the smallest MSFE.<sup>3</sup>

#### 4.2. Simulation results

The simulation outcomes for Design A and B are reported in Table 1. For each design we report the predictive performance in terms of the RMSFE and the estimation accuracy as measured by the mean estimation error (MEE). Whenever applicable, the standard errors of the latter quantity is written in parentheses.

First, we consider the predictive performance under Design A. For all methods, we observe a monotonic decrease in RMSFE when  $T$  increases. Unsurprisingly, the best performance is obtained by the ST-AR method which only needs to estimate two non-zero parameters to recover the DGP. More impressively, with a maximum RMSFE of 1.013, all three SPLASH implementations are close to oracle performance. This demonstrates that the cost incurred by estimating the DGP in a fully data-driven way, as opposed to relying on a correctly specified pre-determined weight matrix, is small. There is little difference across the three SPLASH implementations, with the effect of group-level regularization providing noteworthy improvements only in the most high-dimensional setting ( $N = 100$  and  $T = 500$ ). LaSo obtains a strong predictive accuracy as well, albeit slightly worse than SPLASH when  $T$  is smaller. This is somewhat surprising as the LaSo specification has all the correct spatial interactions readily included in its spatial weight matrix and simply needs to estimate the redundant adjustment matrix as zero. However, this requires a sufficient high penalty to be selected via time series cross-validation which was not the case across simulation trials. Additionally, the LaSo estimator does not always converge and remains rather sensitive to its initialization. The RMSFEs for the GMWY estimators are strikingly high. In the setting  $N = 25$  and  $T = 500$ , the GMWY estimator frequently underestimates the bandwidth by choosing it to equal 1. This causes inferior performance across all metrics. The GMWY( $k_0$ ) estimator, on the other hand, is based on the correct bandwidth, but its forecast performance is worse. Upon closer inspection, we find that these high RMSFEs are caused by several extreme prediction errors. These prediction outliers occur during simulation trials in which the smallest eigenvalue of the estimated matrix  $I - \hat{A}$  is close to zero. As the forecasts are based on the reduced form representation and hence  $(I - \hat{A})^{-1}$ , the GMWY estimator is prone to stability issues when the bandwidth is large relative to the dimension. For  $N = 25$  and  $T = 2000$ , the bandwidth selection in GMWY improves, while its forecast performance ironically worsens as a result of increasing stability issues. Finally, the penalized VAR produces stable forecasts but is outperformed by ST-AR and SPLASH across all settings. A correctly specified spatio-temporal model can thus outperform a reduced form VAR in terms of forecast performance.

Paralleling the results on the predictive performance, the ST-AR attains the lowest (variability in) MEE, followed by the SPLASH estimator. The MEE shows the added benefit of group-level regularization in SPLASH as the implementations for  $\alpha = 0$  and  $\alpha = 0.5$  attain substantially lower estimation errors. For the SPLASH(0.5,  $\lambda$ ) estimator, Fig. 3 illustrates the similarity between the true  $A$  and its estimated counterpart. Indeed, owing to the imposed group-level regularization, the true sparsity pattern in  $A$  is clearly present in the (averaged) SPLASH estimates. Continuing the comparison, the LaSo method displays mixed performance, varying from very large ( $N = 25$  and  $T = 500$ ) to very small ( $N = 100$  and  $T = 200$ ) MEEs. This again reflects the inability to consistently select a sufficiently large penalty parameter, combined with the aforementioned issues with numerical convergence. Finally, GMWY( $k$ ) typically obtains a lower MEE than GMWY( $k_0$ ) as a result of choosing a lower bandwidth. Apparently, estimating a misspecified model by choosing  $k < k_0$  results in a favourable bias–variance trade-off, highlighting the importance of regularization in high-dimensional settings.

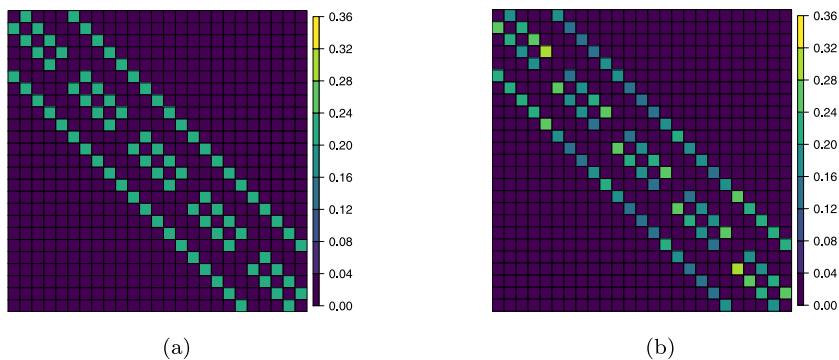
Finally, we look at Design B. The methods that rely on a prespecified spatial weight matrix no longer have a competitive advantage as its interactions are misspecified. Consequently, the SPLASH estimators and the infeasible GMWY( $k_0$ ) estimator now exhibit the best overall forecast performance. The feasible GMWY( $k$ ) estimator shows good performance for  $N = 100$  and  $T \in \{1000, 2000\}$  because the increase in spatial units and time points improves bandwidth selection and diminishes invertibility issues. Performance differences between the SPLASH estimators are generally small as no specific type of sparsity dominates. That is, the group penalty will help in setting to zero the elements outside the bandwidth  $k_0$  and the individual coefficient penalties are

<sup>3</sup> We also tried to select the penalty as the sparsest solution whose prediction error lies within one standard error of the minimum prediction error. This selection rule, however, did not lead to an improvement in forecast or estimation accuracy.

**Table 1**  
Simulation results for Design A and B.

Metric	N	T	SPLASH			ST-AR	LaSo	GMWY		PVAR	
			$\alpha = 0$	$\alpha = 0.5$	$\alpha = 1$			k	$k_0$		
Design A (Spatial grid with neighbour interactions)											
RMSFE	25	500	1.012	1.011	1.012	1.001	1.109	9.924	539.490	1.108	
		1000	1.004	1.004	1.005	1.000	1.040	29.538	66.157	1.067	
		2000	1.005	1.005	1.004	1.000	1.003	894.224	527.024	1.042	
	100	500	1.013	1.014	1.022	1.000	1.526	1.151	463.714	1.158	
		1000	1.010	1.011	1.012	1.000	1.000	1.104	137.901	1.109	
		2000	1.005	1.004	1.004	1.000	1.000	1.087	2.300	1.077	
MEE	25	500	0.334 (0.015)	0.342 (0.015)	0.477 (0.022)	0.019 (0.001)	3.381 (0.164)	1.162 (0.052)	4.448 (0.212)		
		1000	0.283 (0.013)	0.281 (0.013)	0.385 (0.017)	0.013 (0.001)	1.416 (0.081)	1.173 (0.053)	4.062 (0.187)		
		2000	0.243 (0.011)	0.232 (0.010)	0.295 (0.013)	0.009 (0.000)	0.207 (0.022)	1.212 (0.054)	3.956 (0.183)		
		500	0.393 (0.018)	0.411 (0.019)	0.573 (0.026)	0.009 (0.000)	0.768 (0.082)	1.119 (0.050)	2.174 (0.098)		
	100	1000	0.367 (0.017)	0.379 (0.017)	0.537 (0.024)	0.007 (0.000)	0.027 (0.002)	1.066 (0.048)	2.082 (0.094)		
		2000	0.373 (0.017)	0.362 (0.016)	0.504 (0.023)	0.005 (0.000)	0.019 (0.001)	1.046 (0.047)	2.059 (0.093)		
		500	1.025	1.026	1.029	1.154	1.495	5.485	1.048	1.125	
		1000	1.011	1.011	1.012	1.159	1.444	1.387	1.016	1.116	
	RMSFE	100	2000	1.007	1.007	1.008	1.156	2.699	1.008	1.008	1.109
			500	1.052	1.058	1.074	1.140	1.316	1.051	1.037	1.110
			1000	1.025	1.028	1.035	1.136	14.308	1.017	1.016	1.104
			2000	1.017	1.018	1.021	1.135	6.929	1.009	1.009	1.102
MEE	25	500	0.723 (0.033)	0.761 (0.035)	0.883 (0.040)	0.750 (0.034)	1.338 (0.071)	24.662 (5.656)	1.222 (0.056)		
		1000	0.656 (0.030)	0.687 (0.032)	0.798 (0.037)	0.744 (0.034)	1.386 (0.073)	4.272 (0.923)	1.134 (0.052)		
		2000	0.600 (0.028)	0.617 (0.029)	0.709 (0.033)	0.750 (0.034)	1.426 (0.076)	1.829 (0.651)	1.064 (0.049)		
		500	0.831 (0.038)	0.894 (0.041)	1.075 (0.049)	0.735 (0.033)	1.203 (0.064)	1.027 (0.049)	0.799 (0.036)		
	100	1000	0.830 (0.038)	0.889 (0.041)	1.051 (0.048)	0.733 (0.033)	1.143 (0.060)	0.755 (0.035)	0.707 (0.032)		
		2000	0.795 (0.037)	0.849 (0.039)	0.994 (0.046)	0.735 (0.033)	1.203 (0.064)	0.621 (0.028)	0.617 (0.028)		

**Note:** The relative mean-squared forecast error (RMSFE) and mean estimation errors (MEE) are defined in (14) and (16), respectively. In general, lower numbers indicate better performance. As PVAR estimates a reduced form VAR, there are no model errors for **A** and **B** to report for this method. Standard errors of the model errors are listed in parentheses.



**Fig. 3.** Visualizations of the true and estimated spatial weight matrix **A** for Design B. **(a)** The true spatial weight matrix **A** implied by the (5 × 5) spatial grid design (Design B with  $m = 5$ ). **(b)** The average absolute values of the entries in  $\hat{\mathbf{A}}$  as computed by SPLASH(0.5,  $\lambda$ ) for  $N = 25$  and  $T = 1000$ . That is, the  $(i, j)$ th entry in the matrix on the right equals  $\frac{1}{N_{sim}} \sum_{k=1}^{N_{sim}} |\hat{a}_{ij}^k|$  with  $\hat{a}_{ij}^k$  being the estimated  $(i, j)$ th entry of **A** in the  $k$ th Monte Carlo replication.



useful to reveal the random zero pattern that occurs close to the main diagonals of  $A$  and  $B$ . Only being able to indirectly incorporate contemporaneous interactions, the reduced form PVAR estimator performs neither badly nor among the top contenders. Finally, in terms of estimation accuracy, SPLASH appears to obtain the lowest MEE for  $N = 25$ , but is surpassed by  $GMWY(k_0)$  for  $N = 100$ . Surprisingly, ST-AR performs very competitive in terms of estimation accuracy, reflecting a beneficial bias–variance trade-off by estimating a misspecified low-dimensional approximation to the DGP.

## 5. Empirical application

In this section we consider the application of SPLASH to satellite-measured data. Section 5.1 serves as a general guide to practitioners. It explains the important steps in the analysis and provides additional details on the satellite data collection and processing. A concrete illustration is presented in Section 5.2 where we predict daily nitrogen dioxide ( $\text{NO}_2$ ) concentrations over Greater London.

### 5.1. Applying SPLASH: A guide to practitioners

To facilitate the adoption of our SPLASH modelling strategy, we explain in detail how to download, process and model data from the Copernicus satellite. A readily available implementation of each step, together with the raw data used in the preceding application, is available at <https://sites.google.com/view/etiennewijler/code> in the form of separate code scripts and an R package.

**Step 1: Data collection** Satellite images are downloadable from the Copernicus Open Access Hub.<sup>4</sup> First-time users are required to register a free (at the time of writing) account. Single files can be downloaded interactively via a user interface, whereas repeated requests are best handled via the API hub. After defining the area of interest (AOI) based on GPS coordinates, all satellite images that overlap with this AOI can be downloaded. The AOI is subsequently cropped out and the original image is deleted to free storage space. Sentinel-5P Copernicus data products include ozone, sulphur dioxide, nitrogen dioxide, carbon monoxide, formaldehyde and methane. The user might also use satellite data from a different data source.

**Step 2: Data processing** The cropped raw satellite data is mapped onto the spatial grid. All pixels within a given grid cell are averaged into a single value. Enumerating the cells consecutively (by rows or by columns), the vector  $y_t$  with daily observations is constructed. This vector is likely to have some missing values as the satellite collects data while orbiting the earth. These missing observations should be imputed, e.g. using *Multivariate Time Series Data Imputation (mtsdi)* R package. This imputation method in [Junger and Ponce De Leon \(2015\)](#) was designed to impute missing values in air pollution data. Finally, one should assess/model the low-frequency seasonality in the data. That is, applying model (1) or (10) on daily data does not capture the slow concentration variation with the seasons due to the difference in frequency. We provide a possible approach in the next section.

**Step 3: Data modelling** The SPLASH R package fully automates the model-building process by estimating the bandwidth and banded covariance matrices, calculating the quantities of interest, and selecting the optimal penalty parameter by time series cross-validation (all steps in Algorithm 1). The sparsely estimated coefficients give an idea of the important interactions and the implied reduced form VAR can be used for out-of-sample predictions.

