

PA

Online Appendix for "A Fast Estimator for Binary Choice Models with Spatial, Temporal, and Spatio-Temporal Interdependence"

**Julian Wucherpfennig¹, Aya Kachi², Nils-Christian
Bormann³, and Philipp Hunziker⁴**

¹*Centre for International Security, Hertie School, Friedrichstraße 180, 10117 Berlin, Germany.*

²*Faculty of Business and Economics, University of Basel, Peter Merian-Weg 6, 4052, Basel,
Switzerland. Email: aya.kachi@unibas.ch*

³*Department of Philosophy, Politics & Economics, Witten/Herdecke University, Alfred-Herrhausen-Str.
50, Witten, Germany*

⁴*Network Science Institute, Northeastern University, 177 Huntington Ave, Boston, MA 02115. Now at
Google.*

Political Analysis (2020)

DOI: 10.1017/pan.xxxx.xx

Corresponding author

Aya Kachi

Edited by

John Doe

© The Author(s) 2020. Published
by Cambridge University Press
on behalf of the Society for
Political Methodology.

Online Appendix

1 Binary choice models with spatial, temporal, and spatio-temporal interdependence (Full description)

This section specifies and derives a mathematical expression of binary choice models, for which we develop a pseudo maximum likelihood estimator later. We do so for a model with spatial, temporal, and spatio-temporal interdependence, respectively. Note that, in this specification, we try to maintain general mathematical expressions without assuming a specific marginal distribution (such as logistic vs. normal). In fact, the PMLE's estimation feasibility *regardless* of the error-term probability distribution is one of the strengths of this approach. In our view, this strength goes beyond the probit-vs.-logit consideration. This can become useful when one might need to develop an estimator, for instance, for a hybrid of a binary spatial model and another model from a different model class such as duration and count.

1.1 Spatial interdependence

We consider the following model

$$y_i^* = \rho \sum_{j=1}^N w_{ij} y_j^* + \mathbf{x}_i \boldsymbol{\beta} + u_i \quad (1)$$

$$y_i = \begin{cases} 1 & \text{if } y_i^* > 0 \\ 0 & \text{otherwise,} \end{cases} \quad (2)$$

where y_i^* is a continuous latent outcome variable, w_{ij} is a spatial lag between unit i and j indicating how closely the two units are connected in a given space (e.g, geographical proximity, membership in the same organizations etc.), ρ is the spatial autocorrelation parameter, \mathbf{x}_i is a $1 \times k$ vector of covariates with parameter vector $\boldsymbol{\beta}$, and u_i is a zero-mean iid error term with fixed variance. We call this specification the binary spatial autoregressive model (or binary spatial model as we sometimes mention interchangeably). Note that in this specification, spatial dependence occurs on the level of the latent (i.e. not observed) outcome y_i^* . This specification follows Franzese, Hays, and Cook (2016), implying actors of our interest can observe or know more or less what others' latent characteristics are, and not only their revealed binary actions.

It is useful to write the latent equation in matrix notation, yielding

$$\mathbf{y}_{(N \times 1)}^* = \rho \mathbf{W} \mathbf{y}^* + \mathbf{X} \boldsymbol{\beta} + \mathbf{u}, \quad (3)$$

where

$$\mathbf{W}_{N \times N} = \begin{pmatrix} 0 & w_{12} & \cdots & w_{1N} \\ w_{21} & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & w_{N-1,N} \\ w_{N1} & \cdots & w_{N,N-1} & 0 \end{pmatrix}. \quad (4)$$

\mathbf{W} is commonly referred to as the *spatial weights matrix*. Throughout the paper we assume that \mathbf{W} is row-standardized. Doing so ensures that the spatial process defined in (3) is stationary as long as $|\rho| < 1$ (Kelejian and Prucha 2010). Given (3) we can derive the reduced form as

$$\begin{aligned} \mathbf{y}^* &= (\mathbf{I} - \rho \mathbf{W})^{-1} \mathbf{X} \boldsymbol{\beta} + (\mathbf{I} - \rho \mathbf{W})^{-1} \mathbf{u} \\ &= (\mathbf{I} - \rho \mathbf{W})^{-1} \mathbf{X} \boldsymbol{\beta} + \mathbf{v}, \end{aligned} \quad (5)$$

where vector \mathbf{v} contains the reduced-form error terms with non-spherical covariance matrix structure due to the multiplier $(\mathbf{I} - \rho\mathbf{W})^{-1}$.

The main component of the (pseudo) likelihood function of our interest will be the joint probability for the observed random variable Y given the model parameters and regressors. This leads to the following expression:

$$\begin{aligned}
 P(y = 1) &= P(y_i^* > 0) \\
 &= P\left(\left[(\mathbf{I} - \rho\mathbf{W})^{-1}\mathbf{X}\boldsymbol{\beta}\right]_i + v_i > 0\right) \\
 &= P\left(v_i > -\left[(\mathbf{I} - \rho\mathbf{W})^{-1}\mathbf{X}\boldsymbol{\beta}\right]_i\right) \\
 &= 1 - P\left(v_i \leq -\left[(\mathbf{I} - \rho\mathbf{W})^{-1}\mathbf{X}\boldsymbol{\beta}\right]_i\right) \\
 &= 1 - F_{V_i}\left(-\left[(\mathbf{I} - \rho\mathbf{W})^{-1}\mathbf{X}\boldsymbol{\beta}\right]_i\right).
 \end{aligned} \tag{6}$$

where $[\cdot]_i$ indicates the i 'th element of the vector $[\cdot]$. $F_{V_i}(\cdot)$ is the marginal CDF of the random variable V_i (the reduced form error term for unit i). Therefore, expression $F_{V_i}\left(-\left[(\mathbf{I} - \rho\mathbf{W})^{-1}\mathbf{X}\boldsymbol{\beta}\right]_i\right)$ is the marginal CDF of V_i evaluated at $-\left[(\mathbf{I} - \rho\mathbf{W})^{-1}\mathbf{X}\boldsymbol{\beta}\right]_i$. By definition the marginal CDF of V_i is

$$\begin{aligned}
 &F_{V_i}\left(-\left[(\mathbf{I} - \rho\mathbf{W})^{-1}\mathbf{X}\boldsymbol{\beta}\right]_i\right) \\
 &= \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \int_{-\infty}^{-\left[(\mathbf{I} - \rho\mathbf{W})^{-1}\mathbf{X}\boldsymbol{\beta}\right]_i} \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} f_{\mathbf{V}}(s_1, \dots, s_i, \dots, s_N) ds_1 \cdots ds_i \cdots ds_N,
 \end{aligned} \tag{7}$$

where $f_{\mathbf{V}}(s_1, \dots, s_N)$ is the joint PDF of the reduced-form error. The estimation challenge for binary choice models arises when evaluating F_{V_i} at $-\left[(\mathbf{I} - \rho\mathbf{W})^{-1}\mathbf{X}\boldsymbol{\beta}\right]_i$ is analytically intractable (as long as $\rho \neq 0$) (Anselin 2002). As a consequence, direct maximum likelihood estimation of $\boldsymbol{\beta}$ and ρ is not always feasible. Of course, one common exception is spatial probit, where the marginal probability has a closed-form expression.

Using this expression for the choice probability, $P(y = 1)$, we have the following expression that is proportional to the (pseudo) likelihood function for a binary spatial autoregressive model.

$$\begin{aligned}
 L(\rho, \boldsymbol{\beta} | \mathbf{X}, \mathbf{y}) &= \left[\prod_{i=1}^N P(y_i = 1)^{y_i} \right] \left[\prod_{i=1}^N P(y_i = 0)^{(1-y_i)} \right] \\
 &= \left[\prod_{i=1}^N P(y_i = 1)^{y_i} \right] \left[\prod_{i=1}^N \left(1 - P(y_i = 1)\right)^{(1-y_i)} \right].
 \end{aligned} \tag{8}$$

1.2 Temporal dependence

As an intermediate step toward the binary spatio-temporal model—for which our proposed estimator would eventually be useful—we first illustrate a binary temporal autoregressive model, where the latent outcome exhibits a first-order temporal autoregressive process governed by the temporal autocorrelation parameter γ with $|\gamma| < 1$. The structural form error term u_t is a zero-mean iid error term with fixed variance.¹

$$y_t^* = \mathbf{X}_t\boldsymbol{\beta} + \gamma y_{t-1}^* + u_t \tag{9}$$

$$y_t = \begin{cases} 1 & \text{if } y_t^* > 0 \\ 0 & \text{otherwise.} \end{cases} \tag{10}$$

1. The following results generalize trivially to higher-order processes.

As it falls out of the main contribution of this paper, we are grossly skipping over the rich time-series methods literature here and we are aware of it. For a discussion of this class of models in a political science context, see Beck *et al.* 2001, for example.

