

# Shrinkage Estimation of Network Spillovers with Factor Structured Errors<sup>\*</sup>

Ayden Higgins<sup>†</sup>  
University of Surrey, UK

Federico Martellosio<sup>‡</sup>  
University of Surrey, UK

July 5, 2021

## Abstract

This paper explores the estimation of a panel data model with cross-sectional interaction that is flexible both in its approach to specifying the network of connections between cross-sectional units, and in controlling for unobserved heterogeneity. It is assumed that there are different sources of information available on a network, which can be represented in the form of multiple weights matrices. These matrices may reflect observed links, different measures of connectivity, groupings or other network structures, and the number of matrices may be increasing with sample size. A penalised quasi-maximum likelihood estimator is proposed which aims to alleviate the risk of network misspecification by shrinking the coefficients of irrelevant weights matrices to exactly zero. Moreover, controlling for unobserved factors in estimation provides a safeguard against the misspecification that might arise from unobserved heterogeneity. The estimator is shown to be consistent and selection consistent as both  $n$  and  $T$  tend to infinity, and its limiting distribution is characterised. Finite sample performance is assessed by means of a Monte Carlo simulation and the method is applied to study the prevalence of network spillovers in determining growth rates across countries.

**Keywords:** interactive fixed effects, high-dimensional estimation, panel models, penalised quasi-likelihood, social network models.

**JEL classification:** C13, C23, C51.

---

<sup>\*</sup>We are grateful to three anonymous referees and to the Associate Editor for insightful comments that helped us to improve the paper. We would also like to thank Xun Lu and Liangjun Su for sharing their data with us, as well as Valentina Corradi, João Santos Silva and Sorawoot Srisuma for their comments and suggestions.

<sup>†</sup>Email: a.higgins@surrey.ac.uk

<sup>‡</sup>Email: f.martellosio@surrey.ac.uk

# 1 Introduction

Increased attention is being given over to panel data models which take into account cross-sectional interaction. These models have proven to be empirically relevant in a diverse range of economic settings, such as social interactions between individuals, business connections between firms, trading relations between nations, and dependencies between financial assets. At the heart of many econometric models of this kind lies a weights matrix, which summarises the network of connections between interacting cross-sectional units. Yet networks are rarely fully observed, and the uncertainty in how a weights matrix should be specified has been a common critique of this growing literature (see, e.g., Blume et al., 2015; de Paula et al., 2020; Lewbel et al., 2021). In practice, situations in which networks are partially observed are more frequent, with some information being available on cross-sectional links, or their absence, as well as information on other network structures such as groupings. As an example, within a school one might observe family, friendship, classroom and cohort groupings, each of which provide information on the network of connections between different students. In other settings, such as international networks, there are multiple ways to quantify connectivity between nations, including economic measures such as trade and foreign direct investment flows, physical distance, and infrastructure links. Nevertheless, it is not usually obvious how these pieces of the jigsaw fit together, and this uncertainty inevitably increases the risk of model misspecification.

Typical methods to inform the choice of weights matrix include sequential specification testing, or model selection with reference to an information criterion (e.g., Zhang and Yu, 2018). These approaches have largely been focused on the problem of discerning a single best weights matrix from a set of mutually exclusive competitors. In contrast, there are many cases in which weights matrices manifest equally relevant, rather than competing, specifications and, in cases such as these, a model that includes multiple weights matrices may be preferable. This presents a more challenging model selection problem since prospective model specifications may be nested in one another, generating a large number of alternative models. In order to tackle this empirically important issue the current paper uses penalised estimation methods, which retain relevant weights matrix specifications, while at the same time, shrink the coefficients of irrelevant matrices to exactly zero.

A related concern in models of this kind is unobserved heterogeneity. Intuitively, there are likely to be many common factors which are unobserved, and yet have an influence on the outcomes of cross-sectional units; for example exposure to common shocks or a common environment. The presence of common factors can make the identification of model parameters difficult in the event that these are correlated with covariates. The typical

approach in dealing with unobserved heterogeneity is to transform the model in a way which purges the unobserved factors (see, e.g., Yu et al., 2008; Lee and Yu, 2010). Nonetheless, a transform risks purging the very variation needed to identify network spillovers and therefore identification remains a delicate issue, with variation in the regressors, the structure of the weights matrices, and that of the unobserved heterogeneity, each having a part to play. An additional challenge in transforming the model is that prior knowledge on the nature of the unobserved heterogeneity is needed to specify a transform. Traditional examples of this include time, unit and group effect models in which case information on time, unit and group identities is used. Yet with a complex structure of cross-sectional interaction, it is desirable to go beyond these models and to allow for more general forms of heterogeneity. The present paper models a factor structure in the error, which provides this flexibility since common factors may vary across time and have a fully heterogeneous effect on the cross-section. By way of principal component methods, a transform to purge these factors is, in effect, estimated alongside model parameters, removing the reliance on prior knowledge to specify a transform. Taken together, multiple weights matrices, penalisation, and a factor structure error, provide a means of estimating various network spillovers which addresses some of the empirical concerns raised in models of cross-sectional interaction.

The present paper lies in the intersection of several literatures, including social and spatial network models, high-dimensional estimation, and models with factor structured errors. In the social network literature, estimation and identification of network spillovers has been extensively discussed; e.g., Lee (2007), Bramoullé et al. (2009) and Lee et al. (2010). These papers each devote attention to the challenges which may arise in the presence of unobserved heterogeneity, in models where a time dimension is absent. Elsewhere, panel data models which combine interaction and factor structures in the error term have been considered; see, for example, Shi and Lee (2017), Bai and Li (2021) and Kuersteiner and Prucha (2020). In a likelihood framework, Shi and Lee (2017) studies the estimation of a dynamic spatial model with interactive fixed effects and use a single weights matrix to represent dependencies between outcomes. Bai and Li (2021) do similarly, though explicitly allowing for cross-sectional heteroskedasticity. The present paper also pursues likelihood based estimation, and generalises these papers to allow for multiple weights matrices and the possibility that the number may be increasing with sample size. Kuersteiner and Prucha (2020) consider estimation of a model with multiple potentially endogenous weights matrices alongside a factor structure in the error, by way of a method of moments estimator. The approach of Shi and Lee (2017) is partly inspired by Moon and Weidner (2015), who derive the properties of an estimator using an eigenvalue perturbation approach. On the other

hand, Bai and Li (2021) more closely follow Bai (2009), who derives results using first order conditions as a starting point for analysis. In terms of theory, this paper follows the latter approach, and proceeds from first order conditions in similar fashion to Bai (2009).

In the high-dimensional estimation literature, Lu and Su (2016) examine a model with interactive fixed effects and an increasing number of parameters, but without cross-sectional interaction. They make use of the adaptive Lasso penalty of Zou (2006) to induce sparsity amongst both estimated coefficients and factor loadings, assuming that many of these are redundant. Their procedure yields efficiency gains when compared to estimating the model with the number of factors overestimated. The present paper also uses the adaptive Lasso, which penalises the  $\ell_1$  norm of the estimated parameter vector, encouraging sparsity amongst coefficient estimates. High-dimensional spatial models have also been studied elsewhere, such as Lam and Souza (2019), who consider a model which allows for an increasing number of spatial weights matrices, and also use the adaptive Lasso as a penalty, though do not consider unobserved heterogeneity beyond standard fixed effect approaches. Liu (2017) similarly uses penalised estimation in a cross-sectional model with many spatial weights matrices. Gupta and Robinson (2015, 2018) consider estimation of a cross-sectional spatial model, with the number of weights matrices increasing with sample size, by using instrumental variables and quasi-maximum likelihood respectively. The authors carefully study the asymptotic behaviour of these estimators, but do not pursue penalised estimation nor discuss unobserved heterogeneity.

Some recent works have also considered the case where the network is entirely unobserved, such as de Paula et al. (2020) and Lewbel et al. (2021). This situation is especially relevant in the context of social interactions, where connections between individuals might be particularly hard to observe or to quantify. The approach taken in de Paula et al. (2020) involves estimating an entire weights matrix using observations on the same set of individuals across multiple time periods. This can be seen as an extreme case of the current paper in which each weights matrix consists of a single nonzero element taking a value of one. Lewbel et al. (2021) takes a different perspective whereby multiple groups of individuals are observed, a special case of which is when each group consists of the same individuals observed in different time periods. In contrast, the focus of the present paper is where the network is partially observed, which in practice may be quite common. Moreover, establishing identification of the entire weights matrix once a factor term is introduced into the error may be a nontrivial matter.

**Outline:** The model of interest is introduced in Section 2, alongside some basic assumptions and the estimation method. This is followed by asymptotic results in Section 3,

and a discussion on implementation in Section 4. In Section 5 finite sample performance is assessed by means of a small Monte Carlo study, followed by an empirical application of the method to consider whether network spillovers are prevalent in determining growth rates across countries. Section 6 concludes. Proofs of the main results can be found in Appendix A. For further discussion, proofs of lemmas and additional simulation output, see the Supplementary Material.

**Notation:** Throughout the paper, all vectors and matrices are real. For an  $n \times 1$  vector  $\mathbf{b}$  with elements  $b_i$ ,  $\|\mathbf{b}\|_1 := \sum_{i=1}^n |b_i|$ ,  $\|\mathbf{b}\|_2 := \sqrt{\sum_{i=1}^n b_i^2}$ ,  $\|\mathbf{b}\|_\infty := \max_{1 \leq i \leq n} |b_i|$ . Let  $\mathbf{B}$  be an  $n \times m$  matrix with elements  $B_{ij}$ . When  $m = n$ , and the eigenvalues of  $\mathbf{B}$  are real, they are denoted by  $\mu_n(\mathbf{B}) \leq \dots \leq \mu_1(\mathbf{B})$ . The following matrix norms are those induced by their vector counterparts:  $\|\mathbf{B}\|_1 := \max_{1 \leq j \leq m} \sum_{i=1}^n |B_{ij}|$  which is the maximum absolute column sum of  $\mathbf{B}$ ,  $\|\mathbf{B}\|_2 := \sqrt{\mu_1(\mathbf{B}'\mathbf{B})}$ , and  $\|\mathbf{B}\|_\infty := \max_{1 \leq i \leq n} \sum_{j=1}^m |B_{ij}|$  which is the maximum absolute row sum of  $\mathbf{B}$ . The Frobenius norm of  $\mathbf{B}$  is denoted  $\|\mathbf{B}\|_F := \sqrt{\sum_{i=1}^n \sum_{j=1}^m B_{ij}^2} = \sqrt{\text{tr}(\mathbf{B}'\mathbf{B})}$ . Let  $\mathbf{P}_B := \mathbf{B}(\mathbf{B}'\mathbf{B})^+ \mathbf{B}'$  and  $\mathbf{M}_B := \mathbf{I}_n - \mathbf{P}_B$ , where  $\mathbf{I}_m$  is the  $m \times m$  identity matrix and the superscript  $+$  denotes the Moore-Penrose generalised inverse. A sequence of  $n \times n$  matrices  $\mathbf{C}_n$  is said to be uniformly bounded in absolute row and column sums (UB) if both the sequences  $\|\mathbf{C}_n\|_1$  and  $\|\mathbf{C}_n\|_\infty$  are bounded. Throughout  $c$ , potentially indexed when there are many, is used to denote some arbitrary positive constant and ‘w.p.a.1’ indicates ‘with probability approaching 1’.