### 5.2. Predicting $\text{NO}_2$ concentrations based on satellite data

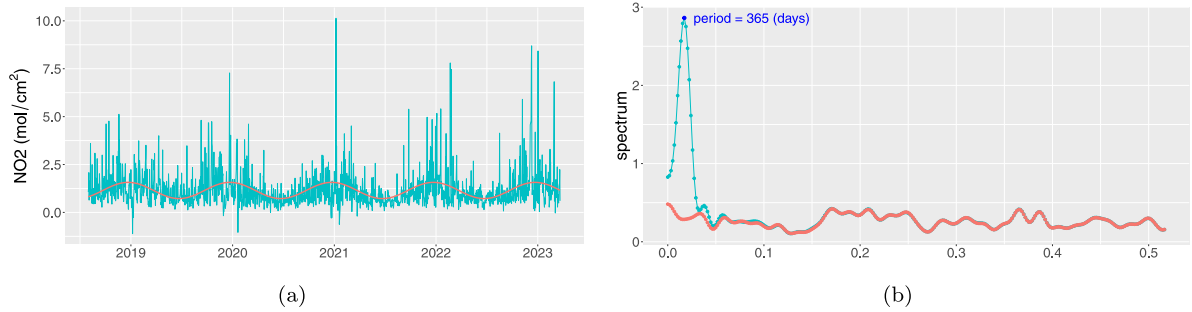
$\text{NO}_2$  is emitted during combustion of fossil fuels (e.g. by motor vehicles) and it has been associated with adverse effects on the respiratory system.<sup>5</sup> The Air Quality Standards Regulations 2010 requires a regular monitoring of  $\text{NO}_2$  concentration levels in the UK.<sup>6</sup> Using satellite data, we examine the empirical performance of the SPLASH estimator when predicting daily  $\text{NO}_2$  concentrations in Greater London. The data is available at the Copernicus Open Access Hub and we consider the time span from 1 August 2018 to 31 March 2023. The original  $\text{NO}_2$  concentrations are reported in  $\text{mol}/\text{m}^2$ , which we convert to  $\text{mol}/\text{cm}^2$  to avoid numerical instabilities caused by small-scale numbers. The far majority of measurements are captured between 11:00 and 14:00 UTC. As displayed in Fig. 5(c), the area of interest is mapped onto a  $(6 \times 7)$  grid with even straight-line distances between the centres of horizontally (8.57 km) and vertically (8.53 km) neighbouring grid cells. We average pixels and impute missing observations as in Step 2 of our guide to practitioners.

A first inspection of the data reveals several noteworthy features. First, the average  $\text{NO}_2$  concentrations (averaged over space and time) is  $1.1689 \text{ mol}/\text{cm}^2$ . Furthermore, there is substantial regional variation in the average daily  $\text{NO}_2$  concentrations. The lowest daily average concentration of  $0.9671 \text{ mol}/\text{cm}^2$  is observed over Chipping Ongar (a small market town North-East of the Greater London Area) and the highest concentration of  $1.3479 \text{ mol}/\text{cm}^2$  is measured at Westminster (the bustling government area near Buckingham Palace).

<sup>4</sup> <https://www.copernicus.eu/en/access-data/conventional-data-access-hubs>

<sup>5</sup> The direct health effect of nitrogen dioxide is difficult to determine because its emission process is typically accompanied with the emission of other air pollutants (see, e.g. [Brunekreef and Holgate, 2002](#)).

<sup>6</sup> Source: <https://www.legislation.gov.uk/uksi/2010/1001/contents/made>.



**Fig. 4.** The seasonality in the data. (a) The daily spatially averaged NO<sub>2</sub> concentration (green) and its first-order Fourier approximation at the angular frequency of  $\omega_{yearly} = \frac{2\pi}{355} \approx 0.0172$  (red). (b) The spectral density estimate of the spatially averaged NO<sub>2</sub> concentration before (green) and after (red) subtracting the first-order Fourier approximation. The spectral densities are computed using the Bartlett kernel (e.g. (6.2.15) in Hamilton, 1994) with a bandwidth of  $q = 400$ . (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

We proceed to the final part of Step 2 in our guide to practitioners and investigate the seasonality in the data. First, in Fig. 4(a) we visualize the spatially averaged NO<sub>2</sub> concentration over time. The plot reveals that this concentration is typically highest in December and lowest in June thereby pointing towards seasonality. To get an idea about the relevant frequencies, we look at the (estimated) spectral density. The clear peak at the angular frequency of  $\omega_{yearly} = \frac{2\pi}{355} \approx 0.0172$  motivates us to include sine and cosine signals with this baseline period of 365 days. As the residuals after the inclusion of a single Fourier term already produce a rather flat spectrum (Fig. 4(b)) and further Fourier terms hardly change the spectrum, we decide to detrend the data with this first-order Fourier approximation throughout this empirical application. The Fourier approximation itself is displayed as the red line in Fig. 4(a).

A rolling-window approach is used to assess the predictive power of the SPLASH estimator. Motivated by our seasonality analysis, we assume that the data is generated by the following model:

$$\begin{aligned}
 z_t &= d_t + y_t, \\
 d_t &= \mu + \iota_N \left( \sum_{j=1}^M a_j \cos(j\omega_{yearly}t) + b_j \sin(j\omega_{yearly}t) \right), \\
 y_t &= \mathbf{A}y_t + \mathbf{B}y_{t-1} + \epsilon_t,
 \end{aligned}$$

where  $\iota_N$  is an  $N$ -dimensional vector of ones. As mentioned in our spectral analysis, we select  $M = 1$  (but also remind the reader that a different data set might require additional Fourier terms). Hence, each spatial unit is assigned its own mean and a common seasonal term. Each window contains 2 years of the data (730 days) allowing 964 one-step ahead forecasts to be made. For each window, we proceed along the following four steps: (i) regress out the deterministic component  $d_t$  to obtain the de-seasonalized data  $\hat{y}_t = z_t - \hat{d}_t$ , (ii) estimate the (hyper)parameters of each model based on  $\hat{y}_t$ , (iii) produce separate one-step ahead forecasts for the de-seasonalized data and the deterministic components, and (iv) add the forecasts together to obtain a forecast for the NO<sub>2</sub> concentration observed directly after the window.

We consider the same type of estimators as previously seen in the simulations. As before, we include SPLASH(0,  $\lambda$ ), SPLASH(0.5,  $\lambda$ ), and SPLASH(1,  $\lambda$ ). With the current spatial grid of  $N = 6 \times 7 = 42$  spatial units, these methods potentially require  $2N^2 - N = 3486$  parameters. However, for identifiability (see Assumption 3), we band the spatial matrix  $\mathbf{A}$  and autoregressive matrix  $\mathbf{B}$  such that  $a_{ij} = b_{ij} = 0$  for  $|i - j| > \lfloor (N - 1)/4 \rfloor = 10$ . Similar to the simulation section, we select the penalty parameter as  $\lambda = m\lambda_{max}$  with the multiplicative factor  $m$  being determined based on a preliminary set of out-of-sample forecasts. Having no a priori information concerning the spatial interactions, we implement the ST-AR and LaSo methods with two spatial weight matrices: the “first-nearest neighbour” specification is identical to the spatial grid (SG) in Design A of Section 4, and the (normalized) inverse distance (ID) spatial weight matrix  $\mathbf{W} = (w_{ij})_{i,j=1}^N$  with each element  $w_{ij}$  representing the inverse distance in kilometres between spatial units  $i$  and  $j$ . The GMWY estimator is implemented with bandwidth  $k = 10$ , as this provides the best (i.e. least unstable) predictions. Penalty selection for LaSo and PVAR follow the same TSCV scheme discussed in Section 4.1.

The forecast performance is measured along three metrics. We report: (i) the average loss relative to PVAR, (ii) the number of spatial units that are predicted more accurately than the PVAR benchmark (#wins), and (iii) the inclusion in the Model Confidence Set (MCS) by Hansen et al. (2011) at a 10% significance level. These three metrics are calculated using two loss functions for the forecast errors: the mean squared forecast error (MSFE) and the mean robust forecast error (MRFE). That is, letting  $\hat{\epsilon}_{i,T_0+h}^{(j)}$  denote the one-step ahead forecast error of estimator  $j$  for the outcome of the  $i$ th spatial unit at period  $T_0 + h$ , we calculate

$$\text{MSFE}_j = \frac{1}{NH} \sum_{i=1}^N \sum_{h=1}^H \left( \hat{\epsilon}_{i,T_0+h}^{(j)} \right)^2, \text{ and } \text{MRFE}_j = \frac{1}{NH} \sum_{i=1}^N \sum_{h=1}^H \frac{\left( \hat{\epsilon}_{i,T_0+h}^{(j)} \right)^2}{1 + \left( \hat{\epsilon}_{i,T_0+h}^{(j)} \right)^2}, \tag{17}$$

**Table 2**  
Forecast performance of various methods for NO<sub>2</sub> satellite data on a (6 × 7) grid of observations.

	MSFE			MRFE		
	RMSFE	#wins	MCS	RMRFE	#wins	MCS
SPLASH(0, $\lambda$ )	<b>0.448</b>	41	Yes	0.885	41	Yes
SPLASH(0.5, $\lambda$ )	0.448	41	Yes	0.886	41	Yes
SPLASH(1, $\lambda$ )	0.450	41	Yes	0.889	40	Yes
ST-AR(SG)	0.452	40	Yes	0.885	41	Yes
ST-AR(ID)	0.450	42	Yes	<b>0.885</b>	39	Yes
LaSo(SG)	0.503	27	Yes	0.892	36	Yes
LaSo(ID)	0.464	33	Yes	0.890	37	Yes
GMWY	36.056	0	No	1.759	0	No
PVAR	1.000	–	No	1.000	–	No

**Note:** Number of grid points (out of  $N = 45$ ) with lower prediction errors (#wins) and significantly lower prediction errors (#sign. wins) compared to the  $L_1$ -penalized reduced form VAR(1) estimator (PVAR). Relative MSFE and MRFE are abbreviated by RMSFE and RMRFE, respectively. Values below (above) 1 indicate superior (inferior) performance compared to PVAR.

with  $H = T - T_0$  the maximum forecast horizon. The MRFE is included to reduce the importance of incidental outliers caused by several abrupt spikes in the NO<sub>2</sub> column densities and/or the stability issues that continued to plague the GMWY method (see Section 4 for explanations). In Table 2 we report the *relative* MSFE (RMSFE) and *relative* MRFE (RMRFE) as the performance in comparison to the PVAR estimator.

We first look at the columns under MSFE. The RMSFEs show that all spatio-temporal methods, with the exception of GMWY, improve substantially over the purely penalized VAR benchmark. Incorporating spatial effects is thus important while modelling NO<sub>2</sub> concentrations empirically. The lowest MSFE is obtained by SPLASH(0,  $\lambda$ ) but the differences in forecasting performance in comparison to ST-AR and LaSo are small. For GMWY, we again observed significant stability issues when inverting the estimated spatial weight matrix to obtain the reduced form representation that is necessary to construct the forecast, resulting in inferior forecasts. We were not able to resolve this by choosing a smaller bandwidth, possibly because such specifications may start to omit important spatial interactions. The SPLASH implementations predict 41 out of 42 spatial units more accurately than the PVAR benchmark, second only to the ST-AR (ID) approach which manages to improve the predictions for all spatial units on the grid. The LaSo approach performs less favourable here, beating the benchmark for only 27 or 33 spatial units based on the SG and ID spatial weight matrix, respectively. Finally, the MCS finds statistically significant evidence that the GMWY and PVAR methods obtain inferior forecasts, but fails to find any significant performance differences among the remaining methods.

Next, we focus on the mean robust forecast error (MRFE). The results are qualitatively similar. The RMRFEs for SPLASH, ST-AR and LaSo lie even closer together and are practically indistinguishable. The GMWY still displays the worst forecast performance, although the difference appears less dramatic as a result of the robust loss metric. SPLASH( $\alpha$ ,  $\lambda$ ) with  $\alpha = 0, 0.5$  and the ST-AR(SG) method are able to beat the PVAR benchmark most often by obtaining lower MRFES for 41 out of 42 spatial units. The MCS again includes all methods except GMWY and PVAR.

We end this application with a visual analysis of the estimated spatial interactions between neighbouring regions in Greater London in Fig. 5. To have both individual and group sparsity, we focus on the SPLASH(0.5,  $\lambda$ ) estimates. First, the window-averaged absolute magnitude of the spatial interactions is displayed in Fig. 5(a). A clear diagonal pattern emerges with the largest interactions clustering on the two diagonals closest to the main diagonal and the two outer diagonals that start at the (1, 7) and (7, 1) elements. These four diagonals correspond to first-order vertical and horizontal interactions, respectively. For each element in  $A$ , the heat map in Fig. 5(b) indicates the fraction of rolling windows with a nonzero estimated coefficient. The aforementioned interactions between first-order horizontal and vertical neighbours are selected in all windows. Additional interactions between diagonal neighbours are relevant in approximately half of the rolling windows. To facilitate the interpretation of this sparsity pattern, we provide a spatial plot of our region of interest with the spatial grid overlaid in Fig. 5(c). For spatial unit 21 specifically, we visualize the spatial interactions that are estimate to be nonzero in at least 20% of the rolling windows and vary the thickness of the arrows to indicate their (average) estimated magnitudes. The strongest interactions observed are between first-order vertical neighbours, followed by those between first-order horizontal neighbours. Interactions between diagonal and second-order vertical neighbours are observed as well, although their average effects are substantially weaker. The fact that the diagonal interactions only seem to occur along the diagonal from South-West to North-East may be a consequence of the predominant South-West winds over London. Overall, the intuitive sparsity patterns that arise, in combination with the improvement in forecast performance, are encouraging and provide empirical validation for the use of SPLASH on spatial data, especially when the spatial units follow a natural ordering on a spatial grid. The knowledge about the relevant interactions as well as their magnitudes might be exploited by policymakers to study the effects of for example stricter emission regulations and/or traffic decrease on regional air quality. However, the analysis of these spatial impulse response would require additional results on statistical inference (see Section 3.2.4).