Next, note that we can rewrite the model in matrix notation as follows (equation(11)). One might argue that matrix notation of a time-series model is not the most useful expression in terms of estimating model parameters; and yet, as a stepping stone toward the binary spatio-temporal model, it is an analytically appealing expression.

$$\mathbf{y}_{(T \times 1)}^* = \mathbf{X}\boldsymbol{\beta} + \gamma\mathbf{T}\mathbf{y}^* + \mathbf{u}, \tag{11}$$

where \mathbf{T} , called the *temporal weights matrix*, is defined as

$$\mathbf{T} = \begin{pmatrix} 0 & 0 & 0 & \dots & 0 \\ 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 0 \end{pmatrix}. \tag{12}$$

It is evident that this model is mathematically comparable to the binary spatial model, the sole difference being that now we impose a weights matrix where the first subdiagonal (all the 1's) maps y_{t-1}^* to y_t^* . The reduced form of the autoregressive model is given by

$$\mathbf{y}_{(T \times 1)}^* = (\mathbf{I} - \gamma\mathbf{T})^{-1}\mathbf{X}\boldsymbol{\beta} + (\mathbf{I} - \gamma\mathbf{T})^{-1}\mathbf{u} \tag{13}$$

As one might see it already, this gives rise to a similar difficulties in ML estimation as the binary spatial model described above.

1.3 Spatio-temporal interdependence

So far, we have illustrated that spatial and temporal interdependence give rise to the same reduced form expression for the latent outcome vector \mathbf{y}^* , and are thus all subject to the same estimation challenge whenever the joint probability $P(y = 1)$ does not have a closed-form expression. This similarity in functional form allows us to combine different dependency structures relatively straightforwardly, yielding models exhibiting multiple types of dependencies among observations. In the following, we consider the binary spatio-temporal autoregressive model (STAR), which combines the binary spatial autoregressive model with the temporal autoregressive model, yielding a panel setup (see e.g. Franzese, Hays, and Cook 2016). The binary STAR model is given by

$$\mathbf{y}_{(NT \times 1)}^* = \mathbf{Q}\mathbf{y}^* + \mathbf{X}\boldsymbol{\beta} + \mathbf{u}, \tag{14}$$

where $\mathbf{y}^* = [y_{\cdot 1}^*, \dots, y_{\cdot T}^*]'$ and $y_{\cdot t}^* = [y_{1t}^*, \dots, y_{Nt}^*]'$. Hence, the cross-sectional $y_{\cdot t}^*$ vectors are stacked “on top of each other”. The \mathbf{X} matrix is constructed analogously. \mathbf{Q} is given by

$$\mathbf{Q}_{NT \times NT} = \rho\mathbf{W}^* + \gamma\mathbf{T}^*, \tag{15}$$

where \mathbf{W}^* is the block-diagonal *panel spatial weights matrix* given by

$$\mathbf{W}^* = \begin{pmatrix} \mathbf{W} & 0 & 0 & \dots & 0 \\ 0 & \mathbf{W} & 0 & \dots & 0 \\ 0 & 0 & \mathbf{W} & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & \mathbf{W} \end{pmatrix}, \quad (16)$$

and \mathbf{T}^* is the *panel temporal weights matrix* given by

$$\mathbf{T}^* = \begin{pmatrix} 0 & 0 & 0 & \dots & 0 \\ I_N & 0 & 0 & \dots & 0 \\ 0 & I_N & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 0 \end{pmatrix}, \quad (17)$$

where I_N is the $N \times N$ identity matrix.

The reduced form of the spatio-temporal autoregressive model is given by

$$\mathbf{y}_{(NT \times 1)}^* = (\mathbf{I} - \mathbf{Q})^{-1} \mathbf{X} \boldsymbol{\beta} + (\mathbf{I} - \mathbf{Q})^{-1} \mathbf{u}, \quad (18)$$

which again gives rise to the familiar complication.

2 A pseudo maximum likelihood estimator for interdependent binary outcomes (Full description)

This section describes the PMLE estimator to tackle spatial, temporal, and spatio-temporal forms of interdependence for binary outcome data. Our analytical point of departure is a pseudo maximum likelihood estimator (PMLE) for binary spatially autoregressive models described in Smirnov (2010), for which the remaining computational burden amounts to inverting an N -dimensional matrix we refer to as the “interdependence multiplier.” We extend the PMLE to cases of temporal and spatio-temporal interdependence, which is a tool that is so far only offered by the Franzese’ et al.’s RIS estimator (2016). We further reduce the estimation costs by proposing an implementation strategy that avoids direct matrix inversion, and instead relies on a combination of iterative gradient procedures and approximations that yield an estimation algorithm with almost linear complexity in N . This additional procedure we propose will be detailed separately in the following section.

When direct ML estimation is infeasible for binary models featuring interdependence of the outcome variables (due to the lack of a closed-form cdf that goes into $P(y = 1)$), it is clear that we require an alternative approach. One option is simulation. Franzese, Hays, and Cook (2016) and Calabrese and Elkink (2014b) provide extensive reviews of the spatial probit literature, and useful comparisons of several simulation-based estimation methods such as recursive-importance-sampling (RIS) and Bayesian MCMC approaches (see also Calabrese and Elkink (2014a) for cases with asymmetric link functions accommodating rare events). Similarly, Beck *et al.* 2001 discuss a Bayesian estimation strategy for the binary temporal autoregressive model. However, simulation-based approaches place a number of burdens on the researchers. First, they are computationally intensive and it usually takes a long time to estimate them. Estimation time can be prohibitive if researchers work with big data and do not have access to high-performance computing clusters. Second, convergence problems in MCMC simulations often require tedious hyperparameter tuning

and exacerbate the estimation-time problem. Third, and as perhaps the most broadly relevant point, currently, applied researchers do not have access to more than the most basic tools, for example, cross-sectional spatial probit estimators. For these reasons, we now introduce a pseudo maximum likelihood (PML) method as a feasible way to reduce estimation time, minimize convergence problems, and enable applied researchers to run models that more clearly address their research problems. Our estimator builds on Smirnov’s (2010) spatial PML estimator and extends it to temporal and spatio-temporal interdependence.

2.1 PMLE for the binary spatial model

Recall the reduced form for the binary spatial model is given by

$$\begin{aligned} \mathbf{y}^* &= (\mathbf{I} - \rho\mathbf{W})^{-1}\mathbf{X}\boldsymbol{\beta} + (\mathbf{I} - \rho\mathbf{W})^{-1}\mathbf{u} \\ &= (\mathbf{I} - \rho\mathbf{W})^{-1}\mathbf{X}\boldsymbol{\beta} + \mathbf{v}. \end{aligned} \tag{19}$$

Denote the spatial multiplier by \mathbf{Z} ,

$$\mathbf{Z}_{(N \times N)} = (\mathbf{I} - \rho\mathbf{W})^{-1}, \tag{20}$$

and, by \mathbf{D} , an $N \times N$ matrix that contains diagonal elements of \mathbf{Z} . All off-diagonal elements of \mathbf{D} are zero. The spatial multiplier indicates the degree of local and global spillovers of an exogenous shock that unit i receives (Anselin 2003); in other words, $z_{ij} = \frac{\partial y_i^*}{\partial u_j}$, where z_{ij} is the ij th element of \mathbf{Z} . The diagonal matrix \mathbf{D} indicates “private effects,” borrowing Smirnov’s (2010) term, of exogenous shocks on the individual latent outcomes. The relative effect captured by \mathbf{D} is “private” in that it indicates the magnitude of the effect that unit i receives from an exogenous shock that occurred to unit i itself; in other words, $d_i = \frac{\partial y_i^*}{\partial u_i}$.