## 2 Model and Estimation

### 2.1 Model

The model considered in this paper supposes that, amongst  $n$  cross-sectional units in time period  $t = 1, \dots, T$ , outcomes are generated according to

$$\mathbf{y}_t = \sum_{q=1}^{Q_{nT}} \rho_q \mathbf{W}_q \mathbf{y}_t + \sum_{k=1}^{K_{nT}} \beta_k \mathbf{x}_{kt} + \boldsymbol{\eta}_t, \quad (1)$$

where  $\mathbf{y}_t$ ,  $\mathbf{x}_{kt}$  and  $\boldsymbol{\eta}_t$  are  $n \times 1$  vectors of outcomes, covariates and error terms respectively, and  $\mathbf{W}_q$  is an  $n \times n$  weights matrix specified in advance. Both the number  $Q_{nT}$  of potentially relevant weights matrices and the number  $K_{nT}$  of potentially relevant regressors can increase with sample size. The covariates may be subdivided into various types, such that

$$\sum_{k=1}^{K_{nT}} \beta_k \mathbf{x}_{kt} = \sum_{\kappa=1}^{K_{nT}^*} \delta_\kappa \mathbf{x}_{\kappa t}^* + \phi_1 \mathbf{y}_{t-1} + \sum_{q=1}^{Q_{nT}} \phi_{q+1} \mathbf{W}_q \mathbf{y}_{t-1}. \quad (2)$$

The first  $K_{nT}^*$  regressors may be either primitive exogenous covariates, or formed by the interaction of weights matrices and primitive exogenous covariates. It is assumed that there is at least one relevant exogenous covariate, i.e. this paper does not study the case of a pure network autoregression. Moreover, lagged outcomes and the interaction of these with weights matrices can provide additional covariates of the form  $\mathbf{W}_q \mathbf{y}_{t-1}$ . It may be that many of the parameters  $\rho_q$ ,  $\delta_\kappa$  and  $\phi_q$  are truly zero since many of the covariates or weights matrix specifications may be irrelevant. Such restrictions need not be imposed a priori, since penalised estimation induces the estimates of these parameters to take values of exactly zero.

The weights matrices  $\mathbf{W}_q$  contain information about the connections between the cross-sectional units, with larger elements – positive or negative – measuring a stronger connection strength. The literature often assumes that the weights matrices have positive elements and are row normalised such that each of the rows of  $\mathbf{W}_q$  sum to 1. These assumptions lend products of the form  $\mathbf{W}_q \mathbf{b}$  the interpretation of a weighted average of the entries of a vector  $\mathbf{b}$ . While these two assumptions are not necessary in this paper, the assumption that the weights matrices have zero diagonals, which forbids self-links, is retained. The coefficients  $\rho_q$  on  $\mathbf{W}_q \mathbf{y}_t$  capture endogenous spillovers; that is, the impact on the outcome of each unit, generated by the units that are neighbours according to the  $q$ -th weights matrix. Analogously, those  $\delta_\kappa$  coefficients on covariates of the form  $\mathbf{W}_q \mathbf{x}_{\kappa t}^*$  capture exogenous spillovers, also referred to as contextual effects in the social interaction literature. The coefficients  $\phi_{q+1}$  on products  $\mathbf{W}_q \mathbf{y}_{t-1}$  capture dynamic spillovers. Combined, the endogenous, exogenous and dynamic spillovers, allow model (1) to quantify a breadth of different network spillovers.

It is assumed that the error term has a factor structure of the form

$$\boldsymbol{\eta}_t = \boldsymbol{\Lambda} \mathbf{f}_t + \boldsymbol{\varepsilon}_t, \quad (3)$$

where  $\boldsymbol{\Lambda}$  is an  $n \times R$  matrix of time-invariant loadings,  $\mathbf{f}_t$  is an  $R \times 1$  vector of unit-invariant factors, and  $\boldsymbol{\varepsilon}_t$  is an  $n \times 1$  vector of idiosyncratic error terms. In addition, the rows of  $\boldsymbol{\Lambda}$  are denoted by  $\boldsymbol{\lambda}_i$ , for  $i = 1, \dots, n$ , and the factors are arranged in the  $T \times R$  matrix  $\mathbf{F} := (\mathbf{f}_1, \dots, \mathbf{f}_T)'$ . Throughout, the number of factors  $R$  is assumed to be a constant independent of sample size. Following a fixed effects approach, both factors and loadings are treated as (nuisance) parameters in estimation. Thus, in the model, either is allowed to be arbitrarily correlated with covariates. The framework is very general; for instance  $\mathbf{f}_t$  could be aggregate shocks affecting the entire network at time  $t$ , with a heterogeneous effect on each individual. Moreover, this factor structure nests more traditional fixed effect models as special cases.

It is worth stressing that unobserved heterogeneity may arise from various sources. Consider, as a simple example, a model with a single exogenous regressor and no endogenous spillovers, i.e.,

$$\mathbf{y}_t = \beta^* \mathbf{x}_t^* + \sum_{q=1}^{Q_{nT}} \alpha_q \mathbf{W}_q \mathbf{x}_t^* + \delta \mathbf{W}_L \mathbf{x}_t^* + \varepsilon_t, \quad (4)$$

with  $\mathbf{W}_q$  being the  $q$ -th observed weights matrix, and  $\beta^*, \alpha, \delta$  being scalars. Suppose that  $\mathbf{W}_L$  is unobserved and is either low rank or well approximated by a low rank matrix and represents, for example, low rank measurement error in some  $\mathbf{W}_q$ , or unobserved connections between cross-sectional units arising due to network sampling; see, for example, Wang (2018). Defining  $\mathbf{\Lambda}^* \mathbf{f}_t^* := \delta \mathbf{W}_L \mathbf{x}_t^*$ , it is clear that (4) is nested in model (1) and highlights that the decomposition of the unobserved term into factors  $\mathbf{\Lambda}^*$  and loadings  $\mathbf{f}_t^*$  is arbitrary; it is the low rank restriction on  $\delta \mathbf{W}_L \mathbf{x}_t^*$  that allows this term to be distinguished and controlled for.

Going forward, it is convenient to introduce some new notation. The subscript  $nT$  used previously is suppressed from  $Q_{nT}$ ,  $K_{nT}$ ,  $K_{nT}^*$ , and the following parameter vectors and covariate matrices are defined:  $\boldsymbol{\rho} := (\rho_1, \dots, \rho_Q)'$ ,  $\boldsymbol{\delta} := (\delta_1, \dots, \delta_{K^*})'$ ,  $\boldsymbol{\phi} := (\phi_1, \dots, \phi_{Q+1})'$ ,  $\boldsymbol{\beta} := (\beta_1, \dots, \beta_K)' := (\boldsymbol{\delta}', \boldsymbol{\phi}')'$ ,  $\boldsymbol{\theta} := (\boldsymbol{\rho}', \boldsymbol{\beta}')'$ , and  $\mathbf{X}_t := (\mathbf{X}_t^*, \mathbf{y}_{t-1}, \mathbf{W}_1 \mathbf{y}_{t-1}, \dots, \mathbf{W}_Q \mathbf{y}_{t-1})$ , where  $\mathbf{X}_t^* := (\mathbf{x}_{1t}^*, \dots, \mathbf{x}_{K^*t}^*)$ , and  $\mathbf{S}(\boldsymbol{\rho}) := \mathbf{I}_n - \sum_{q=1}^Q \rho_q \mathbf{W}_q$ . Given these, model (1) can be restated more succinctly as

$$\mathbf{S}(\boldsymbol{\rho}) \mathbf{y}_t = \mathbf{X}_t \boldsymbol{\beta} + \mathbf{\Lambda} \mathbf{f}_t + \varepsilon_t. \quad (5)$$

Throughout, the superscript 0 is used to distinguish the true values of the factors, loadings, and parameters, as well as the true numbers of these, and the framework is one in which  $n$  and  $T$  diverge simultaneously. The total number of parameters in the vector  $\boldsymbol{\theta}$  is  $P := Q + K$ , of which only  $P^0$  are truly nonzero. In fact, one might often expect that the vector  $\boldsymbol{\theta}$  is sparse in the sense that many of its components are zero, particularly in cases with a large number of weights matrices and covariates. Accordingly  $\boldsymbol{\theta}$  may be reordered as  $\boldsymbol{\vartheta} = (\boldsymbol{\theta}'_{(1)}, \boldsymbol{\theta}'_{(2)})'$ , where  $\boldsymbol{\theta}_{(1)}$  is the  $P^0 \times 1$  vector of nonzero parameters, and  $\boldsymbol{\theta}_{(2)}^0 = \mathbf{0}_{(P-P^0) \times 1}$ . Sparsity, however, is not necessary and indeed the results of this paper equally allow for the possibility that all of the weights matrices and covariates may be relevant. The  $n \times T$  data matrix for the  $\kappa$ -th exogenous covariate is denoted  $\mathbf{X}_\kappa^* := (\mathbf{x}_{\kappa 1}^*, \dots, \mathbf{x}_{\kappa T}^*)$  for  $\kappa = 1, \dots, K^*$ , and the  $n \times T$  data matrix for the lagged outcomes is denoted  $\mathbf{Y}_{-1} := (\mathbf{y}_{t-1}, \dots, \mathbf{y}_{t-T})$ . The data matrix for the generic  $k$ -th covariate of any type,  $\mathbf{X}_\kappa^*$ ,  $\mathbf{Y}_{-1}$  or  $\mathbf{W}_q \mathbf{Y}_{-1}$ , is denoted  $\mathbf{X}_k := (\mathbf{x}_{k1}, \dots, \mathbf{x}_{kT})$ , for  $k = 1, \dots, K$ . Also,  $\mathbf{A}(\boldsymbol{\rho}, \boldsymbol{\phi}) := \mathbf{S}^{-1}(\boldsymbol{\rho})(\phi_1 \mathbf{I}_n + \sum_{q=1}^Q \phi_{q+1} \mathbf{W}_q)$ ,  $\mathbf{A} := \mathbf{A}(\boldsymbol{\rho}^0, \boldsymbol{\phi}^0)$ ,  $\mathbf{S} := \mathbf{S}(\boldsymbol{\rho}^0)$ ,  $\mathbf{G}_q(\boldsymbol{\rho}) := \mathbf{W}_q \mathbf{S}^{-1}(\boldsymbol{\rho})$ , and  $\mathbf{G}_q := \mathbf{G}_q(\boldsymbol{\rho}^0)$ .