## 6. Conclusion

In this paper, we develop the Spatial Lasso-type Shrinkage (SPLASH) estimator, a novel estimation procedure for high-dimensional spatio-temporal models. The SPLASH estimator is designed to promote the recovery of structured forms of sparsity

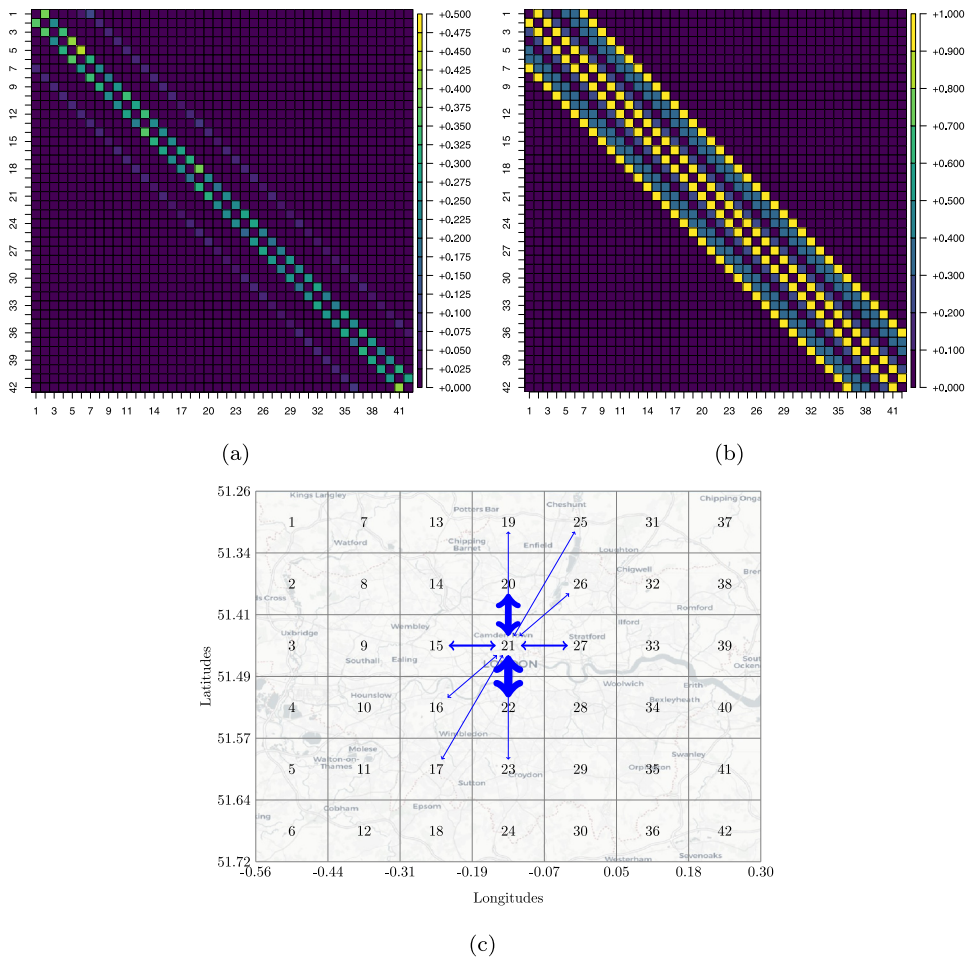


Fig. 5. Illustrations of the sparsity patterns in the estimated  $(42 \times 42)$  matrix  $A$  based on rolling window samples. (a) The average absolute value of the entries in  $A$  as averaged across each one-step ahead forecast. (b) The proportion of rolling window forecast with an estimated coefficient being unequal to zero. (c) The spatial interaction for spatial unit 21 that are non-zero in at least 20% of the forecasts.

without imposing such structure a priori. We derive consistency of our estimator in an asymptotic framework in which the number of both spatial units and temporal observations diverge jointly. To solve the identifiability issue, we rely on a relatively non-restrictive assumption that the coefficient matrices in the spatio-temporal model are sufficiently banded. Based on this assumption, we consider banded estimation of high-dimensional spatio-temporal autocovariance matrices, for which we derive novel convergence rates that are likely to be of independent interest. The SPLASHX extension explains how to include exogenous variables. As an application, we use SPLASH to predict satellite-measured  $\text{NO}_2$  concentrations in London. We find evidence for spatial interactions between neighbouring regions. In addition, our estimator obtains superior forecast accuracy compared to a number of competitive benchmarks, including the recently introduced spatio-temporal estimator by Gao et al. (2019) that inspired the development of SPLASH.

### Acknowledgements

This paper (or earlier versions hereof) has been presented during the internal seminar of the Quantitative Economics department of Maastricht University, the Econometrics Internal Seminar (EIS) at Erasmus University Rotterdam, the Bernoulli-IMS One World Symposium, the workshop on Dimensionality Reduction and Inference in High-dimensional Time Series at Maastricht University, the 2021 Annual Conference of the International Association for Applied Econometrics (IAAE), the 5th Conference on Econometric Models of Climate Change, the internal seminar at Tor Vergata University of Rome, and the 2021 (EC)<sup>2</sup> Conference. We gratefully acknowledge comments and feedback from the participants. Suggestions by Stephan Smeekes and Ines Wilms were particularly helpful, so we thank them explicitly. In addition, we thank the two anonymous referees for their constructive feedback that has resulted in major improvements in our work. All remaining errors are our own.

**Appendix A. Lemmas**

The proofs of the lemmas in this section are available in the supplementary material.

**Lemma 1.** Define the quantities  $N_c = |c|$ ,  $S = \{j : c_j \neq 0\}$ ,  $\mathcal{G}_S = \{g \in \mathcal{G} : c_g \neq \mathbf{0}\}$ ,  $\mathcal{G}_S^c = \{g \in \mathcal{G} : c_g = \mathbf{0}\}$ ,  $\bar{\omega}_\alpha = \max \left\{ (1 - \alpha) \sum_{g \in \mathcal{G}_S} \sqrt{|g|}, \alpha \sqrt{|S|} \right\}$  and consider

$$\mathbf{A} \in \mathcal{C}_{N_c}(\mathcal{G}, S) := \{ \mathbf{A} \in \mathbb{R}^{N_c} : P_{\alpha, S^c}(\mathbf{A}) \leq 3P_{\alpha, S}(\mathbf{A}) \},$$

where

$$P_{\alpha, S^c}(\mathbf{A}) = (1 - \alpha) \sum_{g \in \mathcal{G}_S^c} \sqrt{|g|} \|A_g\|_2 + \alpha \|A_{S^c}\|_1, \text{ and}$$

$$P_{\alpha, S}(\mathbf{A}) = (1 - \alpha) \sum_{g \in \mathcal{G}_S} \sqrt{|g|} \|A_g\|_2 + \alpha \|A_S\|_1.$$

Then, under Assumption 4, it holds that

$$\min_{\mathbf{A} \in \mathcal{C}_{N_c}(\mathcal{G}, S)} \frac{\bar{\omega}_\alpha \|V^{(d)} \mathbf{A}\|_2}{P_{\alpha, S}(\mathbf{A})} \geq \frac{\phi_0}{2}. \tag{A.1}$$

**Lemma 2.** Define the set  $\mathcal{V}(x) := \left\{ \left\| \hat{V}_h - V \right\|_2 \leq x \right\}$ . Then, under Assumption 4, it holds on  $\mathcal{V} \left( \frac{\phi_0}{4} \right)$  that

$$\min_{\mathbf{x} \in \mathcal{C}_{N_c}(\mathcal{G}, S)} \frac{\bar{\omega}_\alpha \left\| \hat{V}_h \mathbf{x} \right\|_2}{P_{\alpha, S}(\mathbf{x})} \geq \frac{\phi_0}{4}.$$

**Lemma 3.** Define the  $N(p + 1)$  vector  $\xi_t = (y'_t, y'_{t-1}, \dots, y'_{t-p})'$ . For all  $j = 1, \dots, p$ , the elements of  $\hat{\Sigma}_j = \frac{1}{T} \sum_{t=p+1}^T y_t y'_{t-j}$  can be expressed as  $\frac{1}{T} \sum_{t=p+1}^T \xi_{it} \xi_{jt}$  after an appropriate choice of  $(i, j)$ .

(a) If Assumptions 1 and 2(b1) (polynomial tails) hold, then we have:

$$\mathbb{P} \left( \left| \sum_{t=p+1}^T \xi_{it} \xi_{jt} - \mathbb{E}(\xi_{it} \xi_{jt}) \right| > T\epsilon \right) \leq \left[ b_1 T^{(1-\delta)/3} + \frac{b_3}{\epsilon} \right] \exp \left( -\frac{T^{(1-\delta)/3}}{2b_1^2} \right) + \frac{b_2}{\epsilon^d T^{\frac{\delta}{2}(d-1)}},$$

for some  $0 < \delta < 1$ .

(b) If Assumptions 1 and 2(b2) (subexponential tails) hold, then we have:

$$\mathbb{P} \left( \left| \sum_{t=p+1}^T \xi_{it} \xi_{jt} - \mathbb{E}(\xi_{it} \xi_{jt}) \right| > T\epsilon \right) \leq \left[ \frac{\kappa_1}{\epsilon} + \frac{2(T\epsilon^2)^{1/7}}{\kappa_2} \right] \exp \left( -\frac{(T\epsilon^2)^{1/7}}{\kappa_3} \right).$$

Explicit expressions for the constants  $b_1, b_2, b_3, \kappa_1, \kappa_2$  and  $\kappa_3$  are provided in the proof.

**Appendix B. Proofs of main results**

**Proof of Theorem 1.** We first prove various intermediate results, see (a)–(d) below. We afterwards combine these results and recover Theorem 1.

(a) The matrix  $\tilde{C}_s = \left( \sum_{j=0}^{s-1} A^j \right) \mathbf{B} =: \tilde{D}_s \mathbf{B}$  has a maximum bandwidth of  $(s + 1)(k_0 - 1) + 1$  and satisfies

$$\left\| \tilde{C}_s - \mathbf{C} \right\|_{\text{F}} \leq \delta_C \delta_A^s.$$

(b) Define  $\Sigma_0^{r,s} = \sum_{j=0}^r \tilde{C}_s^j \tilde{D}_s \Sigma_\epsilon \tilde{D}_s' (\tilde{C}_s^j)'$  with  $\tilde{C}_s$  as in Theorem 1(a). The matrix  $\Sigma_0^{r,s}$  is a banded matrix with bandwidth no larger than  $2(rs + r + s)(k_0 - 1) + 2l_0 + 1$ . Moreover,

$$\left\| \Sigma_0^{r,s} - \Sigma_0 \right\|_{\text{F}} \leq \frac{8C_\epsilon \delta_A^s}{(1 - \delta_A)^2 (1 - \delta_C^2 (1 + \delta_A^s))^2} + \frac{C_\epsilon \delta_C^{2(r+1)}}{(1 - \delta_A)^2 (1 - \delta_C^2)},$$

whenever  $s$  is large enough such that  $\delta_C (1 + \delta_A^s) < 1$ .

(c) Define  $\Sigma_1^{r,s} = \tilde{C}_s \Sigma_0^{r,s}$  with  $\tilde{C}_s$  and  $\Sigma_0^{r,s}$  as in Theorems 1(a) and 1(b), respectively. The matrix  $\Sigma_1^{r,s}$  is a banded matrix with bandwidth no larger than  $(2rs + 2r + 3s + 1)(k_0 - 1) + 2l_0 + 1$ . Moreover,

$$\left\| \Sigma_1^{r,s} - \Sigma_1 \right\|_{\text{F}} \leq \frac{9C_\epsilon \delta_A^s}{(1 - \delta_A)^2 (1 - \delta_C^2 (1 + \delta_A^s))^2} + \frac{C_\epsilon \delta_C^{2(r+1)}}{(1 - \delta_A)^2 (1 - \delta_C^2)},$$

whenever  $s$  is large enough such that  $\delta_C (1 + \delta_A^s) < 1$ .

(d) Take any  $h_1 \geq 2(rs + r + s)(k_0 - 1) + 2l_0 + 1$ , then

$$\left\| \mathcal{B}_{h_1}(\hat{\Sigma}_0) - \Sigma_0 \right\|_2 \leq \epsilon + \frac{16C_\epsilon \delta_A^s}{(1 - \delta_A)^2 (1 - \delta_C^2 (1 + \delta_A^s)^2)^2} + \frac{2C_\epsilon \delta_C^{2(r+1)}}{(1 - \delta_A)^2 (1 - \delta_C^2)},$$

with a probability of at least  $1 - \mathcal{P}_1(\epsilon, N, T)$  (for polynomial tail decay) or  $1 - \mathcal{P}_2(\epsilon, N, T)$  (for exponential tail decay).

(e) Take any  $h_2 \geq (2rs + 2r + 3s + 1)(k_0 - 1) + 2l_0 + 1$ , then

$$\left\| \mathcal{B}_{h_2}(\hat{\Sigma}_1) - \Sigma_1 \right\|_2 \leq \epsilon + \frac{18C_\epsilon \delta_A^s}{(1 - \delta_A)^2 (1 - \delta_C^2 (1 + \delta_A^s)^2)^2} + \frac{2C_\epsilon \delta_C^{2(r+1)}}{(1 - \delta_A)^2 (1 - \delta_C^2)},$$

with a probability of at least  $1 - \mathcal{P}_1(\epsilon, N, T)$  (for polynomial tail decay) or  $1 - \mathcal{P}_2(\epsilon, N, T)$  (for exponential tail decay).

Explicit expressions for  $C_1, C_2, C_3$ , and  $0 \leq \delta_C < 1$  are provided in the proofs below.