On the other hand, the off-diagonal elements of \mathbf{Z} , i.e. $\mathbf{Z} - \mathbf{D}$, represent “aggregate spatial effects” of an exogenous shock. Note that all diagonal elements of $\mathbf{Z} - \mathbf{D}$ are zero. One could interpret it as an aggregate spillover effects that unit i receives from an exogenous shock through all the other units.

The reduced form can now be re-written as

$$\mathbf{y}_{(N \times 1)}^* = \mathbf{Z}\mathbf{X}\boldsymbol{\beta} + \underbrace{(\mathbf{Z} - \mathbf{D})\mathbf{u}}_{\text{“Social effects”}} + \underbrace{\mathbf{D}\mathbf{u}}_{\text{“Private effects”}}, \tag{21}$$

or, for each unit i ,

$$y_i^* = \sum_j \beta z_{ij} x_j + \sum_j [\mathbf{Z} - \mathbf{D}]_{ij} u_j + d_i u_i. \tag{22}$$

We can now rewrite the probability of unit i seeing a positive outcome as

$$\begin{aligned} P(y_i = 1) &= P(y_i^* \geq 0) \\ &= P\left(\sum_j \beta z_{ij} x_j + \sum_j [\mathbf{Z} - \mathbf{D}]_{ij} u_j + d_i u_i \geq 0\right) \\ &= P\left(u_i \leq \frac{\sum_j \beta z_{ij} x_j}{d_i} + \frac{\sum_j [\mathbf{Z} - \mathbf{D}]_{ij} u_j}{d_i}\right). \end{aligned} \tag{23}$$

Note that there is a stochastic element left in the argument of the probability in the above expression: $\sum_j [\mathbf{Z} - \mathbf{D}]_{ij} u_j$. In order to allow for an analytical formulation of a (pseudo) likelihood, we assume that higher-order effects can be “ignored”. Behaviorally, this means that observations may simplify their choice by not worrying about aggregate spatial effects of a random shock that are experienced by other (connected) observations. That is, mathematically, we do not expect a

systematic effect of a random shock on unit i that is carried through the off-diagonal elements of the spatial multiplier; i.e., it does not affect the choice probability systematically. The assumption is warranted because u_{it} are i.i.d with mean 0. Smirnov's (2010) key proposal is to approximate $\sum_j [\mathbf{Z} - \mathbf{D}]_{ij} u_j$ in (23) by its expectation, i.e., zero. This step simplifies the likelihood function. To see why, note that $P(y_{it} = 1)$ can now be written as follows:

$$\begin{aligned} P(y_i = 1) &= P(y_i^* \geq 0) \\ &= P\left(u_i \leq \frac{\sum_j \beta z_{ij} x_j}{d_i}\right) \\ &= F_u\left(\frac{\sum_j \beta z_{ij} x_j}{d_i}\right), \end{aligned} \tag{24}$$

where $F_u(\cdot)$ is the cdf of the *univariate* distribution of u_i , which is typically the standard normal (yielding a Probit model) or a standard logistic (yielding a Logit model).

With this approximation, we can write the pseudo likelihood in closed form. If u_i follows the standard logistic distribution, for instance, we have

$$\begin{aligned} PL(\rho, \boldsymbol{\beta} | \mathbf{X}, \mathbf{y}) &= \left[\prod_{i=1}^N P(y_i = 1)^{y_i} \right] \left[\prod_{i=1}^N (1 - P(y_i = 1))^{(1-y_i)} \right] \\ &\propto \left[\prod_{i=1}^N \frac{\exp((\sum_j \beta z_{ij} x_j)/d_i)}{1 + \exp((\sum_j \beta z_{ij} x_j)/d_i)} \right]^{y_i} \\ &\times \left[\prod_{i=1}^N \frac{1}{1 + \exp((\sum_j \beta z_{ij} x_j)/d_i)} \right]^{(1-y_i)}. \end{aligned} \tag{25}$$

2.2 PMLE for the temporal autoregressive model

Recall the reduced form for the binary temporal autoregressive model, given by

$$\mathbf{y}_{(T \times 1)}^* = (\mathbf{I} - \gamma \mathbf{T})^{-1} \mathbf{X} \boldsymbol{\beta} + (\mathbf{I} - \gamma \mathbf{T})^{-1} \mathbf{u}. \tag{26}$$

Next, let

$$\mathbf{Z}_{(T \times T)} = (\mathbf{I} - \gamma \mathbf{T})^{-1}, \tag{27}$$

denote the dependency multiplier. Applying the logic of the previous section, we can decompose the reduced-form error term into two parts

$$\begin{aligned} \mathbf{y}^* &= \mathbf{Z} \mathbf{X} \boldsymbol{\beta} + \mathbf{Z} \mathbf{u} \\ &= \mathbf{Z} \mathbf{X} \boldsymbol{\beta} + \underbrace{(\mathbf{Z} - \mathbf{D}) \mathbf{u}}_{\text{distributed}} + \underbrace{\mathbf{D} \mathbf{u}}_{\text{contemporaneous}}. \end{aligned} \tag{28}$$

The *distributed* effect captures the effect of exogenous shocks that occurred in the past and were carried over to the outcome of time t . These are distributed because this term focuses on the effect that is carried across multiple time periods ("neighbors" in time). On the other hand, the contemporaneous effects capture the effect of an exogenous shock that occurred in the current time period on the current outcome. Note that due to the lower-diagonal structure of \mathbf{T} , $\mathbf{D} = \mathbf{I}$, and thus $d_i = 1$. Again substituting $(\mathbf{Z} - \mathbf{D}) \mathbf{u}$ with its expectation and given that u is i.i.d., we arrive at

the following expression for the probability of a positive outcome:

$$\begin{aligned}
 Pr(y_t = 1) & \\
 &= Pr(y_t^* > 0) \\
 &= Pr(u_t < [\mathbf{Z}\mathbf{X}\boldsymbol{\beta}]_t),
 \end{aligned} \tag{29}$$

and the pseudo likelihood function, for instance with a logit link function, is given by

$$\begin{aligned}
 PL(\gamma, \boldsymbol{\beta} | \mathbf{X}, \mathbf{y}) &= \left[\prod_{t=1}^T P(y_t = 1)^{y_t} \right] \left[\prod_{t=1}^T (1 - P(y_t = 1))^{(1-y_t)} \right] \\
 &\propto \left[\left(\prod_{t=1}^T 1 - \frac{1}{1 + \exp(-[\mathbf{Z}\mathbf{X}\boldsymbol{\beta}]_t/d_t)} \right)^{y_t} \right] \\
 &\quad \times \left[\left(\prod_{t=1}^T \frac{1}{1 + \exp(-[\mathbf{Z}\mathbf{X}\boldsymbol{\beta}]_t/d_t)} \right)^{(1-y_t)} \right].
 \end{aligned} \tag{30}$$

2.3 PMLE for the spatio-temporal autoregressive model

Similarly to the above model, recall the reduced form:

$$\mathbf{y}_{(NT \times 1)}^* = (\mathbf{I} - \mathbf{Q})^{-1} \mathbf{X}\boldsymbol{\beta} + (\mathbf{I} - \mathbf{Q})^{-1} \mathbf{u}. \tag{31}$$

We denote the spatio-temporal multiplier $(\mathbf{I} - \mathbf{Q})^{-1}$ again as $\mathbf{Z}_{(NT \times NT)}$ and define the matrix $\mathbf{D}_{NT \times NT}$ as a matrix that captures the diagonal elements of \mathbf{Z} with all other elements being zeros.

$$\begin{aligned}
 \mathbf{y}_{(NT \times 1)}^* &= \mathbf{Z}\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} \\
 &= \mathbf{Z}\mathbf{X}\boldsymbol{\beta} + \underbrace{(\mathbf{Z} - \mathbf{D})\mathbf{u}}_{\text{higher-order effects}} + \underbrace{\mathbf{D}\mathbf{u}}_{\text{zero-order effects}}.
 \end{aligned} \tag{32}$$