## 2.2 Assumptions

The first set of assumptions concerns the idiosyncratic error term  $\varepsilon_{it}$ .

### Assumption 1.

- 1.1 *The errors  $\varepsilon_{it}$  are identically and independently distributed over  $i$  and  $t$  with  $\mathbb{E}[\varepsilon_{it}] = 0$ ,  $\mathbb{E}[\varepsilon_{it}^2] = \sigma_0^2 \geq c > 0$ , and fourth moments uniformly bounded over  $i$  and  $t$ .*
- 1.2 *The errors  $\varepsilon_{it}$  are independent of the elements of the matrices  $\mathbf{\Lambda}^0$ ,  $\mathbf{F}^0$ , and  $\mathbf{X}_\kappa^*$ , for  $\kappa = 1, \dots, K^*$ .*

These assumptions have been employed across various papers. Cross-sectional homoskedasticity and independence is commonly assumed, though this can be relaxed by estimation of a more general  $n \times n$  covariance matrix  $\mathbf{\Sigma}^0$ , at the expense of additional parameters; see for example Bai and Liao (2017) and Bai and Li (2021). Additional structure in the error term could also be considered as is commonplace throughout the spatial econometrics literature. Yet since the factor structure provides a mechanism for capturing such correlation, Assumption 1.1 assumes  $\mathbf{\Sigma}^0 = \sigma_0^2 \mathbf{I}_n$ . Differing assumptions concerning the relationship between the errors, the factors, and the loadings appear across the literature; these are comprehensively surveyed by Hsiao (2018). Assumption 1.2 imposes independence of the factors and the loadings from the error term as in Bai (2009).

Some additional assumptions are required regarding the other components of the model. Let  $|\cdot|$  denote the entrywise absolute value of a vector or matrix,  $\mathbf{\Theta}$  denote the parameter space for  $\boldsymbol{\theta}$ , and  $\mathbf{\Theta}_\rho$  and  $\mathbf{\Theta}_\phi$  denote the parameter spaces for  $\boldsymbol{\rho}$  and  $\boldsymbol{\phi}$ , respectively. Since the matrices  $\mathbf{W}_1, \dots, \mathbf{W}_Q$  depend on  $n$ , as well as the number of weights matrices  $Q$ , the matrices  $\mathbf{S}(\boldsymbol{\rho})$ ,  $\mathbf{A}(\boldsymbol{\rho}, \boldsymbol{\phi})$ , as well as the parameter spaces, depend on  $n$  and  $T$ . Therefore, it is understood that the following set of assumptions are to hold for any  $(n, T)$ .

### Assumption 2.

- 2.1 *The parameter vector  $\boldsymbol{\theta}^0$  is in the interior of  $\mathbf{\Theta}$ , with  $\mathbf{\Theta}$  being a compact subset of  $\mathbb{R}^P$ .*
- 2.2 *The weights matrices  $\mathbf{W}_1, \dots, \mathbf{W}_Q$  are nonstochastic and UB uniformly over  $q$ .*
- 2.3 *For all  $\boldsymbol{\rho} \in \mathbf{\Theta}_\rho$  and  $\boldsymbol{\phi} \in \mathbf{\Theta}_\phi$ ,  $\mathbf{S}(\boldsymbol{\rho})$  is invertible,  $\mathbf{S}(\boldsymbol{\rho}), \mathbf{S}^{-1}(\boldsymbol{\rho})$  and  $\sum_{h=1}^{\infty} |\mathbf{A}^h(\boldsymbol{\rho}, \boldsymbol{\phi})|$  are UB,  $\|\mathbf{A}(\boldsymbol{\rho}, \boldsymbol{\phi})\|_2 < 1 - c$  for some  $c > 0$ , and  $\liminf_{n, T \rightarrow \infty} \inf_{\boldsymbol{\rho} \in \mathbf{\Theta}_\rho} \det(\mathbf{S}(\boldsymbol{\rho})) \neq 0$ .*



2.4 The elements of the matrices  $\mathbf{X}_\kappa^*$  have fourth moments uniformly bounded over  $i, t$  and  $\kappa$ , and elements of the matrix  $\sum_{k=1}^K \beta_k^0 \mathbf{X}_k$  have fourth moments uniformly bounded over  $i$  and  $t$ .

2.5 The true number of factors  $R^0$  is constant.

2.6 The elements of the matrices  $\mathbf{F}^0$  and  $\mathbf{\Lambda}^0$  have eighth moments uniformly bounded over  $i$  and  $t$ .

Assumption 2.1 considers a sequence of compact parameter spaces over which to maximise the objective function. The condition in Assumption 2.2 that the weights matrices are UB is standard and serves to limit interactions to a manageable degree. Here, uniform boundedness over  $q$  is also required, due to the possibility that  $Q$  increases with sample size. Assumption 2.3 ensures that the model admits a reduced form, and the dynamic process is stationary. Restrictions on the parameter space of  $\boldsymbol{\rho}$  which ensure that  $\mathbf{S}(\boldsymbol{\rho})$  is invertible have been discussed elsewhere in the literature, particularly in the case  $Q = 1$ . A general condition sufficient for the invertibility of  $\mathbf{S}(\boldsymbol{\rho})$  is  $\|\sum_{q=1}^Q \rho_q \mathbf{W}_q\| < 1$  for some matrix norm  $\|\cdot\|$ , though with  $Q > 1$  more informative conditions can be difficult to obtain outside of exceptional cases.<sup>1</sup> However, as noted by Gupta and Robinson (2018), even when it is possible to characterise inadmissible values of  $\boldsymbol{\rho}$  and exclude these, the resulting parameter space is unlikely to be compact. It is therefore commonplace in the literature to restrict attention to a region around the origin in which  $\mathbf{S}(\boldsymbol{\rho})$  can be guaranteed to be invertible. This is where  $\sum_{q=1}^Q |\rho_q| < (\max_{1 \leq q \leq Q} \|\mathbf{W}_q\|)^{-1}$ .<sup>2</sup> Yet while the set of  $\boldsymbol{\rho}$  which satisfy this is bounded, it is also open. Therefore to ensure the existence of a maximiser over this space, a closed subset can be considered such that  $\sum_{q=1}^Q |\rho_q| \leq (1 - \tau)(\max_{1 \leq q \leq Q} \|\mathbf{W}_q\|)^{-1}$ , with  $\tau \in (0, 1)$ . Row normalisation of the matrices  $\mathbf{W}_q$  further simplifies this condition since it implies  $\max_{1 \leq q \leq Q} \|\mathbf{W}_q\|_\infty = 1$ . Model (5) can be rewritten as  $\mathbf{y}_t = \mathbf{A}\mathbf{y}_{t-1} + \mathbf{S}^{-1}(\mathbf{X}_t^* \boldsymbol{\delta} + \mathbf{\Lambda} \mathbf{f}_t + \boldsymbol{\varepsilon}_t)$ , or, after recursive substitution,  $\mathbf{y}_t = \sum_{h=0}^{\infty} \mathbf{A}^h \mathbf{S}^{-1}(\mathbf{X}_{t-h}^* \boldsymbol{\delta} + \mathbf{\Lambda} \mathbf{f}_{t-h} + \boldsymbol{\varepsilon}_{t-h})$ ; Assumption 2.3 guarantees that this series converges. Further discussion of parameter restrictions ensuring convergence of this series can be found in Lee and Yu (2014) and Shi and Lee (2017).<sup>3</sup> The first part of Assumption 2.4 ensures that  $\|\mathbf{X}_k^*\|_F = O_P(\sqrt{nT})$ , for  $k = 1, \dots, K$ . For the

<sup>1</sup>One such case is where the matrices  $\mathbf{W}_1, \dots, \mathbf{W}_Q$  are simultaneously diagonalisable, for example where they consist of powers of a single weights matrix. Another example is where the weights matrices consist of nonoverlapping blocks.

<sup>2</sup>This inequality is obtained from the condition  $\|\sum_{q=1}^Q \rho_q \mathbf{W}_q\| < 1$  and the fact that  $\|\sum_{q=1}^Q \rho_q \mathbf{W}_q\| \leq \sum_{q=1}^Q |\rho_q| \max_{1 \leq q \leq Q} \|\mathbf{W}_q\|$  for any matrix norm  $\|\cdot\|$ .

<sup>3</sup>For example, where the weights matrices consist of nonoverlapping blocks,  $\sum_{q=1}^Q |\rho_q| + \sum_{q=1}^{Q+1} |\phi_q| < 1$  is sufficient for  $\|\mathbf{A}(\boldsymbol{\rho}, \boldsymbol{\phi})\|_2 < 1$ .

second part, notice that  $\mathbf{G}_q \mathbf{X}_t \boldsymbol{\beta}^0$  can be used as an instrument in the estimation of  $\rho_q$ .<sup>4</sup> With a diverging number of parameters, the second part of Assumption 2.4 assures that for these instruments  $\|\sum_{k=1}^K \beta_k^0 \mathbf{G}_q \mathbf{X}_k\|_F = O_P(\sqrt{nT})$ . An alternative condition sufficient for this is  $\|\boldsymbol{\beta}^0\|_1 < c$ , which follows by Hölder's inequality or, alternatively, Assumption 2.4 could be replaced by one restricting the growth of  $K^0$  and  $n, T$ . Assumption 2.5 is common throughout the literature, but could be relaxed at the expense of slower rates of convergence. Several differing assumptions concerning the moments of the factors and the loadings appear in the literature. Given the possible presence of lagged outcomes as covariates, Assumption 2.6 serves the same purpose as Assumption 5(vi) in Moon and Weidner (2017), and ensures that  $y_{it}$  has uniformly bounded fourth moments.