(a) The proof builds upon results from Guo et al. (2016) on banded vector autoregressions. Recall that  $C = DB$  with  $D = (I_N - A)^{-1}$ .  $\tilde{D}_s = \sum_{j=0}^{s-1} A^j$  has a bandwidth of at most  $s(k_0 - 1) + 1$  and satisfies<sup>7</sup>

$$\left\| \tilde{D}_s - D \right\|_{\perp} = \left\| \sum_{j=0}^{s-1} A^j - (I_N - A)^{-1} \right\|_{\perp} = \left\| - \sum_{j=s}^{\infty} A^j \right\|_{\perp} \leq \sum_{j=s}^{\infty} \|A\|_{\perp}^j \leq \frac{\delta_A^s}{1 - \delta_A}.$$

The product  $\tilde{C}_s = \tilde{D}_s B$  has a maximal bandwidth of  $(s + 1)(k_0 - 1) + 1$ . Since  $\|B\|_{\perp} \leq \delta_B$  and  $\frac{\delta_B}{1 - \delta_A} =: \delta_C < 1$  (Assumption 3(b)), we also have

$$\left\| \tilde{C}_s - C \right\|_{\perp} = \left\| (\tilde{D}_s - D) B \right\|_{\perp} \leq \left\| \tilde{D}_s - D \right\|_{\perp} \|B\|_{\perp} \leq \delta_C \delta_A^s.$$

(b) Iterating on the observation in footnote 7, we conclude that the bandwidth of  $\tilde{C}_s^r$  is at most  $r[(s + 1)(k_0 - 1) + 1] - (r - 1) = r(s + 1)(k_0 - 1) + 1$ . The bandwidth of  $\Sigma_0^{r,s}$  therefore does not exceed

$$2[r(s + 1)(k_0 - 1) + 1] + 2(s(k_0 - 1) + 1) + (2l_0 + 1) - 4 = 2(rs + r + s)(k_0 - 1) + 2l_0 + 1.$$

We now bound  $\left\| \Sigma_0^{r,s} - \Sigma_0 \right\|_1$ . Because  $\Sigma_0 = \sum_{j=0}^{\infty} C^j D \Sigma_\epsilon D'(C')^j$ , it holds that

$$\left\| \Sigma_0^{r,s} - \Sigma_0 \right\|_{\perp} \leq \left\| \sum_{j=0}^r [\tilde{C}_s^j \tilde{D}_s \Sigma_\epsilon \tilde{D}_s' (\tilde{C}_s')^j - C^j D \Sigma_\epsilon D'(C')^j] \right\|_{\perp} + \left\| \sum_{j=r+1}^{\infty} C^j D \Sigma_\epsilon D'(C')^j \right\|_{\perp} \tag{B.1}$$

Decomposing the first RHS term in (B.1) as

$$\begin{aligned} & \sum_{j=0}^r [\tilde{C}_s^j \tilde{D}_s \Sigma_\epsilon \tilde{D}_s' (\tilde{C}_s')^j - C^j D \Sigma_\epsilon D'(C')^j] \\ &= \sum_{j=0}^r [(\tilde{C}_s^j \tilde{D}_s - C^j D) \Sigma_\epsilon \tilde{D}_s' (\tilde{C}_s')^j - C^j D \Sigma_\epsilon (D'(C')^j - \tilde{D}_s' (\tilde{C}_s')^j)] \\ &= \sum_{j=0}^r (\tilde{C}_s^j \tilde{D}_s - C^j D) \Sigma_\epsilon (\tilde{D}_s' (\tilde{C}_s')^j - D'(C')^j) \\ & \quad + \sum_{j=0}^r [(\tilde{C}_s^j \tilde{D}_s - C^j D) \Sigma_\epsilon D'(C')^j - C^j D \Sigma_\epsilon (D'(C')^j - \tilde{D}_s' (\tilde{C}_s')^j)], \end{aligned}$$

it follows that

$$\begin{aligned} & \left\| \sum_{j=0}^r [\tilde{C}_s^j \tilde{D}_s \Sigma_\epsilon \tilde{D}_s' (\tilde{C}_s')^j - C^j D \Sigma_\epsilon D'(C')^j] \right\|_{\perp} \\ & \leq C_\epsilon \sum_{j=0}^r \left\| \tilde{C}_s^j \tilde{D}_s - C^j D \right\|_{\perp}^2 + 2C_\epsilon \sum_{j=0}^r \left\| \tilde{C}_s^j \tilde{D}_s - C^j D \right\|_{\perp} \|C^j D\|_{\perp}. \end{aligned} \tag{B.2}$$

Furthermore, noting that

$$\left\| \tilde{C}_s^j \tilde{D}_s - C^j D \right\|_{\perp} \leq \left\| \tilde{C}_s^j - C^j \right\|_{\perp} \left\| \tilde{D}_s \right\|_{\perp} + \|C^j\|_{\perp} \left\| D - \tilde{D}_s \right\|_{\perp}, \tag{B.3}$$

<sup>7</sup> If matrices  $F_1$  and  $F_2$  are banded matrices with bandwidths  $k_1$  and  $k_2$ , respectively, then the product  $F_1 F_2$  is again a banded matrix with a bandwidth of at most  $k_1 + k_2 - 1$ .



we proceed by bounding each norm in (B.3). First, expanding the matrix powers provides

$$\begin{aligned} \|\tilde{C}_s^j - C^j\|_{\perp} &= \left\| (\tilde{C}_s - C + C)^j - C^j \right\|_{\perp} \\ &= \left\| (\tilde{C}_s - C)^j + (\tilde{C}_s - C)^{j-1}C + (\tilde{C}_s - C)^{j-2}C(\tilde{C}_s - C) + \dots + C(\tilde{C}_s - C)^{j-1} \dots + C^j \right\|_{\perp} \\ &\leq \sum_{k=1}^j \binom{j}{k} \|C\|_{\perp}^{j-k} \|\tilde{C}_s - C\|_{\perp}^k = \sum_{k=0}^{j-1} \binom{j}{k+1} \|C\|_{\perp}^{(j-1)-k} \|\tilde{C}_s - C\|_{\perp}^{k+1} \\ &= \|\tilde{C}_s - C\|_{\perp} \sum_{k=0}^{j-1} \frac{j}{k+1} \binom{j-1}{k} \|C\|_{\perp}^{(j-1)-k} \|\tilde{C}_s - C\|_{\perp}^k \leq \delta_C \delta_A^s j \left[ \|\tilde{C}_s - C\|_{\perp} + \|C\|_{\perp} \right]^{j-1} \\ &\leq \delta_C \delta_A^s j \left[ \delta_C (1 + \delta_A^s) \right]^{j-1} = \delta_C^j j \delta_A^s (1 + \delta_A^s)^{j-1}. \end{aligned}$$

Next, it holds that  $\|D - \tilde{D}_s\|_{\perp} = \left\| \sum_{j=s}^{\infty} A^j \right\|_{\perp} \leq \frac{\delta_A^s}{1 - \delta_A}$ ,  $\|\tilde{D}_s\|_{\perp} \leq \frac{1}{1 - \delta_A}$  and  $\|C^j\|_{\perp} \leq \delta_C^j$ , such that we may bound (B.3) as

$$\|\tilde{C}_s^j \tilde{D}_s - C^j D\|_{\perp} \leq \frac{\delta_A^s j \delta_C^j (1 + \delta_A^s)^j}{(1 + \delta_A^s)(1 - \delta_A)} + \frac{\delta_A^s \delta_C^j}{1 - \delta_A} \leq \frac{2\delta_A^s j \delta_C^j (1 + \delta_A^s)^j}{(1 - \delta_A)}. \tag{B.4}$$

Plugging (B.4) into (B.2), we obtain

$$\begin{aligned} &\left\| \sum_{j=0}^r \left[ \tilde{C}_s^j \tilde{D}_s \Sigma_{\epsilon} \tilde{D}_s' (\tilde{C}_s^j)' - C^j D \Sigma_{\epsilon} D' (C^j)' \right] \right\|_{\perp} \\ &\leq \frac{4C_{\epsilon} \delta_A^{2s}}{(1 - \delta_A)^2} \sum_{j=0}^r j^2 \delta_C^{2j} (1 + \delta_A^s)^{2j} + \frac{4C_{\epsilon} \delta_A^s}{(1 - \delta_A)^2} \sum_{j=0}^r j \delta_C^{2j} (1 + \delta_A^s)^j \\ &\leq \frac{4C_{\epsilon} \delta_A^{2s} \delta_C^2 (1 + \delta_A^s)^2}{(1 - \delta_A)^2 (1 - \delta_C^2 (1 + \delta_A^s)^2)^2} + \frac{4C_{\epsilon} \delta_A^s \delta_C^2 (1 + \delta_A^s)}{(1 - \delta_A)^2 (1 - \delta_C^2 (1 + \delta_A^s)^2)^2} \leq \frac{8C_{\epsilon} \delta_A^s}{(1 - \delta_A)^2 (1 - \delta_C^2 (1 + \delta_A^s)^2)^2}, \end{aligned} \tag{B.5}$$

where we have assumed that  $s$  is sufficiently large, such that  $\delta_C(1 + \delta_A^s) < 1$ , and applied standard results for geometric series, i.e.  $\sum_{j=1}^{\infty} jz^j = \frac{z}{(1-z)^2}$  and  $\sum_{j=1}^{\infty} j^2 z^j = \frac{z(1+z)}{(1-z)^3}$  for  $|z| < 1$ . Next, we move to the second RHS term of (B.1), which is bounded by

$$\left\| \sum_{j=r+1}^{\infty} C^j D \Sigma_{\epsilon} D' (C^j)' \right\|_{\perp} \leq \|\Sigma_{\epsilon}\|_{\perp} \|D\|_{\perp}^2 \sum_{j=r+1}^{\infty} \|C\|_{\perp}^{2j} \leq \frac{C_{\epsilon}}{(1 - \delta_A)^2} \sum_{j=r+1}^{\infty} \delta_C^{2j} = \frac{C_{\epsilon} \delta_C^{2(r+1)}}{(1 - \delta_A)^2 (1 - \delta_C^2)}. \tag{B.6}$$

Plugging both (B.5) and (B.6) into (B.1), results in the final bound

$$\|\Sigma_0^{r,s} - \Sigma_0\|_{\perp} \leq \frac{8C_{\epsilon} \delta_A^s}{(1 - \delta_A)^2 (1 - \delta_C^2 (1 + \delta_A^s)^2)^2} + \frac{C_{\epsilon} \delta_C^{2(r+1)}}{(1 - \delta_A)^2 (1 - \delta_C^2)}$$

thereby proving the intermediate result in part (b).

(c) We have  $\Sigma_1 = (I_N - A)^{-1} B \Sigma_0 = C \Sigma_0$ , and hence

$$\begin{aligned} \|\Sigma_1^{r,s} - \Sigma_1\|_{\perp} &= \|\tilde{C}_s (\Sigma_0^{r,s} - \Sigma_0) + (\tilde{C}_s - C) \Sigma_0\|_{\perp} \leq \|\tilde{C}_s\|_{\perp} \|\Sigma_0^{r,s} - \Sigma_0\|_{\perp} + \|\tilde{C}_s - C\|_{\perp} \|\Sigma_0\|_{\perp} \\ &\leq \frac{8C_{\epsilon} \delta_A^s}{(1 - \delta_A)^2 (1 - \delta_C^2 (1 + \delta_A^s)^2)^2} + \frac{C_{\epsilon} \delta_C^{2(r+1)}}{(1 - \delta_A)^2 (1 - \delta_C^2)} + \frac{C_{\epsilon} \delta_C \delta_A^s}{(1 - \delta_A^2) (1 - \delta_C^2)} \\ &\leq \frac{9C_{\epsilon} \delta_A^s}{(1 - \delta_A)^2 (1 - \delta_C^2 (1 + \delta_A^s)^2)^2} + \frac{C_{\epsilon} \delta_C^{2(r+1)}}{(1 - \delta_A)^2 (1 - \delta_C^2)} \end{aligned} \tag{B.7}$$

where we have used that  $\|\tilde{C}_s - C\|_{\perp} \leq \delta_C \delta_A^s$  for the intermediate result (a),  $\|\tilde{C}_s\|_{\perp} \leq \delta_C (1 + \delta_A^s) < 1$  since  $s$  is taken sufficiently large, and  $\|\Sigma_0\|_{\perp} \leq \sum_{j=0}^{\infty} \|C^j D \Sigma_{\epsilon} D' (C^j)'\|_{\perp} \leq \frac{C_{\epsilon}}{(1 - \delta_A^2) (1 - \delta_C^2)}$ .

(d) We have

$$\begin{aligned} \|\mathcal{B}_{h_1}(\hat{\Sigma}_0) - \Sigma_0\|_{\perp} &\leq \|\mathcal{B}_{h_1}(\hat{\Sigma}_0) - \mathcal{B}_{h_1}(\Sigma_0)\|_{\perp} + \|\mathcal{B}_{h_1}(\Sigma_0) - \mathcal{B}_{h_1}(\Sigma_0^{r,s})\|_{\perp} + \|\mathcal{B}_{h_1}(\Sigma_0^{r,s}) - \Sigma_0\|_{\perp} \\ &= \|\mathcal{B}_{h_1}(\hat{\Sigma}_0 - \Sigma_0)\|_{\perp} + \|\mathcal{B}_{h_1}(\Sigma_0 - \Sigma_0^{r,s})\|_{\perp} + \|\Sigma_0^{r,s} - \Sigma_0\|_{\perp} \\ &\leq \|\mathcal{B}_{h_1}(\hat{\Sigma}_0 - \Sigma_0)\|_{\perp} + 2 \|\Sigma_0^{r,s} - \Sigma_0\|_{\perp} \end{aligned} \tag{B.8}$$

because  $\Sigma_0^{r,s}$  is a banded matrix already and banding can only decrease the norm difference between  $\Sigma_0^{r,s}$  and  $\Sigma_0$ . We consider the two terms in the RHS of (B.8) separately. Note that  $\|\hat{\Sigma} - \sigma_0\|_{\perp} = \|\hat{\Sigma} - \sigma_0\|_1$ , by symmetry. Then, letting  $\xi_t$  be as defined in Lemma 3,