Substituting $(\mathbf{Z} - \mathbf{D})\mathbf{u}$ with its expectation, we arrive at the following expression for the probability of a positive outcome:

$$\begin{aligned}
 Pr(y_{it} = 1) & \\
 &= Pr(y_{it}^* > 0) \\
 &= Pr(u_{it} < [\mathbf{Z}\mathbf{X}\boldsymbol{\beta}]_{it})
 \end{aligned} \tag{33}$$

and the pseudo likelihood function again with a logit link function, is given by

$$\begin{aligned}
 PL(\rho, \gamma, \boldsymbol{\beta} | \mathbf{X}, \mathbf{y}) &= \left[\prod_{i=1}^N \prod_{t=1}^T P(y_{it} = 1)^{y_{it}} \right] \left[\prod_{i=1}^N \prod_{t=1}^T (1 - P(y_{it} = 1))^{(1-y_{it})} \right] \\
 &\propto \left[\left(\prod_{i=1}^N \prod_{t=1}^T 1 - \frac{1}{1 + \exp(-[\mathbf{Z}\mathbf{X}\boldsymbol{\beta}]_{ij,t}/d_{it})} \right)^{y_{it}} \right] \\
 &\quad \times \left[\left(\prod_{i=1}^N \prod_{t=1}^T \frac{1}{1 + \exp(-[\mathbf{Z}\mathbf{X}\boldsymbol{\beta}]_{ij,t}/d_{it})} \right)^{(1-y_{it})} \right].
 \end{aligned} \tag{34}$$

Alternatively, for any binomial link function $g(\cdot)$, we have

$$PL(\rho, \gamma, \boldsymbol{\beta} | \mathbf{X}, \mathbf{y}) \propto \prod_{i=1}^N \prod_{t=1}^T \left[g^{-1} \left(\frac{[\mathbf{Z}\mathbf{X}\boldsymbol{\beta}]_{ij,s}}{d_{it}} \right)^{y_{it}} \left[1 - g^{-1} \left(\frac{[\mathbf{Z}\mathbf{X}\boldsymbol{\beta}]_{ij,s}}{d_{it}} \right) \right]^{(1-y_{it})} \right]. \tag{35}$$

Note that this expression requires an estimate for the values for \mathbf{y}_{i0}^* , i.e. the values *preceding* the first observed period in order to calculate the first period \mathbf{y}_{i1}^* .² Assuming mean stationarity, we draw on Kauppi and Saikkonen (2008) and use what can be viewed as the unconditional expectation of \mathbf{y}^* across all time period (and units): $E[\mathbf{y}^*] = (\mathbf{I} - \rho\mathbf{W} - \gamma)^{-1}\bar{\mathbf{X}}\boldsymbol{\beta}$, where $\bar{\mathbf{X}}$ are the sample means.

3 Speeding up computation further

3.1 Why still costly...

In the previous section, we have derived pseudo likelihood functions for binary (inter-)dependence models that can be evaluated directly, thus permitting a pseudo maximum likelihood (PML) strategy that does not require simulation. However, naive implementations of the proposed PML estimator may still be prohibitively costly to run. To see why, let us assume that we attempt to fit a model on data covering N units over T periods with reduced form

$$\mathbf{y}_{NT \times NT}^* = \mathbf{Z}\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}, \quad (36)$$

where $\mathbf{Z} = \mathbf{A}^{-1} = (\mathbf{I} - \mathbf{Q})^{-1}$. This specification yields a pseudo likelihood function consisting of NT terms of the following form

$$\begin{aligned} P(y_j = 1) &= P(y_j^* \geq 0) \\ &= F_u\left(\frac{\mu_j}{d_j}\right), \end{aligned} \quad (37)$$

with $j \in \{1, 2, \dots, NT\}$, $\mu = \mathbf{Z}\mathbf{X}\boldsymbol{\beta}$, and $d_j = \mathbf{Z}_{jj}$. Perhaps the most straightforward implementation of expression (37) is to invert \mathbf{A} directly using a decomposition-based solver, then multiplying \mathbf{Z} with $\mathbf{X}\boldsymbol{\beta}$ to yield μ , and dividing by $\text{diag}(\mathbf{Z})$. However, this strategy is typically very slow, as most decomposition-based solvers operate with near cubic time complexity. Instead, we propose a strategy that avoids the full inversion of \mathbf{A} , but computes μ and d separately.

3.2 Computing μ

To compute μ we solve the linear system $\mathbf{A}\mu = \mathbf{X}\boldsymbol{\beta}$ for μ using an iterative method. In particular, we propose using the Biconjugate gradient stabilized method (Bi-CGSTAB), which yields similar performance to the more widely known conjugate gradient method, but is applicable even if \mathbf{A} is not symmetric (Van der Vorst 1992). Doing so yields a substantial speed-up over decomposition-based solvers, especially when \mathbf{A} is sparse, which will be the case as long as any spatial weights matrices entering \mathbf{A} are neighborhood based.³

If \mathbf{A} is block-diagonal, which is the case for all panel models that do not feature a temporal autoregressive term, then we can make use of the fact that the inverse of a block-diagonal matrix is the block diagonal matrix of block-wise inverses. In other words, instead of solving the full system, we can solve $\mathbf{A}_t\mu_t = \mathbf{X}_t\boldsymbol{\beta}$ for all $t \in \{1, 2, \dots, T\}$, whereas \mathbf{A}_t represents a block in \mathbf{A} .

3.3 Computing d

First, we note that for panel models, d can always be computed in a period-wise fashion. This is obviously the case if \mathbf{Q} does *not* include the panel temporal weights matrix \mathbf{T}^* , because then \mathbf{A} is block-diagonal, and thus d is the concatenation of the period-wise diagonals $d_t = \text{diag}((\mathbf{A}_t)^{-1})$. Crucially, however, d can be computed analogously even if \mathbf{Q} *does* include \mathbf{T}^* , and thus \mathbf{A} is *not* block-diagonal. In the following, we provide a theorem to this end for the case where \mathbf{Q} represents

2. Because \mathbf{y}^* is latent, dropping the first period from the likelihood merely shifts the problem to the next period, rather than solving it (cf. Franzese, Hays, and Cook 2016).

3. Note that the weights matrices for temporal dependence (\mathbf{A}) and outcome-interdependence (\mathbf{M}) are sparse by construction.

a spatial weights matrix; we note, however, that the results extends to the case where \mathbf{Q} represents a panel outcome weights matrix (i.e. an \mathbf{M}^* term), or a mixture of the two.

Theorem 3.1 Let \mathbf{W}^* be a $NT \times NT$ block-diagonal panel spatial matrix as defined in (16), and \mathbf{T}^* be the panel temporal weights matrix as defined in (17). Then $d = \text{diag}((\mathbf{I} - \gamma\mathbf{T}^* - \rho\mathbf{W}^*)^{-1}) = \text{diag}((\mathbf{I} - \rho\mathbf{W}^*)^{-1})$.

Define a *strictly lower block triangular* (SLBT) matrix as any square matrix with the following structure

$$\begin{pmatrix} 0_N & & & & 0 \\ & 0_N & & & \\ & & \ddots & & \\ \neq 0 & & & 0_N & \\ & & & & 0_N \end{pmatrix},$$

whereas 0_N is the $N \times N$ matrix of zeros. It follows that $\gamma\mathbf{T}^*$ is SLBT. Note that $(\mathbf{I} - \gamma\mathbf{T}^* - \rho\mathbf{W}^*)^{-1}$ can be written as a Neumann series (LeSage and Pace 2009, ch. 2):

$$\begin{aligned} (\mathbf{I} - \gamma\mathbf{T}^* - \rho\mathbf{W}^*)^{-1} &= \mathbf{I} + \sum_{l=1}^L (\gamma\mathbf{T}^* + \rho\mathbf{W}^*)^l \\ &= \mathbf{I} + (\gamma\mathbf{T}^* + \rho\mathbf{W}^*) \\ &\quad + (\gamma^2\mathbf{T}^{*2} + \gamma\rho\mathbf{T}^*\mathbf{W}^* + \gamma\rho\mathbf{W}^*\mathbf{T}^* + \rho^2\mathbf{W}^{*2}) \\ &\quad + \dots \end{aligned} \tag{38}$$