### 2.3 Objective Function

The estimation strategy employed in this paper is penalised quasi-maximum likelihood (PQML), using the multivariate standard normal distribution for the error term, i.e.,  $\varepsilon_{it} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma_0^2)$ , and following a fixed effects approach. Maximum likelihood estimation is a standard in the literature for models of this type, since the simultaneity in the determination of outcomes generates an endogeneity problem which results in least squares estimates being biased. The parameter of interest is  $\boldsymbol{\theta}$ , whereas  $\boldsymbol{\Lambda}, \mathbf{F}, \sigma^2$  are treated as nuisance parameters. Since fixing  $\boldsymbol{\theta}$  results in a pure factor model (and  $\boldsymbol{\Lambda}, \mathbf{F}, \sigma^2$  are not penalised), the estimators of  $\boldsymbol{\Lambda}$  and  $\mathbf{F}$  for fixed  $\boldsymbol{\theta}$  are a solution to a standard principal component problem (see, e.g., Bai, 2009). In this subsection  $R$  is fixed such that  $R \geq R^0$ ; this is discussed in greater detail in Section 3.1. With  $R$  fixed, the average (quasi) log-likelihood is

$$\begin{aligned} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\Lambda}, \mathbf{F}, \sigma^2) := & -\frac{1}{2} \log(2\pi) + \frac{1}{n} \log(\det(\mathbf{S}(\boldsymbol{\rho}))) - \frac{1}{2} \log(\sigma^2) \\ & - \frac{1}{2\sigma^2} \frac{1}{nT} \sum_{t=1}^T (\mathbf{S}(\boldsymbol{\rho}) \mathbf{y}_t - \mathbf{X}_t \boldsymbol{\beta} - \boldsymbol{\Lambda} \mathbf{f}_t)' (\mathbf{S}(\boldsymbol{\rho}) \mathbf{y}_t - \mathbf{X}_t \boldsymbol{\beta} - \boldsymbol{\Lambda} \mathbf{f}_t) \end{aligned} \quad (6)$$

and its penalised counterpart is

$$\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\Lambda}, \mathbf{F}, \sigma^2) := \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\Lambda}, \mathbf{F}, \sigma^2) - \varrho(\boldsymbol{\theta}, \boldsymbol{\gamma}, \zeta), \quad (7)$$

where  $\varrho(\boldsymbol{\theta}, \boldsymbol{\gamma}, \zeta)$  is a penalty function and  $\boldsymbol{\gamma}, \zeta$  are regularisation parameters. The specific form of penalty function is introduced in Section 2.4, and the choice of regularisation parameters is discussed in Section 4.1, however for the moment these are both also taken to

---

<sup>4</sup>Observing that  $\mathbf{S}^{-1} = \mathbf{I}_n + \sum_{q=1}^Q \rho_q^0 \mathbf{G}_q$ , then  $\mathbf{y}_t = \mathbf{X}_t \boldsymbol{\beta}^0 + \sum_{q=1}^Q \rho_q^0 \mathbf{G}_q \mathbf{X}_t \boldsymbol{\beta}^0 + \mathbf{S}^{-1} \boldsymbol{\Lambda}^0 \mathbf{f}_t^0 + \mathbf{S}^{-1} \boldsymbol{\varepsilon}_t$ , which makes the role of  $\mathbf{G}_q \mathbf{X}_t \boldsymbol{\beta}^0$  as an instrument for  $\mathbf{W}_q \mathbf{y}_t$  transparent.

be fixed alongside the number of factors. Concentrating out  $\sigma^2$ , as well as the factors, and dropping the constant in (7) yields the concentrated expression

$$\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\Lambda}) = \frac{1}{n} \log(\det(\mathbf{S}(\boldsymbol{\rho}))) - \frac{1}{2} \log(\hat{\sigma}^2(\boldsymbol{\theta}, \boldsymbol{\Lambda})) - \varrho(\boldsymbol{\theta}, \boldsymbol{\gamma}, \zeta), \quad (8)$$

where  $\hat{\sigma}^2(\boldsymbol{\theta}, \boldsymbol{\Lambda}) := \frac{1}{nT} \sum_{t=1}^T \mathbf{e}_t' \mathbf{M}_{\boldsymbol{\Lambda}} \mathbf{e}_t$  and  $\mathbf{e}_t := \mathbf{S}(\boldsymbol{\rho}) \mathbf{y}_t - \mathbf{X}_t \boldsymbol{\beta}$ . Hereafter the terms likelihood and log-likelihood are used synonymously. In order to maximise (8) with respect to  $\boldsymbol{\Lambda}$ , note that

$$\begin{aligned} \min_{\boldsymbol{\Lambda} \in \mathbb{R}^{n \times R}} \frac{1}{nT} \sum_{t=1}^T \mathbf{e}_t' \mathbf{M}_{\boldsymbol{\Lambda}} \mathbf{e}_t &= \frac{1}{nT} \sum_{t=1}^T \mathbf{e}_t' \mathbf{e}_t - \max_{\boldsymbol{\Lambda} \in \mathbb{R}^{n \times R}} \frac{1}{nT} \sum_{t=1}^T \mathbf{e}_t' \mathbf{P}_{\boldsymbol{\Lambda}} \mathbf{e}_t \\ &= \text{tr} \left( \frac{1}{nT} \sum_{t=1}^T \mathbf{e}_t \mathbf{e}_t' \right) - \max_{\boldsymbol{\Lambda} \in \mathbb{R}^{n \times R}: \boldsymbol{\Lambda}' \boldsymbol{\Lambda} = \mathbf{I}_R} \text{tr} \left( \frac{1}{nT} \sum_{t=1}^T \boldsymbol{\Lambda}' \mathbf{e}_t \mathbf{e}_t' \boldsymbol{\Lambda} \right) \\ &= \text{tr} \left( \frac{1}{nT} \sum_{t=1}^T \mathbf{e}_t \mathbf{e}_t' \right) - \sum_{r=1}^R \mu_r \left( \frac{1}{nT} \sum_{t=1}^T \mathbf{e}_t \mathbf{e}_t' \right), \end{aligned} \quad (9)$$

where the second line follows from the fact that any orthogonal projector  $\mathbf{P}_{\boldsymbol{\Lambda}}$  can be written as  $\boldsymbol{\Lambda}' \boldsymbol{\Lambda}$ , with the columns of  $\boldsymbol{\Lambda}$  forming an orthonormal basis for the column space of  $\boldsymbol{\Lambda}$ , and the third line follows from a standard result (e.g., Horn and Johnson, 2012, Corollary 4.3.39).<sup>5</sup> Hence, (9) can be used to concentrate out  $\boldsymbol{\Lambda}$  in (8), whereby the PQML estimator of  $\boldsymbol{\theta}^0$  is characterised as

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta} \in \boldsymbol{\Theta}} \mathcal{Q}(\boldsymbol{\theta}), \quad (10)$$

where

$$\mathcal{Q}(\boldsymbol{\theta}) = \frac{1}{n} \log(\det(\mathbf{S}(\boldsymbol{\rho}))) - \frac{1}{2} \log \left( \sum_{i=R+1}^n \mu_i \left( \frac{1}{nT} \sum_{t=1}^T \mathbf{e}_t \mathbf{e}_t' \right) \right) - \varrho(\boldsymbol{\theta}, \boldsymbol{\gamma}, \zeta). \quad (11)$$

Here it is worth highlighting that both the factors and the loadings have been concentrated out without imposing any of the normalisations typically encountered in the wider factor literature. This is due to the treatment of both the factors and the loadings as nuisance parameters, in which case only the space spanned by the loadings implicitly features in the objective function (11). It would, of course, be possible to consider estimators of the factors and the loadings, however the same fundamental indeterminacy issue would

---

<sup>5</sup>For example, by the QR decomposition,  $\mathbf{B} = \mathbf{V}_B \mathbf{R}$  with  $\mathbf{V}_B \in \mathbb{R}^{n \times m}$  having orthonormal columns and  $\mathbf{R} \in \mathbb{R}^{m \times m}$  being upper triangular. Since  $\mathbf{B}$  has full column rank  $\mathbf{R}$  is invertible (e.g., Horn and Johnson, 2012, Theorem 2.1.14) and therefore  $\mathbf{P}_B := \mathbf{B}(\mathbf{B}'\mathbf{B})^{-1}\mathbf{B} = \mathbf{V}_B \mathbf{V}_B'$ .

arise in separating these as is encountered elsewhere in the factor literature, and therefore some normalisations would typically be required in order to do this. It should also be pointed out that neither the concentrated likelihood nor the penalised objective function  $Q(\boldsymbol{\theta})$  are concave in  $\boldsymbol{\theta}$ . Although subsequent sections establish the desirable asymptotic properties of global maximisers of these objective functions, it is nonetheless the case that local maximisers which do not possess these properties may indeed exist.

## 2.4 Penalty

The present paper adopts the adaptive Lasso, which induces sparsity in parameter estimation by augmenting an objective function with a constraint on the  $\ell_1$  norm of the estimated parameter vector. A desirable feature of this method of penalisation is that it can achieve the oracle property; that is, perform consistent variable selection and, at the same time, possess an optimal rate of convergence. This is done by using an initial consistent estimator of the parameters to weight the penalty. The cost of this is the need to find an initial consistent estimator, which can be difficult in settings where the number of parameters is greater than the number of observations ( $nT < P$  in the present case). This complication is not considered in this paper and attention is restricted to the  $nT > P$  setting. Explicitly, the penalty function employed in this paper has the additive form

$$\varrho(\boldsymbol{\theta}, \boldsymbol{\gamma}, \zeta) := \gamma_\rho \sum_{q=1}^Q \omega_q |\rho_q| + \gamma_\beta \sum_{k=1}^K \omega_{Q+k} |\beta_k|, \quad (12)$$

where  $\omega_p := |\theta_p^\dagger|^{-\zeta}$ , with  $\theta_p^\dagger$  being an initial consistent estimate of the  $p$ -th parameter, and  $\boldsymbol{\gamma} := (\gamma_\rho, \gamma_\beta)'$  and  $\zeta$  are regularisation parameters.<sup>6</sup> The form of the penalty term in (12) allows the penalty parameters  $\gamma_\rho$  and  $\gamma_\beta$  to differ across the two types of parameter,  $\rho_q$  and  $\beta_k$ . In general the penalty term can be easily modified to allow for a greater or lesser degree of heterogeneity, as applications dictate.