$$\begin{aligned}
 \mathbb{P}\left(\left\|\mathcal{B}_{h_1}(\widehat{\Sigma}_0 - \Sigma_0)\right\|_{\perp} \leq x\right) &= \mathbb{P}\left(\max_{i \leq j \leq N} \sum_{i=1 \wedge (j-h_1)}^{N \vee (j+h_1)} \left| \sum_{t=2}^T \xi_{it} \xi_{jt} - \mathbb{E}(\xi_{it} \xi_{jt}) \right| \leq Tx\right) \\
 &\geq 1 - \sum_{j=1}^N \mathbb{P}\left(\sum_{i=1 \wedge (j-h_1)}^{N \vee (j+h_1)} \left| \sum_{t=2}^T \xi_{it} \xi_{jt} - \mathbb{E}(\xi_{it} \xi_{jt}) \right| > Tx\right) \\
 &\geq 1 - \sum_{j=1}^N \sum_{i=1 \wedge (j-h_1)}^{N \vee (j+h_1)} \mathbb{P}\left(\left| \sum_{t=2}^T \xi_{it} \xi_{jt} - \mathbb{E}(\xi_{it} \xi_{jt}) \right| > \frac{Tx}{2h_1 + 1}\right) \tag{B.9} \\
 &\geq \begin{cases} 1 - (2h_1 + 1)N \left[ \left( b_1 T^{(1-\delta)/3} + \frac{(2h_1+1)b_3}{x} \right) \exp\left(-\frac{T^{(1-\delta)/3}}{2b_1^2}\right) + \frac{b_2(2h_1+1)^d}{x^d T^{\frac{\delta}{2}(d-1)}} \right] & \text{(polynomial tails),} \\ 1 - (2h_1 + 1)N \left[ \frac{\kappa_1(2h_1+1)}{x} + \frac{2}{\kappa_2} \left( \frac{Tx^2}{(2h_1+1)^2} \right)^{1/7} \right] \exp\left(-\frac{1}{\kappa_3} \left( \frac{Tx^2}{(2h_1+1)^2} \right)^{1/7}\right) & \text{(exponential tails),} \end{cases}
 \end{aligned}$$

where the last inequality exploits Lemma 3. Note that the probabilities in (B.9) coincide with the probabilities defined as  $1 - \mathcal{P}_1(x, N, T)$  and  $1 - \mathcal{P}_2(x, N, T)$  in Theorem 1. Overall, if  $\left\|\mathcal{B}_{h_1}(\widehat{\Sigma}_0 - \Sigma_0)\right\|_{\perp} \leq \epsilon$  holds, then, applying intermediate result (b) to (B.8), we obtain the bound

$$\left\|\mathcal{B}_{h_1}(\widehat{\Sigma}_0) - \Sigma_0\right\|_{\perp} \leq \epsilon + 2 \left\|\Sigma_0^{r,s} - \Sigma_0\right\|_{\perp} \leq \epsilon + \frac{16C_\epsilon \delta_A^s}{(1 - \delta_A)^2 (1 - \delta_C^2 (1 + \delta_A^s))^2} + \frac{2C_\epsilon \delta_C^{2(r+1)}}{(1 - \delta_A)^2 (1 - \delta_C^2)}.$$

(e) Mimicking the steps from (d), we find

$$\left\|\mathcal{B}_{h_2}(\widehat{\Sigma}_1) - \Sigma_1\right\|_{\perp} \leq \epsilon + 2 \left\|\Sigma_1^{r,s} - \Sigma_1\right\|_{\perp} \leq \epsilon + \frac{18C_\epsilon \delta_A^s}{(1 - \delta_A)^2 (1 - \delta_C^2 (1 + \delta_A^s))^2} + \frac{2C_\epsilon \delta_C^{2(r+1)}}{(1 - \delta_A)^2 (1 - \delta_C^2)}$$

if we use part (c) and if  $\left\{\left\|\mathcal{B}_{h_2}(\widehat{\Sigma}_1 - \Sigma_1)\right\|_{\perp} \leq \epsilon\right\}$  holds. The probability of the latter event is  $1 - \mathcal{P}_1(\epsilon, N, T)$  (polynomial tails) or  $1 - \mathcal{P}_2(\epsilon, N, T)$  (exponential tails).

We now combine all these intermediate results to recover the result from the theorem. Parts (d) and (e) are both applicable since

$$\begin{aligned}
 h &= \max\{h_1, h_2\} = h_2 = (2rs + 2r + 3s + 1)(k_0 - 1) + 2l_0 + 1 \\
 &\leq (2r(s + 1) + 3(s + 1))(k_0 - 1) + 2l_0 + 1 \leq (s + 1)(2r + 3)(k_0 - 1) + 2l_0 + 1. \tag{B.10}
 \end{aligned}$$

Then, by intermediate results (d) and (e), it holds that

$$\begin{aligned}
 \left\|\widehat{V}_h - V\right\|_{\perp} &\leq \left\|\mathcal{B}_{h_2}(\widehat{\Sigma}_1) - \Sigma_1\right\|_{\perp} + \left\|\mathcal{B}_{h_1}(\widehat{\Sigma}_0) - \Sigma_0\right\|_{\perp} \\
 &\leq 2\epsilon + \frac{34C_\epsilon \delta_A^s}{(1 - \delta_A)^2 (1 - \delta_C^2 (1 + \delta_A^s))^2} + \frac{4C_\epsilon \delta_C^{2(r+1)}}{(1 - \delta_A)^2 (1 - \delta_C^2)}
 \end{aligned}$$

with probabilities  $1 - 2\mathcal{P}_1(\epsilon, N, T)$  and  $1 - 2\mathcal{P}_2(\epsilon, N, T)$  in the cases of polynomial and exponential tail decay, respectively.

It remains to determine  $s$  and  $r$  such that  $\delta_C(1 + \delta_A^s) < 1$ ,

$$\frac{34C_\epsilon \delta_A^s}{(1 - \delta_A)^2 (1 - \delta_C^2 (1 + \delta_A^s))^2} = \epsilon, \quad \text{and} \quad \frac{4C_\epsilon \delta_C^{2(r+1)}}{(1 - \delta_A)^2 (1 - \delta_C^2)} = \epsilon.$$

First, note that we can ensure  $\delta_C(1 + \delta_A^s) = K < 1$ , by choosing  $s > \log\left(\frac{\delta_C}{K - \delta_C}\right) / |\log(\delta_A)|$ . Then,

$$\frac{34C_\epsilon \delta_A^s}{(1 - \delta_A)^2 (1 - \delta_C^2 (1 + \delta_A^s))^2} = \frac{34C_\epsilon \delta_A^s}{(1 - \delta_A)^2 (1 - K^2)^2} = \epsilon \Leftrightarrow s = \frac{\log\left(\frac{34C_\epsilon}{(1 - \delta_A)^2 (1 - K^2)^2 \epsilon}\right)}{|\log(\delta_A)|}. \tag{B.11}$$

Furthermore,

$$\frac{4C_\epsilon \delta_C^{2(r+1)}}{(1 - \delta_A)^2 (1 - \delta_C^2)} = \epsilon \Leftrightarrow r = \frac{\log\left(\frac{4C_\epsilon}{(1 - \delta_A)^2 (1 - \delta_C^2) \epsilon}\right)}{2|\log(\delta_C)|} - \frac{1}{2}. \tag{B.12}$$

Define  $C_1 = \frac{34C_\epsilon}{(1 - \delta_A)^2 (1 - K^2)^2}$  and  $C_2 = \frac{4C_\epsilon}{(1 - \delta_A)^2 (1 - \delta_C^2)}$ . Plugging (B.11)–(B.12) into (B.10), we obtain the result that when

$$h = \left( \frac{\max\left\{\log\left(\frac{K - \delta_C}{\delta_C}\right), \log\left(\frac{C_1}{\epsilon}\right)\right\}}{|\log(\delta_A)|} + 1 \right) \left( \frac{\log\left(\frac{C_2}{\epsilon}\right)}{|\log(\delta_C)|} + 2 \right) (k_0 - 1) + 2l_0 + 1,$$

for some  $K \in (\delta_c, 1)$ , it holds that  $\|\hat{V}_h - V\|_- \leq 4\epsilon$  with probability at least  $1 - 2\mathcal{P}_1(\epsilon, N, T)$  and  $1 - 2\mathcal{P}_2(\epsilon, N, T)$  in the cases of polynomial and exponential tail decay, respectively.  $\square$

**Proof of Theorem 2.** The proof of the theorem relies on the properties of dual norms. Recall  $P_\alpha(c) = \alpha \sum_{g \in \mathcal{G}} \sqrt{|g|} \|c_g\|_2 + (1-\alpha) \|c\|_1$ . Exploiting the properties of  $\|\cdot\|_1$  and  $\|\cdot\|_2$ , it is straightforward to verify that  $P_\alpha(\cdot)$  is a norm for any  $0 \leq \alpha \leq 1$ . For any norm  $\|\cdot\|$ , we define its dual norm  $\|\cdot\|^*$  through  $\|c\|^* = \sup_{x \neq 0} \frac{|c'x|}{\|x\|}$ . The dual-norm inequality states that

$$c'x \leq \|c\|^* \|x\| \quad \text{for all conformable vectors } c \text{ and } x. \tag{B.13}$$

For the norm  $P_\alpha(c)$ , its dual norm  $P_\alpha^*(c)$  is bounded by

$$\begin{aligned} P_\alpha^*(c) &= \sup_{x \neq 0} \frac{|c'x|}{P_\alpha(x)} = \sup_{x \neq 0} \frac{|c'x|}{\alpha \sum_{g \in \mathcal{G}} \sqrt{|g|} \|x_g\|_2 + (1-\alpha) \|x\|_1} \\ &\stackrel{(i)}{\leq} \alpha \sup_{x \neq 0} \frac{|c'x|}{\sum_{g \in \mathcal{G}} \sqrt{|g|} \|x_g\|_2} + (1-\alpha) \sup_{x \neq 0} \frac{|c'x|}{\|x\|_1} \stackrel{(ii)}{\leq} \alpha \max_{g \in \mathcal{G}} \frac{\|c_g\|_2}{\sqrt{|g|}} + (1-\alpha) \|c\|_\infty, \end{aligned} \tag{B.14}$$

by convexity of the function  $f(x) = x^{-1}$  in step (i), and using for step (ii) both  $\|c\|_1^* = \|c\|_\infty$  and

$$\begin{aligned} \sup_{x \neq 0} \frac{|c'x|}{\sum_{g \in \mathcal{G}} \sqrt{|g|} \|x_g\|_2} &= \sup_{x \neq 0, \sum_{g \in \mathcal{G}} \sqrt{|g|} \|x_g\|_2 = 1} \left| \sum_{g \in \mathcal{G}} c'_g x_g \right| \\ &\leq \sup_{x \neq 0, \sum_{g \in \mathcal{G}} \sqrt{|g|} \|x_g\|_2 = 1} \sum_{g \in \mathcal{G}} \left\| \frac{c_g}{\sqrt{|g|}} \right\|_2 \|\sqrt{|g|} x_g\|_2 \leq \max_{g \in \mathcal{G}} \frac{\|c_g\|_2}{\sqrt{|g|}}. \end{aligned}$$

We now start the actual proof. Recall  $\hat{\sigma}_h = \text{vec}(\mathcal{B}_h(\hat{\Sigma}_1)')$ ,  $\hat{V}_h = [\mathcal{B}_h(\hat{\Sigma}_1)' \quad \mathcal{B}_h(\hat{\Sigma}_0)]$ ,  $V = [\Sigma_1' \quad \Sigma_0]$ , and  $C = [A \quad B]$ . Exploiting standard properties of  $\text{vec}(\cdot)$ , we find

$$\begin{aligned} \hat{\sigma}_h - \hat{V}_h^{(d)} c &= \text{vec}(\mathcal{B}_h(\hat{\Sigma}_1)' - \hat{V}_h C') = \text{vec} \left( \left[ \mathcal{B}_h(\hat{\Sigma}_1) - \Sigma_1 \right]' - [\hat{V}_h - V] C' + \underbrace{[\Sigma_1' - V C']}_{=0, \text{ see (5)}} \right) \\ &= \underbrace{\text{vec}(\mathcal{B}_h(\hat{\Sigma}_1)' - \Sigma_1')}_{\hat{\Delta}_\Sigma} - \underbrace{[\hat{V}_h^{(d)} - V^{(d)}]}_{\hat{\Delta}_V} c = \hat{\Delta}_\Sigma - \hat{\Delta}_V c. \end{aligned} \tag{B.15}$$

Using (B.15), we rewrite

$$\begin{aligned} \|\hat{\sigma}_h - \hat{V}_h^{(d)} \hat{c}\|_2^2 &= \left\| \left[ \hat{\sigma}_h - \hat{V}_h^{(d)} c \right] - \left[ \hat{V}_h^{(d)} (\hat{c} - c) \right] \right\|_2^2 \\ &= \|\hat{\sigma}_h - \hat{V}_h^{(d)} c\|_2^2 + \|\hat{V}_h^{(d)} (\hat{c} - c)\|_2^2 - 2(\hat{c} - c)' \hat{V}_h^{(d)'} (\hat{\Delta}_\Sigma - \hat{\Delta}_V c). \end{aligned}$$

Recalling the objective function  $\mathcal{L}_\alpha(c; \lambda) = \|\hat{\sigma}_h - \hat{V}_h^{(d)} c\|_2^2 + \lambda P_\alpha(c)$  and noting that  $\mathcal{L}_\alpha(\hat{c}; \lambda) \leq \mathcal{L}_\alpha(c; \lambda)$  by construction, it follows that

$$\begin{aligned} \|\hat{V}_h^{(d)} (\hat{c} - c)\|_2^2 + \lambda P_\alpha(\hat{c}) &\leq 2(\hat{c} - c)' \hat{V}_h^{(d)'} (\hat{\Delta}_\Sigma - \hat{\Delta}_V c) + \lambda P_\alpha(c) \\ &\leq 2P_\alpha(\hat{c} - c) P_\alpha^* \left( \hat{V}_h^{(d)'} (\hat{\Delta}_\Sigma - \hat{\Delta}_V c) \right) + \lambda P_\alpha(c) \\ &\leq P_\alpha(\hat{c} - c) \left[ 2P_\alpha^* \left( \hat{V}_h^{(d)'} \hat{\Delta}_\Sigma \right) + 2P_\alpha^* \left( \hat{V}_h^{(d)'} \hat{\Delta}_V c \right) \right] + \lambda P_\alpha(c), \end{aligned} \tag{B.16}$$

where we used the dual-norm inequality (see (B.13)) and the triangle property of (dual) norms in the second and third inequality, respectively. Define the sets

$$\mathcal{H}_1(x) = \left\{ 2P_\alpha^* \left( \hat{V}_h^{(d)'} \hat{\Delta}_\Sigma \right) \leq x \right\} \quad \text{and} \quad \mathcal{H}_2(x) = \left\{ 2P_\alpha^* \left( \hat{V}_h^{(d)'} \hat{\Delta}_V c \right) \leq x \right\}. \tag{B.17}$$