Note that the product of two SLBT matrices is SLBT. Further note that the product of an SLBT matrix with a block-diagonal (BD) matrix is SLBT, regardless of the order of multiplication. It follows that the only terms in (38) with non-zero diagonals are of the form $\rho^l\mathbf{W}^{*l}$ for $l > 1$. Thus,

$$\begin{aligned} \text{diag}((\mathbf{I} - \gamma\mathbf{T}^* - \rho\mathbf{W}^*)^{-1}) &= \text{diag}(\mathbf{I}) + \sum_{l=1}^L \text{diag}((\gamma\mathbf{T}^* + \rho\mathbf{W}^*)^l) \\ &= \text{diag}(\mathbf{I}) + \sum_{l=1}^L \text{diag}((\rho\mathbf{W}^*)^l) \\ &= \text{diag}((\mathbf{I} - \rho\mathbf{W}^*)^{-1}). \end{aligned}$$

Thanks to the above theorem, we now only need an efficient method for calculating $d_t = \text{diag}((\mathbf{A}_t)^{-1})$. Here we propose two approaches. The first (preferred) one applies whenever \mathbf{A}_t is composed of only a single weights parameter and weights matrix, e.g. $\mathbf{A}_t = \mathbf{I} - \rho\mathbf{W}$ or $\mathbf{A}_t = \mathbf{I} - \lambda\mathbf{M}$. In this case, we make use of the fact that $(\mathbf{A}_t)^{-1}$ can be written as a Neumann series, e.g. for the spatial case

$$(\mathbf{A}_t)^{-1} = \mathbf{I} + \sum_{l=1}^{\infty} (\rho\mathbf{W})^l. \tag{39}$$

Thus, an approximation for d_t may be obtained via

$$d_t \approx \tilde{d}_t = \text{diag}(\mathbf{I}) + \sum_{l=1}^L \text{diag}(\rho^l\mathbf{W}^l), \tag{40}$$

with L suitably large; we use $L = 8$. Note that we can precompute the series $\{\mathbf{W}, \mathbf{W}^2, \dots, \mathbf{W}^L\}$ prior to optimization. Thus, the time complexity of evaluating \tilde{d}_t during optimization is linear in N .

The second approach for computing d_t comes into play when \mathbf{A}_t is composed of multiple weights matrices and parameters, as for instance in the binary simultaneous equation spatial model discussed in Section 1.3. In this case, we use the method of Takahashi, Fagan, and Chin 1973, and examined by Erisman and Tinney 1975, which relies on a recursive algorithm to calculate the diagonal of a matrix inverse. Importantly, using the Takahashi equations to calculate d_t is considerably faster than computing the full decomposition-based inverse if \mathbf{A}_t is sparse, which will generally be the case as long as any spatial weights matrices entering into \mathbf{A}_t are neighborhood-based.

4 Evaluation and comparison of estimation strategies

In the remainder of the Online Appendix, we present all the MC simulations referred to in the main text. In all cases, we compare the performance of our estimator to that of some well-recognized alternatives.

5 Notes on replication material

Tables A3 & A4 compare the spatio-temporal estimation results when using different approximations for y^* for the initial period. Table A4 is equivalent to Table 1 reported in the main paper, and we recommend using this estimation procedure. Replication results for Table A3 are available on request.

Further note that our replication material does not replicate Figure A3, which demonstrates estimation times. Unless the code is rerun on exactly the same hardware setup, estimation times will obviously differ.

Table A1. Simulation results for **spatial** PMLE (500 iterations)

	N=256			N=1,024			N=4,096		
	$\beta_0 = -0.5$	$\beta_1 = 1$	ρ	$\beta_0 = -0.5$	$\beta_1 = 1$	ρ	$\beta_0 = -0.5$	$\beta_1 = 1$	ρ
Experiment #1: $\rho = 0$									
Mean Coefficient Estimate	-0.530	1.019	-0.038	-0.507	1.006	-0.010	-0.501	1.002	-0.001
Mean Bias	0.139	0.101	0.205	0.063	0.049	0.097	0.032	0.024	0.051
RMSE	0.185	0.127	0.269	0.078	0.061	0.122	0.040	0.031	0.065
Actual SD of estimates	0.183	0.126	0.266	0.078	0.061	0.121	0.040	0.031	0.065
Overconfidence	1.089	0.991	1.080	0.961	0.976	0.954	0.990	0.984	1.007
Experiment #2: $\rho = 0.25$									
Mean Coefficient Estimate	-0.534	1.024	0.212	-0.503	1.004	0.245	-0.498	0.996	0.249
Mean Bias	0.140	0.110	0.183	0.064	0.050	0.086	0.032	0.025	0.045
RMSE	0.197	0.140	0.244	0.079	0.062	0.108	0.040	0.031	0.056
Actual SD of estimates	0.194	0.138	0.241	0.079	0.062	0.108	0.040	0.030	0.056
Overconfidence	1.224	1.061	1.178	1.050	0.980	1.060	1.063	0.963	1.090
Experiment #3: $\rho = 0.5$									
Mean Coefficient Estimate	-0.523	1.000	0.472	-0.497	0.987	0.492	-0.488	0.975	0.498
Mean Bias	0.147	0.115	0.139	0.069	0.054	0.069	0.037	0.034	0.036
RMSE	0.201	0.143	0.188	0.086	0.067	0.089	0.046	0.041	0.046
Actual SD of estimates	0.199	0.143	0.186	0.086	0.066	0.088	0.044	0.033	0.046
Overconfidence	1.309	1.051	1.223	1.221	0.985	1.209	1.256	0.991	1.251

Overconfidence is the standard deviation of the estimated parameter divided by the mean of its estimated standard error.

Table A2. Simulation results for **temporal** PMLE (500 iterations; y_0^* is estimated by $E(y_t^*)$)

	N=64 & T=4			N=64 & T=16			N=256 & T=16		
	$\beta_0 = -0.5$	$\beta_1 = 1$	γ	$\beta_0 = -0.5$	$\beta_1 = 1$	γ	$\beta_0 = -0.5$	$\beta_1 = 1$	γ
Experiment #1: $\gamma = 0$									
Mean Coefficient Estimate	-0.503	1.017	-0.001	-0.502	1.006	-0.001	-0.501	1.002	-0.000
Mean Bias	0.089	0.102	0.086	0.045	0.049	0.039	0.021	0.024	0.019
RMSE	0.113	0.128	0.106	0.057	0.061	0.048	0.026	0.031	0.024
Actual SD of estimates	0.113	0.127	0.106	0.057	0.061	0.048	0.026	0.031	0.024
Overconfidence	1.013	1.000	1.025	1.050	0.971	1.008	0.963	0.985	1.004
Experiment #2: $\gamma = 0.25$									
Mean Coefficient Estimate	-0.497	0.999	0.249	-0.488	0.980	0.249	-0.486	0.972	0.249
Mean Bias	0.083	0.108	0.075	0.045	0.052	0.035	0.023	0.034	0.018
RMSE	0.104	0.136	0.095	0.056	0.064	0.044	0.028	0.041	0.022
Actual SD of estimates	0.104	0.137	0.095	0.054	0.061	0.044	0.024	0.030	0.022
Overconfidence	1.049	1.038	1.036	1.126	0.974	0.975	1.005	0.970	0.965
Experiment #3: $\gamma = 0.5$									
Mean Coefficient Estimate	-0.470	0.936	0.495	-0.442	0.882	0.498	-0.440	0.878	0.499
Mean Bias	0.083	0.130	0.060	0.067	0.119	0.034	0.061	0.122	0.017
RMSE	0.101	0.159	0.076	0.079	0.133	0.043	0.065	0.126	0.021
Actual SD of estimates	0.097	0.146	0.076	0.054	0.062	0.043	0.025	0.031	0.021
Overconfidence	1.061	1.019	1.048	1.244	1.001	1.148	1.140	1.011	1.118

Overconfidence is the standard deviation of the estimated parameter divided by the mean of its estimated standard error.