The parameter  $\zeta$  is a positive constant and is used to adjust the weight of penalisation according to the rate of consistency of the initial estimator. Combined,  $\zeta$  and  $\theta_p^\dagger$  generate bespoke weights  $\omega_p$  for each parameter that will increase for truly zero coefficients and tend to a constant for truly nonzero coefficients. The other penalty parameters  $\gamma_\rho$  and  $\gamma_\beta$  are positive sequences which tend towards zero as  $n$  and  $T$  increase. Let  $\underline{\theta}^0$  and  $\bar{\theta}^0$  denote, respectively, the minimum and maximum element of  $|\boldsymbol{\theta}_{(1)}^0|$ . Note that both  $\underline{\theta}^0$  and  $\bar{\theta}^0$  can vary with sample size due to the increasing dimension of  $\boldsymbol{\theta}_{(1)}^0$ . The following are assumed.

---

<sup>6</sup>If  $\theta_p^\dagger = 0$  then  $\omega_p$  is set equal to  $\infty$ .

**Assumption 3.**

$$3.1 \quad 0 < c_1 \leq \underline{\theta}^0 \leq \bar{\theta}^0 \leq c_2 < \infty.$$

$$3.2 \quad \max\{\gamma_\rho, \gamma_\beta\} \min\{n, T\} = O(1).$$

$$3.3 \quad \|\boldsymbol{\theta}^\dagger - \boldsymbol{\theta}^0\|_2 = O_P(c_{nT}), \text{ for some sequence } c_{nT} \rightarrow 0 \text{ as } n, T \rightarrow \infty.$$

In this paper it is assumed that, while the dimension of  $\boldsymbol{\theta}^0$  may be increasing with sample size, the value of each element is fixed. Nonetheless, this does not rule out either the minimum or maximum (in absolute value) nonzero elements in  $\boldsymbol{\theta}^0$  becoming arbitrarily small or large as its dimension increases, and therefore Assumption 3.1 imposes that the nonzero elements in  $\boldsymbol{\theta}^0$  are uniformly bounded away from zero and from infinity. Assumption 3.2 requires the penalty parameters  $\gamma_\rho$  and  $\gamma_\beta$  to converge to zero sufficiently fast that they do not adversely impact the rate of consistency of the estimator. Assumption 3.3 requires consistency of the initial estimator  $\boldsymbol{\theta}^\dagger$  at some rate  $c_{nT}$ . If the speed at which  $c_{nT} \rightarrow 0$  is especially slow, then  $\zeta$  can be adjusted to compensate for this. In the following it is shown that the unpenalised likelihood can be used to produce a initial consistent estimator, though other estimation procedures might equally be considered.

In principle it would also be possible to obtain several of the results in this paper under a ‘moving parameter’ framework, where the values of the nonzero elements in  $\boldsymbol{\theta}^0$  might vary with sample size; in particular, where some may converge to zero asymptotically. However, the rate at which they could be allowed do so would need to be sufficiently slow that a choice of  $\gamma_\rho$  and  $\gamma_\beta$  could still be made to ensure the consistency and model selection consistency of the procedure. Moreover, in Section 3.3 the assumption that the nonzero elements in  $\boldsymbol{\theta}^0$  are fixed is important for the validity of the asymptotic distribution derived in that section. Therefore, the assumption that true parameters are fixed is maintained throughout this paper.

### 3 Asymptotic Results

#### 3.1 Consistency

Mirroring Bai (2009), in this section a preliminary consistency result is established which will be improved upon later. Yet, before proceeding, it is worth providing a few remarks on the identification of model parameters. In the standard consistency argument for an extremum estimator, the essence of the idea is to show that *“the limit of the maximum  $\hat{\boldsymbol{\theta}}$  should be the maximum of the limit”*, with the latter being unique (Newey and McFadden,

1994, p. 2120). In that argument the role that identification plays is transparent, and with identification established, uniform convergence of the sample objective function to the limiting objective function often then appeals to a uniform law of large numbers, and consistency follows thereafter. Yet in models where the number of parameters, nuisance or otherwise, depends on the sample size, there is no fixed population distribution from which a sample is drawn, and therefore uniform convergence must be considered more carefully. In cases such as these, consistency is often shown directly, forgoing an explicit identification result. For these same reasons this paper also proceeds directly to consistency, with further discussion on identification being available in Appendix B of the Supplementary Material.

Before formulating the next assumption, it is necessary to introduce some additional notation. Define the  $n \times P$  matrix of instruments  $\mathbf{Z}_t := (\mathbf{G}_1 \mathbf{X}_t \boldsymbol{\beta}^0, \dots, \mathbf{G}_Q \mathbf{X}_t \boldsymbol{\beta}^0, \mathbf{X}_t)$ . The  $n \times T$  data matrix for the instrument associated with some  $\rho_q$  is  $\sum_{k=1}^K \beta_k^0 \mathbf{G}_q \boldsymbol{\chi}_k$ . The generic  $n \times T$  data matrix of either type,  $\boldsymbol{\chi}_k$  or  $\sum_{k=1}^K \beta_k^0 \mathbf{G}_q \boldsymbol{\chi}_k$ , is denoted  $\mathbf{Z}_p := (z_{p1}, \dots, z_{pT})$ , where  $z_{pt}$  is the  $p$ -th column of  $\mathbf{Z}_t$ , for  $p = 1, \dots, P$ . Finally, let  $\mathcal{H}_1(\boldsymbol{\Lambda}, \mathbf{F}) := \frac{1}{nT} \mathbf{Z}' (\mathbf{M}_F \otimes \mathbf{M}_\Lambda) \mathbf{Z}$  and  $\mathcal{H}_2 := \frac{1}{nT} \mathbf{Z}' \mathbf{Z}$ , where  $\mathbf{Z} := (\mathbf{Z}'_1, \dots, \mathbf{Z}'_T)$  is a  $P \times nT$  matrix.

**Assumption 4.**

4.1  $R \geq R^0$ .

4.2  $\inf_{\boldsymbol{\Lambda} \in \mathbb{R}^{n \times R}, \mathbf{F} \in \mathbb{R}^{T \times R^0}} \mu_P(\mathcal{H}_1(\boldsymbol{\Lambda}, \mathbf{F})) \geq c_1 > 0$  w.p.a.1 as  $n, T \rightarrow \infty$ .

4.3  $\mu_1(\mathcal{H}_2) \leq c_2 < \infty$  w.p.a.1 as  $n, T \rightarrow \infty$ .

4.4  $\frac{P}{\min\{n, T\}} \rightarrow 0$ .

Assumption 4.1 allows for the number of factors  $R^0$  to be unknown, as long as the number of factors  $R$  used in estimation is no less than  $R^0$ ; see Moon and Weidner (2015). Assumption 4.2 demands a certain level of variation in sample data after projecting out arbitrary factors and factor loadings. This condition can intuitively be understood by considering the particular case of individual or time effects, in which case the projections perform between individual and between time period differences to the data. It is also worth noting that Assumptions 4.2 and 4.3 imply that, w.p.a.1,

$$\sup_{\boldsymbol{\Lambda} \in \mathbb{R}^{n \times R}, \mathbf{F} \in \mathbb{R}^{T \times R^0}} \mu_1(\mathcal{H}_1(\boldsymbol{\Lambda}, \mathbf{F})) \leq c_2 < \infty \quad (13)$$

and

$$\mu_P(\mathcal{H}_2) \geq c_1 > 0, \quad (14)$$

which ensures both  $\mathcal{H}_1$  and  $\mathcal{H}_2$  are well defined asymptotically (see Appendix C in the Supplementary Material for details). Assumption 4.4 requires that the number of parameters does not grow too fast in relation to  $n$  and  $T$ . This is necessary since consistency is stated in terms of the  $\ell_2$  norm of a vector with increasing dimension. Unfettered growth in the number of parameters relative to the sample size could thus lead to inconsistency even in the event that an estimator converged pointwise. Recall that  $\hat{\boldsymbol{\theta}}$  denotes the maximiser of the penalised likelihood function and let  $\tilde{\boldsymbol{\theta}}$  denote the maximiser of the unpenalised likelihood function.

**Proposition 1** (Consistency). *Under Assumptions 1–4,  $\|\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^0\|_2 = O_P(a_{nT})$  and  $\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^0\|_2 = O_P(a_{nT})$ , where  $a_{nT} := \sqrt{\frac{P}{\min\{n, T\}}}$ .*

This preliminary result is an important step towards those which follow. Moreover, the result is of interest in and of itself since it applies provided that the number of factors is not underspecified, and irrespective of the relationship between  $n$  and  $T$ , as long as both diverge to infinity. In contrast later in the paper, it will be required that the true number of factors is known, and that  $n$  and  $T$  grow in proportion (see Assumption 6). Despite both the factors and the loadings having been concentrated out, the spaces spanned by both are implicitly estimated by their respective first order conditions, and as a result both  $n$  and  $T$  are required to diverge. The rate  $a_{nT}$  is in line with the existing literature; see for example Theorem 4.1 in Moon and Weidner (2015), where a preliminary  $\sqrt{\min\{n, T\}}$ -consistency rate is established for a fixed number of (non-nuisance) parameters.<sup>7</sup>

### 3.2 Selection Consistency

In addition to the consistency result established in Proposition 1, it is also desirable that the proposed estimator is selection consistent. This requires that, with probability approaching 1, the estimates of the truly zero coefficients are zero, while those of nonzero coefficients are nonzero.

**Assumption 5.**  $\min\{\gamma_\rho, \gamma_\beta\} c_{nT}^{-\zeta} \rightarrow \infty$  as  $n, T \rightarrow \infty$ .