On the set  $\mathcal{H}_1(\frac{\lambda}{4}) \cap \mathcal{H}_2(\frac{\lambda}{4})$ , we can scale (B.16) by a factor 2 to obtain

$$2 \|\hat{V}_h^{(d)} (\hat{c} - c)\|_2^2 + 2\lambda P_\alpha(\hat{c}) \leq \lambda P_\alpha(\hat{c} - c) + 2\lambda P_\alpha(c). \tag{B.18}$$

We subsequently manipulate  $P_\alpha(\hat{c})$  and  $P_\alpha(\hat{c} - c)$ . Using the reverse triangle inequality, we have

$$\begin{aligned} P_\alpha(\hat{c}) &= \alpha \sum_{g \in \mathcal{G}} \sqrt{|g|} \|\hat{c}_g\|_2 + (1 - \alpha) \|\hat{c}\|_1 \\ &\geq \alpha \sum_{g \in \mathcal{G}_S} \sqrt{|g|} \left[ \|c_g\|_2 - \|\hat{c}_g - c_g\|_2 \right] + \alpha \sum_{g \in \mathcal{G}_{S^c}} \sqrt{|g|} \|\hat{c}_g\|_2 + (1 - \alpha) \left[ \|c_S\|_1 - \|\hat{c}_S - c_S\|_1 + \|\hat{c}_{S^c}\|_1 \right] \\ &= P_{\alpha,S}(c) + P_{\alpha,S^c}(\hat{c} - c) - P_{\alpha,S}(\hat{c} - c), \end{aligned} \tag{B.19}$$

where  $\mathcal{G}_S$  and  $\mathcal{G}_{S^c}$  are defined in Lemma 1. Simple rewriting provides

$$\begin{aligned} P_\alpha(\hat{c} - c) &= \alpha \sum_{g \in \mathcal{G}} \sqrt{|g|} \|\hat{c}_g - c_g\|_2 + (1 - \alpha) \|\hat{c}_S - c_S\|_1 + \alpha \sum_{g \in \mathcal{G}_{S^c}} \sqrt{|g|} \|\hat{c}_g\|_2 + (1 - \alpha) \|\hat{c}_{S^c}\|_1 \\ &= P_{\alpha,S}(\hat{c} - c) + P_{\alpha,S^c}(\hat{c} - c). \end{aligned} \tag{B.20}$$

Combining results (B.18)–(B.20) yields

$$2 \|\hat{V}_h^{(d)}(\hat{c} - c)\|_2^2 + \lambda P_{\alpha,S^c}(\hat{c} - c) \leq 3\lambda P_{\alpha,S}(\hat{c} - c), \tag{B.21}$$

and  $\hat{c} - c$  is thus a member of the set  $\mathcal{C}_{N_c}(\mathcal{G}, S)$  as defined in Lemma 1.

In combination with Lemma 2, thus requiring  $\mathcal{H}_1(\frac{\lambda}{4}) \cap \mathcal{H}_2(\frac{\lambda}{4}) \cap \mathcal{V}(\frac{\phi_0}{2})$  to hold, we conclude

$$\begin{aligned} 2 \|\hat{V}_h^{(d)}(\hat{c} - c)\|_2^2 + \lambda P_\alpha(\hat{c} - c) &= 2 \|\hat{V}_h^{(d)}(\hat{c} - c)\|_2^2 + \lambda P_{\alpha,S}(\hat{c} - c) + \lambda P_{\alpha,S^c}(\hat{c} - c) \\ &\stackrel{(i)}{\leq} 4\lambda P_{\alpha,S}(\hat{c} - c) \stackrel{(ii)}{\leq} 16 \|\hat{V}_h^{(d)}(\hat{c} - c)\|_2 \left( \frac{\bar{\omega}_\alpha \lambda}{\phi_0} \right) \stackrel{(iii)}{\leq} \|\hat{V}_h^{(d)}(\hat{c} - c)\|_2^2 + \frac{64\bar{\omega}_\alpha^2 \lambda^2}{\phi_0^2}, \end{aligned} \tag{B.22}$$

where step (i) follows from (B.21), step (ii) is implied by  $P_{\alpha,S}(\hat{c} - c) \leq \frac{4\bar{\omega}_0 \|\hat{V}_h^{(d)}(\hat{c} - c)\|_2}{\phi_0}$  for  $\hat{c} - c \in \mathcal{C}_{N_c}(\mathcal{G}, S)$  (Lemma 2), and step (iii) uses the elementary inequality  $16uv \leq u^2 + 64v^2$  (i.e. manipulating  $(u - 8v)^2 \geq 0$ ). A straightforward rearrangement of (B.22) provides the inequality of Theorem 2.

It remains to determine a lower bound on the probability of  $\mathcal{H}_1(\frac{\lambda}{4}) \cap \mathcal{H}_2(\frac{\lambda}{4}) \cap \mathcal{V}(\frac{\phi_0}{2})$ . We rely on the elementary inequality  $\mathbb{P}\left(\mathcal{H}_1(\frac{\lambda}{4}) \cap \mathcal{H}_2(\frac{\lambda}{4}) \cap \mathcal{V}(\frac{\phi_0}{2})\right) \geq 1 - \mathbb{P}\left(\mathcal{H}_1(\frac{\lambda}{4})^c\right) - \mathbb{P}\left(\mathcal{H}_2(\frac{\lambda}{4})^c\right) - \mathbb{P}\left(\mathcal{V}(\frac{\phi_0}{2})^c\right)$  to bound the individual probabilities.

We start with  $\mathbb{P}\left(\mathcal{H}_1(\frac{\lambda}{4})^c\right) = \mathbb{P}\left(2P_\alpha^*(\hat{V}_h^{(d)'} \hat{\Delta}_\Sigma) > \frac{\lambda}{4}\right) \leq \mathbb{P}\left(\|\hat{V}_h^{(d)'} \hat{\Delta}_\Sigma\|_\infty > \frac{\lambda}{8}\right)$ . The last inequality is true because continuing from (B.14), we have

$$P_\alpha^*(c) \leq \alpha \max_{g \in \mathcal{G}} \frac{\|c_g\|_2}{\sqrt{|g|}} + (1 - \alpha) \|c\|_\infty \leq \alpha \max_{g \in \mathcal{G}} \frac{\sqrt{|g|} \|c_g\|_\infty}{\sqrt{|g|}} + (1 - \alpha) \|c\|_\infty = \|c\|_\infty \tag{B.23}$$

for any vector  $c$ . Subsequently, we have

$$\begin{aligned} \mathbb{P}\left(\|\hat{V}_h^{(d)'} \hat{\Delta}_\Sigma\|_\infty > \frac{\lambda}{8}\right) &= \mathbb{P}\left(\left\| \left[ (\hat{V}_h^{(d)} - V^{(d)}) + V^{(d)} \right]' \hat{\Delta}_\Sigma \right\|_\infty > \frac{\lambda}{8}\right) \\ &\leq \mathbb{P}\left(\|\hat{V}_h - V\|_1 \|\hat{\Delta}_\Sigma\|_\infty + C_V \|\hat{\Delta}_\Sigma\|_\infty > \frac{\lambda}{8}\right) \leq \mathbb{P}\left(\|\hat{V}_h - V\|_1^2 + C_V \|\hat{V}_h - V\|_1 > \frac{\lambda}{8}\right) \\ &\leq \mathbb{P}\left(\|\hat{V}_h - V\|_1 > \frac{\lambda^{1/2}}{4}\right) + \mathbb{P}\left(\|\hat{V}_h - V\|_1 > \frac{\lambda}{16C_V}\right), \end{aligned} \tag{B.24}$$

exploiting block-diagonality of  $\hat{V}_h^{(d)} - V^{(d)}$  and  $V^{(d)}$  such that  $\|\hat{V}_h^{(d)} - V^{(d)}\|_1 = \max_{1 \leq i \leq N} \|\hat{V}_{i,h} - V_i\|_1 \leq \|\hat{V}_h - V\|_1$  and  $\|V^{(d)}\|_1 \leq \|V\|_1 \leq C_V$  (explicitly assumed in Theorem 2). Bounds for the final RHS terms in (B.24) are available from Theorem 1.

Second, we have

$$\begin{aligned} \mathbb{P}\left(\mathcal{H}_2(\frac{\lambda}{4})^c\right) &\leq \mathbb{P}\left(\|\hat{V}_h^{(d)'} \hat{\Delta}_V c\|_\infty > \frac{\lambda}{8}\right) \leq \mathbb{P}\left(\left\| \hat{V}_h^{(d)'} \left[ \hat{V}_h^{(d)} - V^{(d)} \right] \right\|_\infty > \frac{\lambda}{8 \|c\|_\infty}\right) \\ &\leq \mathbb{P}\left(\|\hat{V}_h - V\|_1 \|\hat{V}_h - V\|_\infty + C_V \|\hat{V}_h - V\|_\infty > \frac{\lambda}{8 \|c\|_\infty}\right) \\ &\leq \mathbb{P}\left(\|\hat{V}_h - V\|_1 > \frac{\lambda^{1/2}}{4}\right) + \mathbb{P}\left(\|\hat{V}_h - V\|_1 > \frac{\lambda}{16C_V}\right), \end{aligned} \tag{B.25}$$

where the last line relies on the union bound and the fact that  $\|c\|_\infty < 1$  (implied by Assumption 1). Hence, the sets  $\mathcal{H}_1(\frac{\lambda}{4})^c$  and  $\mathcal{H}_2(\frac{\lambda}{4})^c$  admit the same probability bound.

Finally,  $\mathbb{P}\left(\mathcal{V}\left(\frac{\phi_0}{2}\right)^c\right) = \mathbb{P}\left(\|\hat{\mathbf{V}}_h - \mathbf{V}\|_1 > \frac{\phi_0}{2}\right) \leq \mathbb{P}\left(\|\hat{\mathbf{V}}_h - \mathbf{V}\|_r > \frac{\phi_0}{2}\right)$ . Combining all previous results, we conclude

$$\begin{aligned} &\mathbb{P}\left(\mathcal{A}_1\left(\frac{\lambda}{4}\right) \cap \mathcal{A}_2\left(\frac{\lambda}{4}\right) \cap \mathcal{V}\left(\frac{\phi_0}{2}\right)\right) \\ &\geq 1 - 2\mathbb{P}\left(\|\hat{\mathbf{V}}_h - \mathbf{V}\|_r > \frac{\lambda^{1/2}}{4}\right) - 2\mathbb{P}\left(\|\hat{\mathbf{V}}_h - \mathbf{V}\|_r > \frac{\lambda}{16C_V}\right) - \mathbb{P}\left(\|\hat{\mathbf{V}}_h - \mathbf{V}\|_r > \frac{\phi_0}{2}\right) \\ &\geq 1 - 5\mathbb{P}\left(\|\hat{\mathbf{V}}_h - \mathbf{V}\|_r > 6f(\lambda, \phi_0)\right), \end{aligned} \tag{B.26}$$

where  $f(\lambda, \phi_0) = \min\left(\frac{\lambda^{1/2}}{24}, \frac{\lambda}{96C_V}, \frac{\phi_0}{12}\right)$ . The proof is completed by evaluating the final probability in (B.26) using Theorem 1.  $\square$

**Proof of Corollary 1.** First, we derive the conditions under which the set on which the performance bound in Theorem 2 holds occurs with probability converging to one. Under Assumption 2(b1), along with the remaining assumptions in Theorem 2, this probability is given by  $1 - \mathcal{P}_1(f(\lambda, \phi_0), N, T)$ , where we recall from Theorem 1 that

$$\begin{aligned} \mathcal{P}_1(f(\lambda, \phi_0), N, T) &= (2h(f(\lambda, \phi_0)) + 1)N \left( b_1 T^{(1-\delta)/3} + \frac{[2h(f(\lambda, \phi_0)) + 1]b_3}{f(\lambda, \phi_0)} \right) \exp\left(-\frac{T^{(1-\delta)/3}}{2b_1^2}\right) \\ &\quad + (2h(f(\lambda, \phi_0)) + 1)N \frac{b_2 [2h(f(\lambda, \phi_0)) + 1]^d}{f(\lambda, \phi_0)^d T^{\frac{\delta}{2}(d-1)}}, \end{aligned} \tag{B.27}$$

for some  $0 < \delta < 1$  and  $f(\lambda, \phi_0) = \min\left(\frac{\lambda^{1/2}}{24}, \frac{\lambda}{96C_V}, \frac{\phi_0}{12}\right)$ . Given that  $\lambda = O(T^{-q_\lambda})$  with  $q_\lambda > 0$ , it follows immediately that

$$f(\lambda, \phi_0) = \min\left(\frac{\lambda^{1/2}}{24}, \frac{\lambda}{96C_V}, \frac{\phi_0}{12}\right) = O(T^{-q_\lambda}) \tag{B.28}$$

and, recalling the definition of  $h(f(\lambda, \phi_0))$  in (4),

$$h(f(\lambda, \phi_0)) = O(\log(T)^2 T^{q_k}). \tag{B.29}$$

Since these are polynomial rates in  $T$ , it follows that the first RHS-term in (B.27) converges to zero exponentially in  $T$  for any  $\delta < 1$ . The second RHS-term, however, converges to zero at most at a polynomial rate, such that

$$\begin{aligned} \mathcal{P}_1(f(\lambda, \phi_0), N, T) &= O\left(\frac{(2h(f(\lambda, \phi_0)) + 1)N[h(f(\lambda, \phi_0)) + 1]^d}{\lambda^d T^{\frac{\delta}{2}(d-1)}}\right) \\ &= O\left(\log(T)^{2(d+1)} T^{q_N + (d+1)q_k + dq_\lambda - \frac{\delta(d-1)}{2}}\right). \end{aligned} \tag{B.30}$$