Table A3. Simulation results for **spatio-temporal** PMLE (500 iterations) when $y_1^* = X_1\beta + u_1$ is assumed

	N=64 & T=4				N=64 & T=16				N=256 & T=16					
	$\beta_0 = -0.5$		$\beta_1 = 1$		$\beta_0 = -0.5$		$\beta_1 = 1$		$\beta_0 = -0.5$		$\beta_1 = 1$			
	ρ	γ	ρ	γ	ρ	γ	ρ	γ	ρ	γ	ρ	γ		
Experiment #1: $\gamma = 0.25, \rho = 0.25$														
Mean Coefficient Estimate	-0.515	1.005	0.252	0.234	0.234	0.250	-0.483	0.971	0.250	0.250	-0.483	0.965	0.249	0.250
Mean Bias	0.150	0.112	0.087	0.159	0.159	0.035	0.079	0.056	0.035	0.078	0.042	0.039	0.019	0.040
RMSE	0.203	0.147	0.108	0.206	0.206	0.045	0.100	0.070	0.045	0.097	0.051	0.047	0.024	0.051
Actual SD of estimates	0.203	0.147	0.109	0.206	0.206	0.046	0.098	0.064	0.046	0.098	0.049	0.031	0.024	0.051
Overconfidence	1.266	1.097	1.130	1.171	1.171	1.030	1.210	0.972	1.030	1.161	1.144	0.960	1.068	1.134
Experiment #2: $\gamma = 0.25, \rho = 0.5$														
Mean Coefficient Estimate	-0.512	0.992	0.255	0.482	0.482	0.248	-0.478	0.935	0.248	0.495	-0.465	0.927	0.247	0.501
Mean Bias	0.176	0.133	0.088	0.134	0.134	0.041	0.111	0.088	0.041	0.070	0.059	0.075	0.022	0.035
RMSE	0.241	0.166	0.110	0.178	0.178	0.052	0.142	0.106	0.052	0.093	0.071	0.083	0.028	0.044
Actual SD of estimates	0.241	0.166	0.110	0.177	0.177	0.052	0.141	0.083	0.052	0.093	0.062	0.040	0.027	0.044
Overconfidence	1.490	1.076	1.246	1.292	1.292	1.262	1.653	1.059	1.262	1.361	1.445	1.020	1.276	1.247
Experiment #3: $\gamma = 0.5, \rho = 0.25$														
Mean Coefficient Estimate	-0.486	0.943	0.485	0.239	0.239	0.493	-0.451	0.880	0.493	0.246	-0.441	0.869	0.495	0.247
Mean Bias	0.154	0.122	0.092	0.142	0.142	0.039	0.105	0.125	0.039	0.069	0.070	0.131	0.021	0.035
RMSE	0.200	0.153	0.117	0.183	0.183	0.051	0.127	0.143	0.051	0.088	0.082	0.137	0.026	0.044
Actual SD of estimates	0.200	0.142	0.116	0.183	0.183	0.050	0.118	0.078	0.050	0.088	0.057	0.039	0.026	0.044
Overconfidence	1.387	1.049	1.128	1.255	1.255	1.178	1.575	1.103	1.178	1.354	1.541	1.102	1.190	1.306

Overconfidence is the standard deviation of the estimated parameter divided by the mean of its estimated standard error.

Table A4. Simulation results for **spatio-temporal** PMLE (500 iterations; y_0^* is estimated by $E(y_t^*)$)

	N=64 & T=4				N=64 & T=16				N=256 & T=16			
	$\beta_0 = -0.5$	$\beta_1 = 1$	γ	ρ	$\beta_0 = -0.5$	$\beta_1 = 1$	γ	ρ	$\beta_0 = -0.5$	$\beta_1 = 1$	γ	ρ
Experiment #1: $\gamma = 0.25, \rho = 0.25$												
Mean Coefficient Estimate	-0.523	1.001	0.249	0.229	-0.484	0.970	0.251	0.249	-0.483	0.964	0.249	0.249
Mean Bias	0.171	0.122	0.080	0.154	0.083	0.057	0.036	0.080	0.044	0.039	0.019	0.042
RMSE	0.237	0.155	0.100	0.202	0.105	0.071	0.046	0.099	0.056	0.047	0.024	0.053
Actual SD of estimates	0.236	0.156	0.100	0.201	0.104	0.064	0.046	0.099	0.053	0.031	0.024	0.053
Overconfidence	1.302	1.104	1.086	1.203	1.192	0.971	1.011	1.166	1.161	0.952	1.027	1.167
Experiment #2: $\gamma = 0.25, \rho = 0.5$												
Mean Coefficient Estimate	-0.614	0.997	0.233	0.459	-0.469	0.930	0.250	0.497	-0.464	0.925	0.247	0.501
Mean Bias	0.318	0.162	0.102	0.160	0.136	0.093	0.046	0.075	0.071	0.077	0.023	0.037
RMSE	0.580	0.206	0.141	0.231	0.176	0.113	0.059	0.099	0.087	0.086	0.029	0.047
Actual SD of estimates	0.570	0.206	0.140	0.228	0.174	0.089	0.059	0.099	0.079	0.042	0.029	0.047
Overconfidence	2.027	1.092	1.371	1.474	1.548	1.067	1.210	1.394	1.392	1.011	1.137	1.283
Experiment #3: $\gamma = 0.5, \rho = 0.25$												
Mean Coefficient Estimate	-0.530	0.937	0.493	0.224	-0.451	0.876	0.496	0.245	-0.438	0.860	0.499	0.246
Mean Bias	0.229	0.169	0.079	0.128	0.143	0.130	0.042	0.079	0.086	0.140	0.022	0.039
RMSE	0.331	0.210	0.104	0.176	0.179	0.149	0.053	0.101	0.103	0.146	0.027	0.049
Actual SD of estimates	0.330	0.200	0.103	0.174	0.173	0.083	0.053	0.101	0.082	0.041	0.027	0.049
Overconfidence	1.464	1.069	1.125	1.354	1.514	1.054	1.142	1.393	1.458	1.058	1.141	1.338

Overconfidence is the standard deviation of the estimated parameter divided by the mean of its estimated standard error.

Table A5. Summary statistics for ρ parameter in common spatial probit estimators from 500 Monte Carlo iterations (100 iterations for RIS with N=1,024).

	N	$\rho = 0$			$\rho = 0.25$			$\rho = 0.5$		
		256	1,024	4,096	256	1,024	4,096	256	1,024	4,096
Bayes	Mean Bias	0.154	0.075	0.038	0.139	0.062	0.032	0.111	0.050	0.025
	RMSE	0.197	0.094	0.048	0.183	0.080	0.040	0.149	0.064	0.030
	Overconfidence	1.002	0.985	1.026	1.062	0.974	1.005	1.011	0.955	0.931
	Non-convergence	0	0	0	0	0	0	0	0	0
GMM	Mean Bias	0.205	0.098		0.185	0.088		0.135	0.070	
	RMSE	0.267	0.122		0.243	0.110		0.188	0.088	
	Overconfidence	1.065	0.960		1.160	1.060		1.009	1.072	
	Non-convergence	0	0		2	0		76	16	
Naive Probit	Mean Bias	0.548	0.278	0.139	0.712	0.626	0.654	1.405	1.465	1.493
	RMSE	0.692	0.345	0.176	0.849	0.694	0.673	1.532	1.496	1.500
	Overconfidence	1.271	1.279	1.311	1.275	1.210	1.219	1.188	1.106	1.087
	Non-convergence	0	0	0	0	0	0	0	0	0
RIS	Mean Bias	0.099	0.046		0.147	0.121		0.255	0.229	
	RMSE	0.129	0.057		0.182	0.131		0.280	0.236	
	Overconfidence	3.141	1.468		5.051	4.213		10.756	12.421	
	Non-convergence	0	0		0	0		0	0	
SPMLE	Mean Bias	0.205	0.097	0.051	0.183	0.086	0.045	0.139	0.069	0.036
	RMSE	0.269	0.122	0.065	0.244	0.108	0.056	0.188	0.089	0.046
	Overconfidence	1.080	0.954	1.007	1.178	1.060	1.090	1.223	1.209	1.251
	Non-convergence	1	1	0	0	1	5	3	6	10

Table A6. Summary statistics for γ parameter for RIS and PMLE estimators from 500 Monte Carlo iterations (100 iterations for RIS with N=64 and T=16).