Assumption 5 ensures selection consistency of the estimator by taking advantage of the singularity of the penalty term at zero. Under Assumption 5,  $\min\{\gamma_\rho, \gamma_\beta\} |\theta_p^\dagger|^{-\zeta}$  will be explosive in probability for those truly zero  $\theta_p$  and as a result, asymptotically, the first

---

<sup>7</sup>By imposing sparsity, and, with a judicious and data specific choice of penalty parameters, it may be possible to obtain faster rates of convergence. This may be of particular significance in very high dimensional settings with potentially  $P > nT$ , though such results are not pursued in this paper.

order conditions cannot not be met unless  $\hat{\theta}_p$  takes a value of exactly zero. For the following, recall from the end of Section 2.1 that  $\boldsymbol{\theta}_{(2)}$  contains the truly zero  $\theta_p$ .

**Proposition 2** (Selection Consistency). *Under Assumptions 1–5,*

$$\Pr \left( \|\hat{\boldsymbol{\theta}}_{(2)}\|_2 = 0 \right) \rightarrow 1 \text{ as } n, T \rightarrow \infty. \quad (15)$$

Proposition 2 demonstrates that the estimator will correctly set coefficients with a true value of zero to exactly zero with probability approaching 1. Moreover, the consistency result proved in Proposition 1 implies that, with probability approaching 1, the estimates of nonzero coefficients must be nonzero. Thus together, Propositions 1 and 2 indicate that, with an appropriate choice of regularisation parameters, the PQMLE is model selection consistent.

### 3.3 Asymptotic Distribution

An implication of the model selection consistency result obtained in Proposition 2 is that the asymptotic distribution of nonzero coefficient estimates coincides with that of the infeasible ‘oracle’ estimator, which uses knowledge of which parameters are truly zero. The limiting distribution of the nonzero coefficient estimates is derived appealing to this result, and, in keeping with the high dimensional literature, this is done indirectly, by considering arbitrary linear combinations of parameters. In adopting this approach, the results which are obtained are pointwise and therefore, as remarked in Section 2.4 it is important for the validity of these results that the true parameters have fixed values that are, by Assumption 3.1, well separated from zero. A consequence of the lack of uniformity is that the finite sample distribution of the estimator may be quite different to that derived in Theorem 1; a point made clear by Leeb and Pötscher (2005). However, this is a broader issue in the literature and is particularly difficult to overcome in models of significant complexity, where obtaining uniform results is often challenging.

**Assumption 6.**

$$6.1 \quad \frac{P^4}{\min\{n, T\}} \rightarrow 0 \text{ as } n, T \rightarrow \infty.$$

$$6.2 \quad \frac{1}{n} \boldsymbol{\Lambda}^{0'} \boldsymbol{\Lambda}^0 \xrightarrow{p} \boldsymbol{\Sigma}_{\boldsymbol{\Lambda}^0} \text{ as } n \rightarrow \infty \text{ with } \boldsymbol{\Sigma}_{\boldsymbol{\Lambda}^0} \text{ being a } R^0 \times R^0 \text{ positive definite matrix.}$$

$$6.3 \quad \frac{1}{T} \mathbf{F}^{0'} \mathbf{F}^0 \xrightarrow{p} \boldsymbol{\Sigma}_{\mathbf{F}^0} \text{ as } T \rightarrow \infty \text{ with } \boldsymbol{\Sigma}_{\mathbf{F}^0} \text{ being a } R^0 \times R^0 \text{ positive definite matrix.}$$

$$6.4 \quad \frac{T}{n} \rightarrow c \text{ with } 0 < c < \infty.$$



6.5  $R = R^0$ .

6.6  $\max\{\gamma_\rho, \gamma_\beta\}\sqrt{PnT} = o(1)$ .

Assumption 6.1 ensures that the estimation of the coefficients has a negligible effect on the estimation of the factors and the loadings. Lu and Su (2016), who consider estimation of a standard regression model without interaction, require  $P^2/\min\{n, T\} \rightarrow 0$  for analogous purposes. A stronger condition is needed here to ensure that the estimators of the reduced form factors  $\mathbf{S}^{-1}(\boldsymbol{\rho})\boldsymbol{\Lambda}$  converge sufficiently fast, since the reduced form is implicitly used in instrumenting the endogenous variables. As  $\mathbf{S}(\boldsymbol{\rho}) = \mathbf{I}_n - \sum_{q=1}^Q \rho_q \mathbf{W}_q$  involves an increasing number of weights matrices, the number of these cannot be allowed to increase too quickly. Moreover the convergence of the covariance matrix requires further limits on the growth of  $P$ . Fan and Peng (2004) require  $P^5/n \rightarrow 0$ , which corresponds to Assumption 6.1 in a cross-sectional framework. The condition given in Liu (2017), in a cross-sectional spatial model without a factor structure error effects, also requires  $P^5/n \rightarrow 0$ . Assumptions 6.2 and 6.3 impose that the factors are strong, that is to say that the factors and loadings have a nonnegligible impact on the variance of the unobserved term  $\boldsymbol{\eta} := (\boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_T)$ . Other authors consider models with weak factors however this is not pursued here. Assumption 6.4 requires  $n$  and  $T$  to grow in proportion. Similar asymptotic regimes are assumed in several papers in which biases arise in models with interactive fixed effects, and which use similar estimation approaches. Examples of these include Moon and Weidner (2017) and Shi and Lee (2017). Other papers, such as Bai (2009) and Lu and Su (2016), consider regimes where both  $n/T^2, T/n^2 \rightarrow 0$ , which provide similar limits on the relative growth rates of  $n$  and  $T$ . Assumption 6.5 requires the true number of factors to be known. Nonetheless, Proposition 1 shows that the PQML estimator remains consistent as long as the number of factors is not understated; that is  $R \geq R^0$ . In the absence of interaction, Moon and Weidner (2015) show that the asymptotic distribution of a least squares estimator is unaffected by overstatement of the number of factors, under certain conditions. It might, therefore, be expected that this extends to the present setting, however, since there may be significant complications in obtaining such results, the asymptotic distribution is derived under the assumption  $R = R^0$ . Section 4.2 shows how the number of factors can be chosen consistently with reference to an information criterion. Assumption 6.6 strengthens the restrictions on the penalty term.

Let  $\mathcal{D}$  denote the sigma algebra generated by  $\boldsymbol{\mathcal{X}}_1^*, \dots, \boldsymbol{\mathcal{X}}_K^*, \boldsymbol{\Lambda}^0$  and  $\mathbf{F}^0$ . With a slight abuse of notation, in the following the subscripts  $p$  and  $q$  are also used to refer to an element in the indices  $q = 1, \dots, Q^0$  and  $p = 1, \dots, P^0$  which indexes quantities associated only with nonzero parameter values. Define  $\bar{\mathbf{Z}}_p := \mathbb{E}[\mathbf{Z}_p|\mathcal{D}]$ ,  $\mathbf{Z}_p := \mathbf{M}_{\boldsymbol{\Lambda}^0} \bar{\mathbf{Z}}_p \mathbf{M}_{\mathbf{F}^0} + (\mathbf{Z}_p - \bar{\mathbf{Z}}_p)$ ,

$\mathbf{Z}_{(1)} := (\text{vec}(\mathbf{Z}_1), \dots, \text{vec}(\mathbf{Z}_{P^0}))$ , and  $\mathbf{Z}_{(1)} := (\text{vec}(\mathbf{Z}_1), \dots, \text{vec}(\mathbf{Z}_{P^0}))$ , that is,  $\mathbf{Z}_{(1)}$  and  $\mathbf{Z}_{(1)}$  contain only covariates associated with nonzero parameters. Also, let

$$\mathbf{D} := \frac{1}{\sigma_0^2} \frac{1}{nT} \mathbf{Z}'_{(1)} (\mathbf{M}_{\mathbf{F}^0} \otimes \mathbf{M}_{\mathbf{A}^0}) \mathbf{Z}_{(1)} + \begin{pmatrix} \mathbf{\Omega} & \mathbf{0}_{Q^0 \times K^0} \\ \mathbf{0}_{K^0 \times Q^0} & \mathbf{0}_{K^0 \times K^0} \end{pmatrix}, \quad (16)$$

$$\mathbf{V} := \frac{\mathcal{M}_\varepsilon^3}{\sigma_0^4} (\mathbf{\Phi} + \mathbf{\Phi}') + \frac{\mathcal{M}_\varepsilon^4 - 3\sigma_0^4}{\sigma_0^4} \begin{pmatrix} \mathbf{\Xi} & \mathbf{0}_{Q^0 \times K^0} \\ \mathbf{0}_{K^0 \times Q^0} & \mathbf{0}_{K^0 \times K^0} \end{pmatrix}, \quad (17)$$

where the matrices  $\mathbf{\Omega}$  and  $\mathbf{\Xi}$  are  $Q^0 \times Q^0$  with elements  $\Omega_{qq'} := \frac{1}{n} \text{tr}(\mathbf{G}_q(\mathbf{G}_{q'} + \mathbf{G}_{q'}')) - \frac{2}{n^2} \text{tr}(\mathbf{G}_q) \text{tr}(\mathbf{G}_{q'})$  and  $\Xi_{qq'} := \sum_{t=1}^T \sum_{i=1}^n (\mathbf{G}_q^*)_{ii} (\mathbf{G}_{q'}^*)_{ii}$ , respectively, for  $q, q' = 1, \dots, Q^0$ , and with  $\mathbf{G}_q^* := \mathbf{G}_q - \frac{1}{n} \text{tr}(\mathbf{G}_q) \mathbf{I}_n$ . The matrix  $\mathbf{\Phi}$  is  $P^0 \times P^0$  and has the structure  $\mathbf{\Phi} := (\bar{\mathbf{\Phi}}', \mathbf{0}_{P^0 \times K^0})'$ , with  $\bar{\mathbf{\Phi}}_{qp} := \sum_{t=1}^T \sum_{i=1}^n (\mathbf{Z}_p)_{it} (\mathbf{G}_q^*)_{ii}$ , for  $q = 1, \dots, Q^0$  and  $p = 1, \dots, P^0$ .

**Assumption 7.**

7.1 For some fixed integer  $L$ ,  $\mathbf{S}$  is a nonstochastic  $L \times P^0$  matrix such that  $\mathbf{S}\mathbf{S}'$  converges to a (entrywise) nonnegative matrix with eigenvalues bounded away from zero and infinity as  $n, T \rightarrow \infty$ .