From (B.30), it follows that  $\mathcal{P}_1(f(\lambda, \phi_0), N, T) \rightarrow 0$  as  $T \rightarrow \infty$  if

$$q_N + (d+1)q_k + dq_\lambda - \frac{\delta(d-1)}{2} < 0 \Rightarrow q_\lambda < \frac{\delta(d-1)}{2d} - \frac{(d+1)q_k}{d} - \frac{q_N}{d}$$

In a similar fashion, we derive the conditions under which  $\mathcal{P}_2(f(\lambda, \phi_0), N, T) \rightarrow 0$ , by noting that

$$\begin{aligned} \mathcal{P}_2(f(\lambda, \phi_0), N, T) &= (2h(f(\lambda, \phi_0)) + 1)N \left[ \frac{\kappa_1 [2h(f(\lambda, \phi_0)) + 1]}{f(\lambda, \phi_0)} + \frac{2}{\kappa_2} \left( \frac{Tf(\lambda, \phi_0)^2}{[2h+1]^2} \right)^{\frac{1}{7}} \right] \\ &\quad \times \exp\left(-\frac{1}{\kappa_3} \left( \frac{Tf(\lambda, \phi_0)^2}{[2h(f(\lambda, \phi_0)) + 1]^2} \right)^{\frac{1}{7}}\right) \end{aligned} \tag{B.31}$$

converges to zero exponentially fast in  $T$  if  $\frac{Tf(\lambda, \phi_0)^2}{[2h(f(\lambda, \phi_0)) + 1]^2}$  diverges at a polynomial rate in  $T$ . Making use of (B.28) and (B.29), it follows that

$$\frac{Tf(\lambda, \phi_0)^2}{[2h(f(\lambda, \phi_0)) + 1]^2} = O(\log(T)^{-4} T^{1-2q_\lambda-2q_k}),$$

which translates to the condition  $1 - 2q_\lambda - 2q_k > 0$ , or  $q_\lambda < \frac{1}{2} - q_k$ . This establishes conditions (i) and (ii) in Corollary 1.

We proceed by deriving the order of the performance bound in Theorem 2. Noting that

$$\sum_{g \in \mathcal{G}_S} \sqrt{|g|} \leq |\mathcal{G}_S| \max_{g \in \mathcal{G}_S} \sqrt{|g|} = O\left(T^{q_g + q_N/2}\right),$$

it follows that

$$\bar{\omega} = O\left((1 - \alpha)T^{q_g + q_N/2} + \alpha T^{q_S/2}\right).$$

Then, by Theorem 2,

$$\left\| \hat{\mathbf{V}}_h^{(d)}(\hat{c} - c) \right\|_2^2 \leq \frac{4\bar{\omega}_\alpha^2 \lambda^2}{\phi_0^2} = O\left((1 - \alpha)T^{2q_g + q_N - 2q_\lambda} + \alpha T^{q_S - 2q_\lambda}\right)$$

and

$$(1 - \alpha) \sum_{g \in \mathcal{G}} \sqrt{|g|} \|\hat{c}_g - c_g\|_2 + \alpha \|\hat{c} - c\|_1 \leq \frac{4\hat{\omega}_\alpha^2 \lambda}{\phi_0^2} = O\left((1 - \alpha)T^{2q_g + q_N - q_\lambda} + \alpha T^{q_S - q_\lambda}\right),$$

on a set with probability  $1 - \mathcal{P}_1(f(\lambda, \phi_0), N, T)$  or  $1 - \mathcal{P}_2(f(\lambda, \phi_0), N, T)$ , depending on whether Assumption 2(b1) or 2(b2) applies, respectively. Since we have shown that both  $\mathcal{P}_1(f(\lambda, \phi_0), N, T) \rightarrow 0$  and  $\mathcal{P}_2(f(\lambda, \phi_0), N, T) \rightarrow 0$  under the conditions imposed in Corollary 1, the proof is complete.  $\square$

**Proof of Theorem 3.** Set  $\beta^* = (\beta'_1, \dots, \beta'_K)'$ ,  $U^{(d)} = [\mathbf{W}_1^{*(d)} \dots \mathbf{W}_K^{*(d)}]$ , and  $\mathcal{Q}^{(d)} = [\mathbf{V}^{*(d)} \ U^{(d)}]$  (and similar quantities with “hats”). The residuals entering the  $L_2$  component of the SPLASHX objective function are

$$\begin{aligned} \hat{\sigma}^* - \hat{\mathbf{V}}^{*(d)} c - \sum_{k=1}^K \hat{\mathbf{W}}_k^{*(d)} \beta_k &= \hat{\sigma}^* - \left[ \hat{\mathbf{V}}^{*(d)} \ \hat{U}^{(d)} \right] \begin{bmatrix} c \\ \beta^* \end{bmatrix} \\ &= \underbrace{\hat{\sigma}^* - [\mathbf{V}^{*(d)} \ U^{(d)}] \mathbf{q}}_{=0} + \underbrace{[\hat{\sigma}^* - \sigma^*]}_{:=\hat{\Delta}_1} - \underbrace{[\hat{\mathbf{V}}^{*(d)} - \mathbf{V}^{*(d)}] c}_{:=\hat{\Delta}_2} - \underbrace{[\hat{U}^{(d)} - U^{(d)}] \beta^*}_{:=\hat{\Delta}_3} \\ &= \hat{\Delta}_1 - \hat{\Delta}_2 c - \hat{\Delta}_3 \beta^*. \end{aligned} \tag{B.32}$$

Exploiting the previously defined notation and (B.32), we have

$$\begin{aligned} \left\| \hat{\sigma}^* - \hat{\mathbf{V}}^{*(d)} c - \sum_{k=1}^K \hat{\mathbf{W}}_k^{*(d)} \beta_k \right\|_2^2 &= \left\| \hat{\sigma}^* - \hat{\mathcal{Q}}^{(d)} \hat{\mathbf{q}} \right\|_2^2 = \left\| \left\{ \hat{\sigma}^* - \hat{\mathcal{Q}}^{(d)} \mathbf{q} \right\} - \left\{ \hat{\mathcal{Q}}^{(d)} (\hat{\mathbf{q}} - \mathbf{q}) \right\} \right\|_2^2 \\ &= \left\| \hat{\sigma}^* - \hat{\mathcal{Q}}^{(d)} \mathbf{q} \right\|_2^2 + \left\| \hat{\mathcal{Q}}^{(d)} (\hat{\mathbf{q}} - \mathbf{q}) \right\|_2^2 - 2(\hat{\mathbf{q}} - \mathbf{q})' \hat{\mathcal{Q}}^{(d)'} (\hat{\Delta}_1 - \hat{\Delta}_2 c - \hat{\Delta}_3 \beta^*). \end{aligned} \tag{B.33}$$

We subsequently adjust the penalty function  $P_\alpha(c)$  to incorporate the penalty on the coefficients in front of the exogenous variables. Define the index set  $g_k$  such that  $\mathbf{q}_{g_k} = \beta_k$  ( $k = 1 \dots, K$ ) and enlarge  $\mathcal{G}$  to  $\mathcal{G}^* = \mathcal{G} \cup \bigcup_{k=1}^K g_k$ . As  $|\beta_k| = |g_k| = N$ , we have

$$P_\alpha(c) + \sum_{k=1}^K (1 - \alpha) \sqrt{N} \|\beta_k\|_2 + \alpha \|\beta_k\|_1 = (1 - \alpha) \sum_{g \in \mathcal{G}^*} \sqrt{|g|} \|\mathbf{q}_g\|_2 + \alpha \|\mathbf{q}\|_1 =: \tilde{P}_\alpha(\mathbf{q}).$$

Using this newly defined norm  $\tilde{P}_\alpha(\mathbf{q})$ , a more concise notation for the SPLASHX objective function is  $\mathcal{L}_\alpha^*(\mathbf{q}; \lambda) = \left\| \hat{\sigma}^* - \hat{\mathcal{Q}}^{(d)} \mathbf{q} \right\|_2^2 + \lambda \tilde{P}_\alpha(\mathbf{q})$  and its dual norm  $\tilde{P}_\alpha^*(\mathbf{q})$  satisfies  $\tilde{P}_\alpha^*(\mathbf{q}) \leq \|\mathbf{q}\|_\infty$  (see (B.23) in the main text). From  $\mathcal{L}_\alpha^*(\hat{\mathbf{q}}; \lambda) \leq \mathcal{L}_\alpha^*(\mathbf{q}; \lambda)$ , it follows that

$$\begin{aligned} \left\| \hat{\mathcal{Q}}^{(d)} (\hat{\mathbf{q}} - \mathbf{q}) \right\|_2^2 + \lambda \tilde{P}_\alpha(\hat{\mathbf{q}}) &\leq 2(\hat{\mathbf{q}} - \mathbf{q})' \hat{\mathcal{Q}}^{(d)'} (\hat{\Delta}_1 - \hat{\Delta}_2 c - \hat{\Delta}_3 \beta^*) + \lambda \tilde{P}_\alpha(\mathbf{q}) \\ &\leq \tilde{P}_\alpha(\hat{\mathbf{q}} - \mathbf{q}) \left[ 2\tilde{P}_\alpha^* \left( \hat{\mathcal{Q}}^{(d)'} \hat{\Delta}_1 \right) + 2\tilde{P}_\alpha^* \left( \hat{\mathcal{Q}}^{(d)'} \hat{\Delta}_2 c \right) + 2\tilde{P}_\alpha^* \left( \hat{\mathcal{Q}}^{(d)'} \hat{\Delta}_3 \beta^* \right) \right] + \lambda \tilde{P}_\alpha(\mathbf{q}) \\ &\leq \tilde{P}_\alpha(\hat{\mathbf{q}} - \mathbf{q}) \left[ 2 \left\| \hat{\mathcal{Q}}^{(d)'} \hat{\Delta}_1 \right\|_\infty + 2 \left\| \hat{\mathcal{Q}}^{(d)'} \hat{\Delta}_2 c \right\|_\infty + 2 \left\| \hat{\mathcal{Q}}^{(d)'} \hat{\Delta}_3 \beta^* \right\|_\infty \right] + \lambda \tilde{P}_\alpha(\mathbf{q}). \end{aligned} \tag{B.34}$$

We subsequently define the following three sets:

$$\mathcal{H}_1^*(x) = \left\{ 2 \left\| \hat{\mathcal{Q}}^{(d)'} \hat{\Delta}_1 \right\|_\infty \leq x \right\}, \quad \mathcal{H}_2^*(x) = \left\{ 2 \left\| \hat{\mathcal{Q}}^{(d)'} \hat{\Delta}_2 c \right\|_\infty \leq x \right\}$$

and

$$\mathcal{H}_3^*(x) = \left\{ 2 \left\| \hat{\mathcal{Q}}^{(d)'} \hat{\Delta}_3 \beta^* \right\|_\infty \leq x \right\}.$$

Rescaling (B.34) by a factor 2 and on  $\mathcal{H}_1^*\left(\frac{x}{6}\right) \cap \mathcal{H}_2^*\left(\frac{x}{6}\right) \cap \mathcal{H}_3^*\left(\frac{x}{6}\right)$ , we get

$$2 \left\| \hat{\mathcal{Q}}^{(d)} (\hat{\mathbf{q}} - \mathbf{q}) \right\|_2^2 + 2\lambda \tilde{P}_\alpha(\hat{\mathbf{q}}) \leq \lambda \tilde{P}_\alpha(\hat{\mathbf{q}} - \mathbf{q}) + 2\lambda \tilde{P}_\alpha(\mathbf{q}). \tag{B.35}$$

We first manipulate the term  $\tilde{P}_\alpha(\hat{\mathbf{q}})$  in the LHS of (B.35). As in (B.19) of the main paper, we get

$$\tilde{P}_\alpha(\hat{\mathbf{q}}) \geq \tilde{P}_{\alpha, S^*}(\mathbf{q}) + \tilde{P}_{\alpha, S^{*c}}(\hat{\mathbf{q}} - \mathbf{q}) - \tilde{P}_{\alpha, S^*}(\hat{\mathbf{q}} - \mathbf{q}),$$

where  $S^*$  is the index sets of all nonzero coefficients in  $\mathbf{q}$ . For the term  $\tilde{P}_\alpha(\hat{\mathbf{q}} - \mathbf{q})$  in the RHS of (B.35), it follows from (B.20) that  $\tilde{P}_\alpha(\hat{\mathbf{q}} - \mathbf{q}) = \tilde{P}_{\alpha, S^*}(\hat{\mathbf{q}} - \mathbf{q}) + \tilde{P}_{\alpha, S^{*c}}(\hat{\mathbf{q}} - \mathbf{q})$ . We obtain

$$2 \left\| \hat{\mathcal{Q}}^{(d)} (\hat{\mathbf{q}} - \mathbf{q}) \right\|_2^2 + \lambda \tilde{P}_{\alpha, S^{*c}}(\hat{\mathbf{q}} - \mathbf{q}) \leq 3\lambda \tilde{P}_{\alpha, S^*}(\hat{\mathbf{q}} - \mathbf{q}). \tag{B.36}$$

Accordingly, with  $N^* = N_c + NK$ ,

$$(\hat{\mathbf{q}} - \mathbf{q}) \in \mathcal{C}_{N^*}^*(\mathcal{G}, S) := \left\{ \mathbf{A} \in \mathbb{R}^{N^*} : \tilde{P}_{\alpha, S^{*c}}(\mathbf{A}) \leq 3\tilde{P}_{\alpha, S^*}(\mathbf{A}) \right\}.$$



**Lemma 1** now goes through with  $\bar{\omega}_\alpha^*$  as in **Theorem 3**. For applicability of **Lemma 2**, we define  $\mathcal{V}^*(x) = \{\|\hat{\mathbf{Q}} - \mathbf{Q}\|_2 \leq x\}$  and assume  $\mathcal{V}^*(\frac{\phi_0^*}{2})$ . Derivations identical to those in (B.22) prove the main inequality of **Theorem 3**.