N x T	$\gamma = 0$				$\gamma = 0.25$				$\gamma = 0.5$			
	64 x 4	64 x 16	256 x 16	64 x 4	64 x 16	256 x 16	64 x 4	64 x 16	256 x 16	64 x 4	64 x 16	256 x 16
RIS												
Mean Bias	0.075	0.032		0.085	0.082		0.162	0.197				
RMSE	0.093	0.039		0.106	0.089		0.178	0.200				
Overconfidence	5.909	1.997		11.607	9.399		22.784	16.475				
No convergence	0	0		0	0		0	0				
PMLE												
Mean Bias	0.086	0.039		0.075	0.035		0.060	0.034				
RMSE	0.106	0.048		0.095	0.044		0.076	0.043				
Overconfidence	1.025	1.008		1.036	0.975		1.048	1.148				
No convergence	0	0		0	0		0	0				

Table A7. Summary statistics for ρ parameter for RIS and PMLE estimators from 500 Monte Carlo iterations (100 iterations for RIS with N=64 and T=16).

N x T	$\gamma = 0.25 \ \& \ \rho = 0.25$				$\gamma = 0.25 \ \& \ \rho = 0.5$				$\gamma = 0.5 \ \& \ \rho = 0.25$			
	64 x 4	64 x 16	256 x 16	64 x 4	64 x 16	256 x 16	64 x 4	64 x 16	256 x 16	64 x 4	64 x 16	256 x 16
RIS												
Mean Bias	0.219	0.168		0.317	0.294		0.275	0.201				
RMSE	0.260	0.187		0.351	0.306		0.299	0.214				
Overconfidence	22.108	19.233		24.361	22.673		30.293	29.351				
No convergence	0	0		0	0		0	0				
PMLE												
Mean Bias	0.154	0.080	0.042	0.160	0.075	0.037	0.128	0.079	0.039			
RMSE	0.202	0.099	0.053	0.231	0.099	0.047	0.176	0.101	0.049			
Overconfidence	1.203	1.166	1.167	1.474	1.394	1.283	1.354	1.393	1.338			
No convergence	2	0	0	3	3	2	1	1	0			

Table A8. Summary statistics for γ parameter for RIS and PMLE estimators from 500 Monte Carlo iterations (100 iterations for RIS with N=64 and T=16).

N x T	$\gamma = 0.25 \& \rho = 0.25$				$\gamma = 0.25 \& \rho = 0.5$				$\gamma = 0.5 \& \rho = 0.25$			
	64 x 4	64 x 16	256 x 16	64 x 4	64 x 16	256 x 16	64 x 4	64 x 16	256 x 16	64 x 4	64 x 16	256 x 16
RIS												
Mean Bias	0.090	0.092		0.096	0.094		0.170	0.232				
RMSE	0.112	0.099		0.122	0.105		0.189	0.236				
Overconfidence	12.465	8.739		16.473	12.526		21.284	17.412				
No convergence	0	0		0	0		0	0				
PMLE												
Mean Bias	0.080	0.036	0.019	0.102	0.046	0.023	0.079	0.042	0.022			
RMSE	0.100	0.046	0.024	0.141	0.059	0.029	0.104	0.053	0.027			
Overconfidence	1.086	1.011	1.027	1.371	1.210	1.137	1.125	1.142	1.141			
No convergence	2	0	0	3	3	2	1	1	0			

Figure A1. Distribution of ρ estimates from Monte Carlo simulations for Bayes, GMM, MLE, RIS, and PMLE estimator.

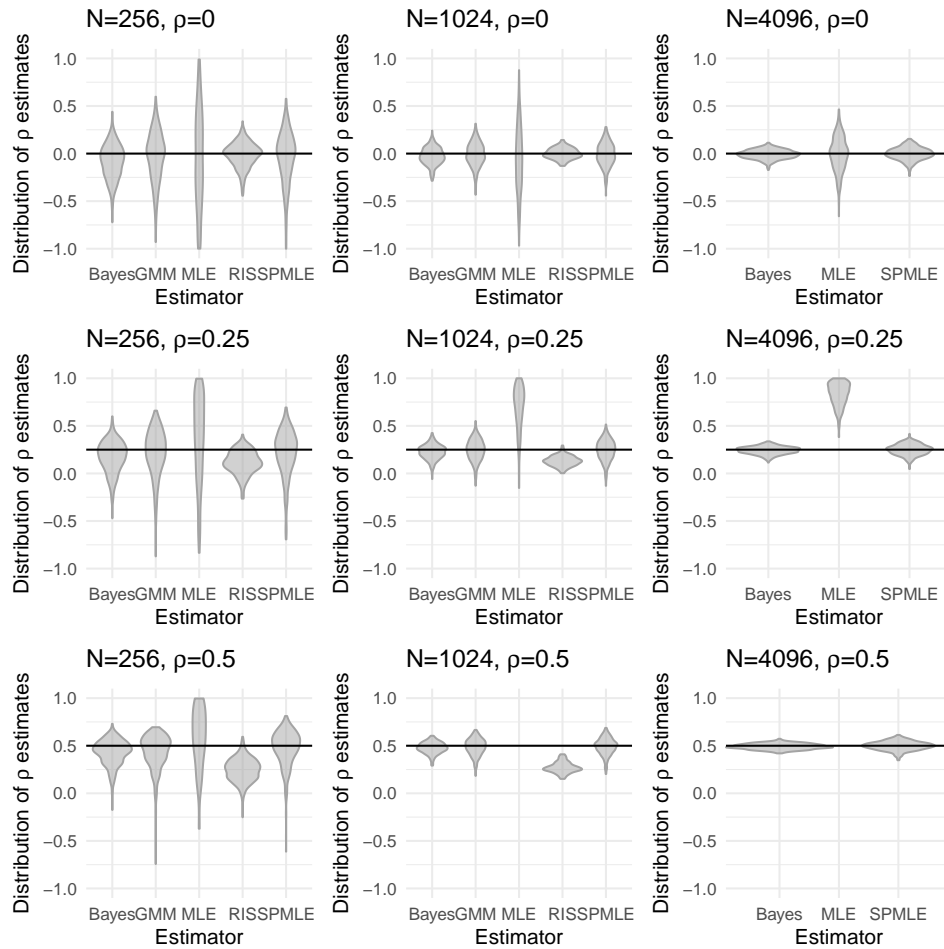


Figure A2. Distribution of γ estimates from Monte Carlo simulations for RIS and PMLE estimator.