7.2 There exist nonstochastic  $P^0 \times P^0$  matrices  $\mathbf{D} := \mathbb{E}[\mathbf{D}]$  and  $\mathbf{V} := \mathbb{E}[\mathbf{V}]$  such that  $\|\mathbf{D} - \mathbf{D}\|_2 = o_P(1)$ ,  $\|\mathbf{V} - \mathbf{V}\|_2 = o_P(1)$ , and the eigenvalues of  $\mathbf{D}$ ,  $\mathbf{V}$  and  $\mathbf{D} + \mathbf{V}$  are bounded from below by zero and from above by a constant.

Since the limiting distribution of the estimator is difficult to derive directly, a selection matrix  $\mathbf{S}$  is introduced with a finite dimension  $L$ . Assumption 7.1 sets out basic properties of this matrix. Assumption 7.2 ensures that the covariance matrix of the PQMLE is well defined asymptotically. Let  $\mathcal{M}_\varepsilon^m$  denote the  $m$ -th raw moment of  $\varepsilon_{it}$ ,  $\mathbf{J}_h := (\mathbf{0}_{T \times (T-h)}, \mathbf{I}_T, \mathbf{0}_{T \times h})'$ , are recall that  $\boldsymbol{\theta}_{(1)}$  contains only those truly nonzero coefficients.

**Theorem 1** (Asymptotic Normality). *Under Assumptions 1–7,*

$$\sqrt{nT} (\mathbf{S}(\mathbf{D} + \mathbf{V})\mathbf{S}')^{-\frac{1}{2}} \mathbf{S}(\mathbf{D}(\hat{\boldsymbol{\theta}}_{(1)} - \boldsymbol{\theta}_{(1)}^0) - \mathbf{b}) \xrightarrow{d} \mathcal{N}(\mathbf{0}_{L \times 1}, \mathbf{I}_L), \quad (18)$$

with

$$\mathbf{b} := \begin{pmatrix} \mathbf{b}^{(1)} \\ \mathbf{0}_{K^0 \times 1} \end{pmatrix} + \begin{pmatrix} \mathbf{b}^{(2)} \\ \mathbf{0}_{K^0 \times 1} \\ \mathbf{b}^{(3)} \end{pmatrix}, \quad (19)$$

where the vector  $\mathbf{b}^{(1)}$  is  $Q^0 \times 1$  with elements  $\mathbb{b}_q^{(1)} := \sqrt{\frac{T}{n}} (\frac{R^0}{n} \text{tr}(\mathbf{G}_q) - \text{tr}(\mathbf{P}_{\Lambda^0} \mathbf{G}_q))$ , the vector  $\mathbf{b}^{(2)}$  is  $Q^0 \times 1$  with elements  $\mathbb{b}_q^{(2)} := -\frac{1}{\sqrt{nT}} \sum_{h=1}^{T-1} \text{tr}(\mathbf{J}_0 \mathbf{P}_{\mathbf{F}^0} \mathbf{J}'_h) \text{tr}(\mathbf{W}_q \mathbf{A}^h \mathbf{S}^{-1})$  and the vector  $\mathbf{b}^{(3)}$  is  $(Q^0 + 1) \times 1$  with first element  $\mathbb{b}_1^{(3)} := -\frac{1}{\sqrt{nT}} \sum_{h=1}^{T-1} \text{tr}(\mathbf{J}_0 \mathbf{P}_{\mathbf{F}^0} \mathbf{J}'_h) \text{tr}(\mathbf{A}^{h-1} \mathbf{S}^{-1})$  and remaining elements  $\mathbb{b}_{q+1}^{(3)} := -\frac{1}{\sqrt{nT}} \sum_{h=1}^{T-1} \text{tr}(\mathbf{J}_0 \mathbf{P}_{\mathbf{F}^0} \mathbf{J}'_h) \text{tr}(\mathbf{W}_q \mathbf{A}^{h-1} \mathbf{S}^{-1})$ .<sup>8</sup>

Theorem 1 describes the asymptotic properties of the estimator for the nonzero coefficients, detailing the asymptotic covariance matrix and the bias terms which arise. Closer inspection reveals the bias  $\mathbf{b}^{(1)}$  is of order  $\sqrt{T/n}$ , while  $\mathbf{b}^{(2)}$  and  $\mathbf{b}^{(3)}$  are of order  $\sqrt{n/T}$ . These biases are a consequence of the incidental parameters in both dimensions of the panel. The bias  $\mathbf{b}^{(1)}$  is comprised to two parts. The first reflects the general loss of information in  $\mathbf{G}_q$  resulting from reducing its rank by  $R^0$  with the projection  $\mathbf{M}_{\Lambda^0}$ . The second depends on the resemblance between the loadings and the network structure; both are sources of cross-sectional dependence and therefore may be conflated. If the column space of  $\mathbf{G}_q$  is orthogonal to the space of loadings, then  $\mathbf{P}_{\Lambda^0} \mathbf{G}_q = \mathbf{0}_{n \times n}$  and the second part of  $\mathbf{b}^{(1)}$  does not feature. The second source of bias is characterised in  $\mathbf{b}^{(2)}$  for the  $\boldsymbol{\rho}$  coefficients, and in  $\mathbf{b}^{(3)}$  for the  $\boldsymbol{\phi}$  coefficients. These two biases arise due to the inclusion of a lagged outcome and are a generalisation of the usual fixed  $T$  bias encountered in dynamic panels with individual fixed effects. As expected, when the number of parameters is fixed, with  $\mathbf{S} = \mathbf{I}_{p^0}$  the distribution collapses to that of the QMLE where the covariance matrix has a typical sandwich form.

### 3.4 Bias Correction

Given the characterisation of the bias term in Theorem 1, it is shown in the following proposition that this can be consistently estimated and the limiting distribution of the PQMLE can be recentred. Let  $\hat{\mathbf{D}}$  and  $\hat{\mathbf{b}}$  denote the analogues of  $\mathbf{D}$  and  $\mathbf{b}$ , respectively, where  $\boldsymbol{\theta}^0, \mathbf{F}^0, \Lambda^0$  and  $\sigma_0^2$  are replaced by their estimates.

**Proposition 3** (Bias Correction). *Under Assumptions 1–7,*

$$\sqrt{nT}(\mathbf{S}(\mathbf{D} + \mathbf{V})\mathbf{S}')^{-\frac{1}{2}} \mathbf{S} \mathbf{D} (\hat{\boldsymbol{\theta}}_{(1)}^c - \boldsymbol{\theta}_{(1)}^0) \xrightarrow{d} \mathcal{N}(\mathbf{0}_{L \times 1}, \mathbf{I}_L), \quad (20)$$

with  $\hat{\boldsymbol{\theta}}_{(1)}^c := \hat{\boldsymbol{\theta}}_{(1)} - \hat{\mathbf{D}}^{-1} \hat{\mathbf{b}}$  being the bias corrected estimator.

---

<sup>8</sup>Note that here it is assumed that  $\phi_1^0$  is nonzero so that  $\mathbb{b}_1^{(3)}$  appears in the bias term.

## 4 Implementation

This section discusses a way in which the estimation procedure proposed in this paper can be implemented and, in particular, describes the approach used to obtain the results in Section 5. This largely concerns how to choose the user specified inputs: the number of factors  $R$ , and the regularisation parameters  $\gamma_\rho, \gamma_\beta$  and  $\zeta$ . Two methods to inform these choices are discussed in Sections 4.1 and 4.2, with the overall suggestion being to proceed in the following way. First, by Proposition 1 the coefficients can be consistently estimated with knowledge only of an upper bound on the number of factors. Thus, with a suitable choice of the penalty parameters (discussed in Section 4.1) penalised estimation can be performed using a large  $R$ , and consistent estimates of the coefficients obtained. Using these coefficient estimates, a pure factor model can be constructed and the true number of factors detected (discussed in Section 4.2). Finally, the model should be re-estimated inputting the detected number of factors to obtain the final estimates. Of course, this multi-step procedure neglects to account for uncertainty at each stage and ideally it would be preferable to select both the penalty parameters and the number of factors jointly, however, the approach adopted here is pragmatic. Additional Monte Carlo results are provided in Appendix J of the Supplementary Material in order to assess the possible impact of varying the number of factors on the properties of the estimator.

### 4.1 Choosing the Penalty Parameters

The fixed regularisation parameter  $\zeta$  can typically be chosen in line with the rate of convergence of the initial estimator, in order to scale the parameter-specific weights  $\omega_p$  appropriately. For example, if  $c_{nT}$  is known to converge to zero slowly,  $\zeta$  can be increased in order to ensure Assumption 5 is satisfied.<sup>9</sup> The other regularisation parameters  $\gamma_\rho$  and  $\gamma_\beta$ , which must converge to zero, could also be chosen simply as sequences which, in combination with  $\zeta$ , ensure Assumptions 3.2, 5 and 6.6 are satisfied.<sup>10</sup> However, as an alternative, this section considers an information criterion that can be used to select  $\gamma_\rho$  and  $\gamma_\beta$ , similar to what is proposed in Lu and Su (2016). This is suggested in order to go some way in tailoring the choice of  $\gamma_\rho$  and  $\gamma_\beta$  to the data. Recalling  $\gamma := (\gamma_\rho, \gamma_\beta)'$ , the information criterion takes

---

<sup>9</sup>In both the simulations and the application  $\zeta$  is set equal to 4, which performs well in practice and, with  $\tilde{\theta}$  as an initial estimate, would also be suitable for a general choice of  $\gamma_\rho$  and  $\gamma_\beta$ ; see footnote 10.