To derive the probability of the inequality being true, we need the probability of the occurrence of the event  $\mathcal{H}_1^*(\frac{\lambda}{6}) \cap \mathcal{H}_2^*(\frac{\lambda}{6}) \cap \mathcal{H}_3^*(\frac{\lambda}{6}) \cap \mathcal{V}^*(\frac{\phi_0^*}{2})$ . This probability is no smaller than

$$1 - \mathbb{P}\left(\mathcal{H}_1^*\left(\frac{\lambda}{6}\right)^c\right) - \mathbb{P}\left(\mathcal{H}_2^*\left(\frac{\lambda}{6}\right)^c\right) - \mathbb{P}\left(\mathcal{H}_3^*\left(\frac{\lambda}{6}\right)^c\right) - \mathbb{P}\left(\mathcal{V}^*\left(\frac{\phi_0^*}{2}\right)^c\right), \tag{B.37}$$

and all probabilities in (B.37) can be retraced to probabilities involving  $\|\hat{\mathbf{Q}} - \mathbf{Q}\|_{\cdot}$ . That is, bounding terms as in (B.24)–(B.25), we find

$$\begin{aligned} \mathbb{P}\left(\mathcal{H}_1^*\left(\frac{\lambda}{6}\right)^c\right) &= \mathbb{P}\left(\|\hat{\mathbf{Q}}^{(d)'} \hat{\mathbf{a}}_1\|_\infty > \frac{\lambda}{12}\right) \leq \mathbb{P}\left(\|\hat{\mathbf{Q}} - \mathbf{Q}\|_{\cdot}^2 + C_Q \|\hat{\mathbf{Q}} - \mathbf{Q}\|_{\cdot} > \frac{\lambda}{12}\right) \\ &\leq \mathbb{P}\left(\|\hat{\mathbf{Q}} - \mathbf{Q}\|_{\cdot} > \frac{\lambda^{1/2}}{\sqrt{24}}\right) + \mathbb{P}\left(\|\hat{\mathbf{Q}} - \mathbf{Q}\|_{\cdot} > \frac{\lambda}{24C_Q}\right), \end{aligned} \tag{B.38}$$

$$\begin{aligned} \mathbb{P}\left(\mathcal{H}_2^*\left(\frac{\lambda}{6}\right)^c\right) &= \mathbb{P}\left(\|\hat{\mathbf{Q}}^{(d)'} \hat{\mathbf{a}}_2 c\|_\infty > \frac{\lambda}{12}\right) \leq \mathbb{P}\left(\|\hat{\mathbf{Q}}^{(d)'} [\hat{\mathbf{Q}}^{(d)} - \mathbf{Q}^{(d)}]\|_\infty > \frac{\lambda}{12}\right) \\ &\leq \mathbb{P}\left(\|\hat{\mathbf{Q}} - \mathbf{Q}\|_{\cdot} > \frac{\lambda^{1/2}}{\sqrt{24}}\right) + \mathbb{P}\left(\|\hat{\mathbf{Q}} - \mathbf{Q}\|_{\cdot} > \frac{\lambda}{24C_Q}\right), \end{aligned} \tag{B.39}$$

and

$$\begin{aligned} \mathbb{P}\left(\mathcal{H}_3^*\left(\frac{\lambda}{6}\right)^c\right) &= \mathbb{P}\left(\|\hat{\mathbf{Q}}^{(d)'} \hat{\mathbf{a}}_3 \beta^*\|_\infty > \frac{\lambda}{12}\right) \\ &\leq \mathbb{P}\left(\|\hat{\mathbf{Q}} - \mathbf{Q}\|_{\cdot} > \frac{\lambda^{1/2}}{\sqrt{24C_\beta}}\right) + \mathbb{P}\left(\|\hat{\mathbf{Q}} - \mathbf{Q}\|_{\cdot} > \frac{\lambda}{24C_Q C_\beta}\right) \end{aligned} \tag{B.40}$$

For (B.39) and (B.40), **Assumption 1**\*(a)–(b) and **Assumption 1**\*(c) are needed to guarantee  $\|c\|_1 \leq 1$  and  $\|\beta^*\|_1 \leq C_\beta$ , respectively.

Also,  $\mathbb{P}\left(\mathcal{V}^*\left(\frac{\phi_0^*}{4}\right)^c\right) = \mathbb{P}\left(\|\hat{\mathbf{Q}} - \mathbf{Q}\|_2 > \frac{\phi_0^*}{2}\right) \leq \mathbb{P}\left(\|\hat{\mathbf{Q}} - \mathbf{Q}\|_{\cdot} > \frac{\phi_0^*}{2}\right)$ . Given these bounds, (B.37) translates to

$$\mathbb{P}\left(\mathcal{H}_1^*\left(\frac{\lambda}{6}\right) \cap \mathcal{H}_2^*\left(\frac{\lambda}{6}\right) \cap \mathcal{H}_3^*\left(\frac{\lambda}{6}\right) \cap \mathcal{V}^*\left(\frac{\phi_0^*}{2}\right)\right) \geq 1 - 7\mathbb{P}\left(\|\hat{\mathbf{Q}} - \mathbf{Q}\|_{\cdot} > 6f^*(\lambda, \phi_0^*)\right),$$

where  $f^*(\lambda, \phi_0^*) = \min\left\{\frac{\lambda^{1/2}}{12\sqrt{6}}, \frac{\lambda}{144C_Q}, \frac{\lambda^{1/2}}{12\sqrt{6C_\beta}}, \frac{\lambda}{144C_Q C_\beta}, \frac{\phi_0^*}{12}\right\}$ .

All that remains is a lower bound for  $\mathbb{P}\left(\|\hat{\mathbf{Q}} - \mathbf{Q}\|_{\cdot} > x\right)$ . We instead derive an upper bound for  $\mathbb{P}\left(\|\hat{\mathbf{Q}} - \mathbf{Q}\|_{\cdot} \leq x\right)$  as follows

$$\begin{aligned} \mathbb{P}\left(\|\hat{\mathbf{Q}} - \mathbf{Q}\|_{\cdot} \leq x\right) &= \mathbb{P}\left(\max\left\{\|\hat{\mathbf{Q}} - \mathbf{Q}\|_1, \|\hat{\mathbf{Q}} - \mathbf{Q}\|_\infty\right\} \leq x\right) \geq \mathbb{P}\left((K+2)N \|\hat{\mathbf{Q}} - \mathbf{Q}\|_{max} \leq x\right) \\ &= 1 - \mathbb{P}\left(\|\hat{\mathbf{Q}} - \mathbf{Q}\|_{max} > \frac{x}{(K+2)N}\right) \\ &\geq 1 - \mathbb{P}\left(\bigcup_{1 \leq i \leq (K+1)N, 1 \leq j \leq (K+2)N} |\hat{Q}_{ij} - Q_{ij}| > \frac{x}{(K+2)N}\right) \\ &\geq 1 - \sum_{i=1}^{(K+1)N} \sum_{j=1}^{(K+2)N} \mathbb{P}\left(\left|\sum_{t=2}^T \xi_{it} \xi_{jt} - \mathbb{E}(\xi_{it} \xi_{jt})\right| > \frac{Tx}{(K+2)N}\right) \\ &\geq \begin{cases} 1 - (K+1)(K+2)N^2 \left[ \left(b_1 T^{(1-\delta)/3} + \frac{(K+2)N b_3}{x}\right) \exp\left(-\frac{T^{(1-\delta)/3}}{2b_1^2}\right) + \frac{b_2(K+2)^d N^d}{x^d T^{\frac{\delta}{2}(d-1)}} \right] \\ \text{(polynomial tails),} \\ 1 - (K+1)(K+2)N^2 \left[ \frac{\kappa_1(K+2)N}{x} + \frac{2}{\kappa_2} \left(\frac{Tx^2}{(K+2)N}\right)^{1/7} \right] \exp\left(-\frac{1}{\kappa_3} \left(\frac{Tx^2}{(K+2)^2 N^2}\right)^{1/7}\right) \\ \text{(exponential tails),} \end{cases} \end{aligned}$$

where  $\xi_{it}$  denotes a generic element of an autocovariance matrix (as in **Lemma 3**). In accordance with **Theorem 3**, these RHS probabilities are equivalent to  $1 - \mathcal{P}_1^*(x, N, T)$  (polynomial tails) and  $1 - \mathcal{P}_2^*(x, N, T)$  (exponential tails).  $\square$

**Appendix C. Supplementary data**

Supplementary material related to this article can be found online at <https://doi.org/10.1016/j.jeconom.2023.105520>.

## References

- Ahrens, A., Bhattacharjee, A., 2015. Two-step lasso estimation of the spatial weights matrix. *Econometrics* 3, 128–155.
- Bickel, P.J., Ritov, Y., Tsybakov, A.B., 2009. Simultaneous analysis of Lasso and Dantzig selector. *Ann. Statist.* 37, 1705–1732.
- Brunekreef, B., Holgate, S.T., 2002. Air pollution and health. *Lancet* 360 (9341), 1233–1242.
- Debarys, N., LeSage, J., 2018. Flexible dependence modeling using convex combinations of different types of connectivity structures. *Reg. Sci. Urban Econ.* 69, 48–68.
- Dou, B., Parrell, M.L., Yao, Q., 2016. Generalized Yule-Walker estimation for spatio-temporal models with unknown diagonal coefficients. *J. Econometrics* 194, 369–382.
- Gao, Z., Ma, Y., Wang, H., Yao, Q., 2019. Banded spatio-temporal autoregressions. *J. Econometrics* 208, 211–230.
- Van de Geer, S., Bühlmann, P., Ritov, Y., Dezeure, R., 2014. On asymptotically optimal confidence regions and tests for high-dimensional models. *Ann. Statist.* 42, 1166–1202.
- Gelper, S., Wilms, I., Croux, C., 2016. Identifying demand effects in a large network of product categories. *J. Retail.* 92, 25–39.
- Guo, S., Wang, Y., Yao, Q., 2016. High-dimensional and banded vector autoregressions. *Biometrika* 103, 889–903.
- Hamilton, J.D., 1994. *Time Series Analysis*, first ed. Princeton University Press.
- Hansen, P.R., Lunde, A., Nason, J.M., 2011. The model confidence set. *Econometrica* 79 (2), 453–497.
- Hyndman, R.J., Athanasopoulos, G., 2018. *Forecasting: Principles and Practice*. OTexts.
- Junger, W.L., Ponce De Leon, A., 2015. Imputation of missing data in time series for air pollutants. *Atmos. Environ.* 102, 96–104.
- Knight, K., Fu, W., 2000. Asymptotics for LASSO-type estimators. *Ann. Statist.* 28, 1356–1378.
- Kock, A.B., Callot, L., 2015. Oracle inequalities for high dimensional vector autoregressions. *J. Econometrics* 186, 325–344.
- Lam, C., Souza, P.C.L., 2014. Regularization for Spatial Panel Time Series using the Adaptive Lasso. Working paper. London School of Economics and Political Science.
- Lam, C., Souza, P.C.L., 2019. Estimation and selection of spatial weight matrix in a spatial lag model. *J. Bus. Econom. Statist.* 38, 1–41.
- Lee, L.-F., 2004. Asymptotic distributions of quasi-maximum likelihood estimators for spatial autoregressive models. *Econometrica* 72, 1899–1925.
- Lee, L.-F., Yu, J., 2010. Estimation of spatial autoregressive panel data models with fixed effects. *J. Econometrics* 154, 165–185.
- Lee, L.-f., Yu, J., 2014. Efficient GMM estimation of spatial dynamic panel data models with fixed effects. *J. Econometrics* 180, 174–197.
- Leeb, H., Pötscher, B.M., 2005. Model selection and inference: Facts and fiction. *Econometric Theory* 21, 21–59.
- Leeb, H., Pötscher, B.M., 2008. Sparse estimators and the oracle property, or the return of Hodges' estimator. *J. Econometrics* 142 (1), 201–211.
- Ma, Y., Guo, S., Wang, H., 2023. Sparse spatio-temporal autoregressions by profiling and bagging. *J. Econometrics* 232 (1), 132–147.
- Masini, R.P., Medeiros, M.C., Mendes, E.F., 2022. Regularized estimation of high-dimensional vector autoregressions with weakly dependent innovations. *J. Time Series Anal.* 43 (4), 532–557.
- Medeiros, M.C., Mendes, E.F., 2016.  $\ell_1$ -regularization of high-dimensional time series models with non-Gaussian and heteroskedastic errors. *J. Econometrics* 191, 255–271.
- Simon, N., Friedman, J., Hastie, T., Tibshirani, R., 2013. A sparse-group lasso. *J. Comput. Graph. Statist.* 22, 231–245.
- Wang, D., Tsay, R.S., 2023. Rate-optimal robust estimation of high-dimensional vector autoregressive models. *Ann. Statist.* 51 (2), 846–877.
- Yu, J., de Jong, R.M., Lee, L.-F., 2008. Quasi-maximum likelihood estimators for spatial dynamic panel data with fixed effects when both  $n$  and  $T$  are large. *J. Econometrics* 146, 118–134.
- Yu, J., de Jong, R., Lee, L.-f., 2012. Estimation for spatial dynamic panel data with fixed effects: The case of spatial cointegration. *J. Econometrics* 167, 16–37.
- Zhang, X., Yu, J., 2018. Spatial weight matrix selection and model averaging for spatial autoregressive models. *J. Econometrics* 203, 1–18.
- Zou, H., 2006. The adaptive lasso and its oracle properties. *J. Amer. Statist. Assoc.* 101, 1418–1429.