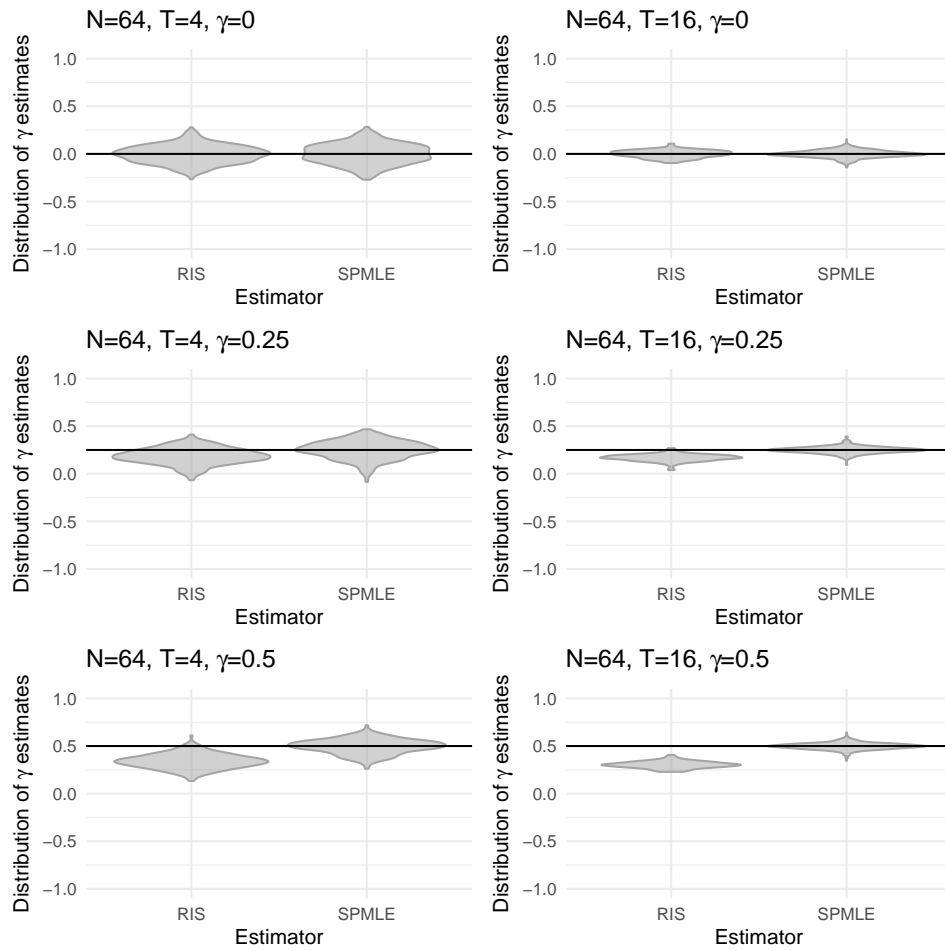
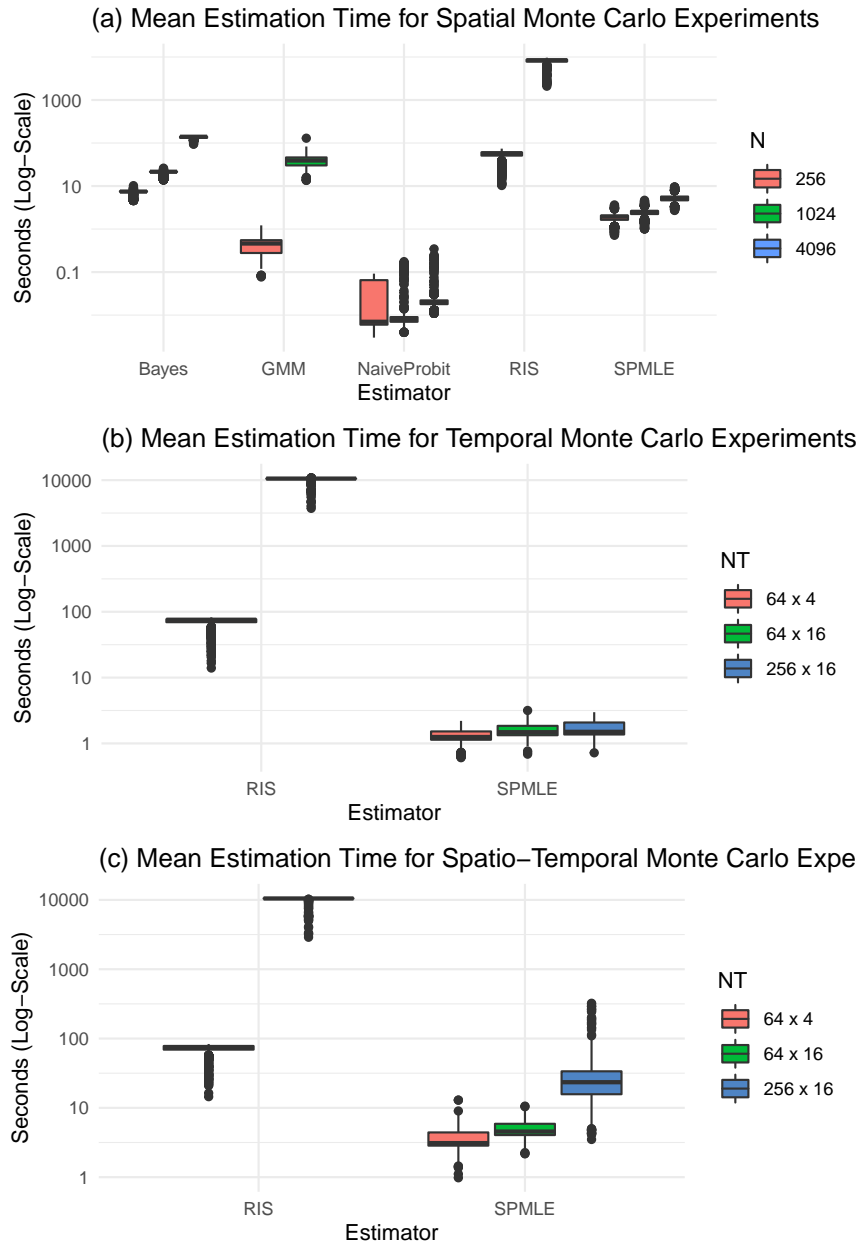


Table A9. Replication of Franzese, Hays, and Cook's (2016) simulation results for spatio-temporal RIS (100 iterations). The experiments here use the DGP (incl. \mathbf{W}) of Franzese et al. (2016). The implementation code is the authors' own—a direct translation of the original MATLAB code to our own R code.

	$\beta_0 = -1.5$	$\beta_1 = 3$	ρ	γ
Experiment #1: $\rho = 0.1, \gamma = 0.3$				
Mean Coefficient Estimate	-1.434	2.859	0.088	0.270
Bias	0.115	0.219	0.037	0.032
RMSE	0.145	0.263	0.045	0.038
Actual SD of estimates	0.130	0.224	0.043	0.022
Experiment #2: $\rho = 0.1, \gamma = 0.5$				
Mean Coefficient Estimate	-1.270	2.487	0.070	0.448
Bias	0.239	0.523	0.060	0.053
RMSE	0.277	0.601	0.107	0.061
Actual SD of estimates	0.156	0.316	0.104	0.033
Experiment #3: $\rho = 0.25, \gamma = 0.3$				
Mean Coefficient Estimate	-1.401	2.798	0.220	0.274
Bias	0.122	0.232	0.042	0.029
RMSE	0.152	0.283	0.051	0.035
Actual SD of estimates	0.116	0.199	0.041	0.023
Experiment #4: $\rho = 0.25, \gamma = 0.5$				
Mean Coefficient Estimate	-1.190	2.374	0.231	0.456
Bias	0.320	0.639	0.046	0.048
RMSE	0.373	0.747	0.064	0.062
Actual SD of estimates	0.209	0.409	0.061	0.044

Figure A3. Mean Estimation Times for Spatial, Temporal and Spatio-Temporal Estimators



References

- Anselin, L. 2002. "Under the Hood: Issues in the Specification and Interpretation of Spatial Regression Models." *Agricultural Economics* 27 (3): 247–267.
- . 2003. "Spatial Externalities, Spatial Multipliers, and Spatial Econometrics." *International Regional Science Review* 26 (2): 153–166.
- Beck, N., D. Epstein, S. Jackman, and S. O'Halloran. 2001. *Alternative Models of Dynamics in Binary Time-Series-Cross-Section Models: The Example of State Failure*.
- Calabrese, R., and J. A. Elkind. 2014a. "Estimating Binary Spatial Autoregressive Models for Rare Events." *Working paper*.
- Calabrese, R., and J. A. Elkind. 2014b. "Estimators of binary spatial autoregressive models: A Monte Carlo study." *Journal of Regional Science* 54 (4): 664–687.
- Erisman, A., and W. Tinney. 1975. "On computing certain elements of the inverse of a sparse matrix." *Communications of the ACM* 18 (3): 177–179.
- Franzese, R. J., J. C. Hays, and S. J. Cook. 2016. "Spatial-and spatiotemporal-autoregressive probit models of interdependent binary outcomes." *Political Science Research and Methods* 4 (1): 151–173.
- Kauppi, H., and P. Saikkonen. 2008. "Predicting US recessions with dynamic binary response models." *The Review of Economics and Statistics* 90 (4): 777–791.
- Kelejian, H. H., and I. R. Prucha. 2010. "Specification and estimation of spatial autoregressive models with autoregressive and heteroskedastic disturbances." *Journal of Econometrics* 157 (1): 53–67.
- LeSage, J., and R. K. Pace. 2009. *Introduction to spatial econometrics*. Chapman / Hall/CRC.
- Smirnov, O. A. 2010. "Modeling spatial discrete choice." *Regional Science and Urban Economics* 40 (5): 292–298.
- Takahashi, K., J. Fagan, and M.-S. Chin. 1973. "Formation of a sparse bus impedance matrix and its application to short circuit study." *Power Industry Computer Applications Conference*: 4–6.
- Van der Vorst, H. A. 1992. "Bi-CGSTAB: A fast and smoothly converging variant of Bi-CG for the solution of nonsymmetric linear systems." *SIAM Journal on scientific and Statistical Computing* 13 (2): 631–644.