<sup>10</sup>For example, if  $c_{nT} = a_{nT}$ , then with  $\zeta = 4$  and  $\gamma_\rho = \gamma_\beta = 1/\min\{n, T\}$  would satisfy Assumptions 3.2 and 5 as long as  $P^2/\min\{n, T\} \rightarrow 0$ . With  $n \propto T$  under Assumption 6.4, and again, with  $c_{nT} = a_{nT}$  and  $\zeta = 4$ , then  $\gamma_\rho = \gamma_\beta = n^{-3/2}$  would satisfy Assumptions 3.2, 5 and 6.6 as long as  $P^4/\min\{n, T\} \rightarrow 0$ .

the form

$$\text{IC}^*(\gamma) := \hat{\sigma}^2(\gamma) + \varrho_\rho |\mathcal{S}_\rho(\gamma)| + \varrho_\beta |\mathcal{S}_\beta(\gamma)|, \quad (21)$$

where the notation  $\hat{\sigma}^2(\gamma)$  is used for  $\hat{\sigma}^2$  to emphasise the dependence on  $\gamma$ ,  $\varrho_\rho$  and  $\varrho_\beta$  are some positive penalty functions of  $(n, T)$ , and  $\mathcal{S}_\rho(\gamma)$ ,  $\mathcal{S}_\beta(\gamma)$  denote the index sets for the nonzero elements of the parameter estimates under  $\gamma$ . Following closely the exposition in Lu and Su (2016), define  $\mathcal{S}_{F,\rho} := \{1, \dots, Q\}$  and  $\mathcal{S}_{F,\beta} := \{1, \dots, K\}$  as the index sets for the full set of weights matrices and for all covariates respectively. Analogous index sets  $\mathcal{S}_{T,\rho} := \{1, \dots, Q^0\}$  and  $\mathcal{S}_{T,\beta} := \{1, \dots, K^0\}$  are used to denote the relevant covariates and weights matrices. Next, define two closed intervals,  $\Gamma_\rho := [0, \bar{\gamma}_\rho]$  and  $\Gamma_\beta := [0, \bar{\gamma}_\beta]$ , with  $\Gamma_\rho, \Gamma_\beta \subset \mathbb{R}_+$  and where  $\bar{\gamma}_\rho, \bar{\gamma}_\beta$  are two upper bounds beyond which all parameters would be set to zero. The space  $\Gamma := \Gamma_\rho \times \Gamma_\beta$  can be subdivided into three regions:

$$\Gamma^0 := \{\gamma \in \Gamma : \mathcal{S}_\rho(\gamma) = \mathcal{S}_{T,\rho} \text{ and } \mathcal{S}_\beta(\gamma) = \mathcal{S}_{T,\beta}\},$$

$$\Gamma^- := \{\gamma \in \Gamma : \mathcal{S}_\rho(\gamma) \not\supset \mathcal{S}_{T,\rho} \text{ or } \mathcal{S}_\beta(\gamma) \not\supset \mathcal{S}_{T,\beta}\},$$

$$\Gamma^+ := \{\gamma \in \Gamma : \mathcal{S}_\rho(\gamma) \supset \mathcal{S}_{T,\rho}, \mathcal{S}_\beta(\gamma) \supset \mathcal{S}_{T,\beta} \text{ and } |\mathcal{S}_\rho(\gamma)| + |\mathcal{S}_\beta(\gamma)| > |\mathcal{S}_{T,\rho}| + |\mathcal{S}_{T,\beta}|\},$$

where  $|\cdot|$  denotes the cardinality of a set. Respectively, these are the sets of  $\gamma$  in which the true model is selected, the model is underfitted and the model is overfitted. The following assumptions are made.

**Assumption 8.**

$$8.1 \quad \frac{P^2}{\min\{n, T\}} \rightarrow 0 \text{ as } n, T \rightarrow \infty.$$

$$8.2 \quad \text{As } n, T \rightarrow \infty, (\sqrt{Q}a_{nT})^{-1}\varrho_\rho \rightarrow \infty, (\sqrt{Q}a_{nT})^{-1}\varrho_\beta \rightarrow \infty, Q^0\varrho_\rho \rightarrow 0, \text{ and } K^0\varrho_\beta \rightarrow 0.$$

$$8.3 \quad \text{For any } \gamma \in \Gamma^-, \text{ there exists } \sigma_-^2 \text{ such that } \hat{\sigma}^2(\gamma) \xrightarrow{P} \sigma_-^2 > \sigma_0^2.$$

Assumption 8 is analogous to Assumptions A.7 and A.8 in Lu and Su (2016). Assumption 8.2 requires that the penalty functions  $\varrho_\rho$  and  $\varrho_\beta$  relax sufficiently fast as sample size increases. In practice, there may be many functions which satisfy Assumption 8.2, though these may have different impacts in finite samples; for further discussion see Bai and Ng (2002). Assumption 8.3 ensures that underfitted models yield a larger mean squared error than a correctly fitted model.

**Proposition 4** (Information Criterion Consistency). *Under Assumptions 1–5 and 8,*

$$\Pr \left( \inf_{\gamma \in \Gamma^- \cup \Gamma^+} \text{IC}^*(\gamma) > \text{IC}^*(\gamma^0) \right) \rightarrow 1 \text{ as } n, T \rightarrow \infty, \quad (22)$$

for any  $\gamma^0 \in \Gamma^0$ .

## 4.2 Choosing the Number of Factors

Following the procedure outlined at the beginning of Section 4, penalised estimation can first be performed with the number of factors  $R$  set to a large enough value, denoted by  $R_{\max}$ , in order to obtain consistent estimates of the parameters, denoted by  $\check{\rho}$  and  $\check{\beta}$ . A pure factor model can then be constructed as

$$\mathbf{S}(\check{\rho})\mathbf{Y} - \sum_{k=1}^K \check{\beta}_k \mathbf{Z}_k = \mathbf{\Lambda}^0 \mathbf{F}^{0'} + \check{\boldsymbol{\varepsilon}}, \quad (23)$$

with  $\check{\boldsymbol{\varepsilon}} := \sum_{q=1}^Q (\rho_q^0 - \check{\rho}_q) \mathbf{G}_q (\sum_{k=1}^K \beta_k^0 \mathbf{X}_k + \mathbf{\Lambda}^0 \mathbf{F}^{0'} + \boldsymbol{\varepsilon}) + \sum_{k=1}^K (\beta_k^0 - \check{\beta}_k) \mathbf{X}_k + \boldsymbol{\varepsilon}$ . Existing information criteria can then be used to detect the number of factors, and this suggested number can be input into a second estimation step. For example, Shi and Lee (2017) consider information criteria of the form

$$\text{IC}(R) := \log \left( \frac{1}{nT} \sum_{i=R+1}^n \mu_i \left( \left( \mathbf{\Lambda}^0 \mathbf{F}^{0'} + \check{\boldsymbol{\varepsilon}} \right) \left( \mathbf{\Lambda}^0 \mathbf{F}^{0'} + \check{\boldsymbol{\varepsilon}} \right)' \right) \right) + \varrho_f R, \quad (24)$$

with  $\varrho_f$  being a positive penalty function of  $(n, T)$ . With minor modification to Theorem 5 in that paper, it can be shown that the information criterion in (24) is consistent in determining the number of factors, in the sense that  $\lim_{n, T \rightarrow \infty} \Pr(R^* = R^0) = 1$ , with  $R^* := \arg \min_{0 \leq R \leq R_{\max}} \text{IC}(R)$  and under the additional assumption that the penalty function  $\varrho_f$  satisfies  $\varrho_f \rightarrow 0$  and  $a_{nT} \varrho_f \rightarrow \infty$ , with  $a_{nT}$  being the preliminary rate established in Proposition 1.

## 5 Illustration

This section demonstrates the finite sample performance and practicability of the procedure through the use of a small Monte Carlo study and an empirical example.

### 5.1 Simulations

In the following design, the data are generated according to model (1), with the number of parameters and weights matrices increasing with sample size. The design is summarised in Table 1 with a little under half of the parameters taking a true value of 0 for each sample size. Dashes in the table indicate that a covariate is absent.

Table 1: True parameter values

$n$	$T$	$\rho_1^0$	$\rho_2^0$	$\rho_3^0$	$\rho_4^0$	$\rho_5^0$	$\delta_1^0$	$\delta_2^0$	$\delta_3^0$	$\delta_4^0$	$\delta_5^0$	$\delta_{11}^0$	$\delta_{12}^0$	$\delta_{13}^0$	$\delta_{14}^0$	$\delta_{15}^0$	$\phi_1^0$	$\phi_2^0$	$\phi_3^0$	$\phi_4^0$	$\phi_5^0$
25	25	0.2	0.2	0	-	-	3	0	-3	-	-	1	0	-1	-	-	0.15	0	-0.15	-	-
	50	0.2	0.2	0	-	-	3	0	-3	0	-	1	0	-1	-	-	0.15	0	-0.15	-	-
	100	0.2	0.2	0	-	-	3	0	-3	0	3	1	0	-1	-	-	0.15	0	-0.15	-	-
50	25	0.2	0.2	0	0.2	-	3	0	-3	-	-	1	0	-1	0	-	0.15	0	-0.15	0	-
	50	0.2	0.2	0	0.2	-	3	0	-3	0	-	1	0	-1	0	-	0.15	0	-0.15	0	-
	100	0.2	0.2	0	0.2	-	3	0	-3	0	3	1	0	-1	0	-	0.15	0	-0.15	0	-
100	25	0.2	0.2	0	0.2	0	3	0	-3	-	-	1	0	-1	0	1	0.15	0	-0.15	0	0
	50	0.2	0.2	0	0.2	0	3	0	-3	0	-	1	0	-1	0	1	0.15	0	-0.15	0	0
	100	0.2	0.2	0	0.2	0	3	0	-3	0	3	1	0	-1	0	1	0.15	0	-0.15	0	0

The error term  $\varepsilon_{it}$ , the loadings  $\lambda_{ir}^0$  and the factors  $f_{tr}^0$  are generated as standard normal variables.<sup>11</sup> Primitive exogenous variables are generated according to  $x_{kit}^* = \nu + \sum_{r=1}^{R^0} \lambda_{ir}^0 f_{rt}^0 + e_{it}$  with  $\nu$  being uniformly drawn from the integers  $\{-10, \dots, 10\}$  and  $e_{it} \sim \mathcal{N}(0, 2)$ . By design these are correlated with the factors and the loadings and have associated coefficients  $\delta_k^0$ . There are also additional covariates formed by interacting the  $q$ -th weights matrix with the first primitive exogenous regressor in the manner of (2). These covariates have associated coefficients  $\delta_{1q}^0$ . The number of weights matrices is increasing with  $n$ , with the first weights matrix being constructed as if the cross-sectional units were arrayed on a line and connected only to the units immediately to the left and right. This is the simplest example of a path and produces a matrix with ones along the diagonals immediately above and below the main diagonal, and zeros elsewhere. The remaining matrices are specified in similar fashion, but now represent neighbours to the  $q$ -th degree. All matrices are then row normalised. Finally, a lag of outcomes is included, as well as interactions of this lagged outcome and the weights matrices.<sup>12</sup>

Table 2 reports bias corrected estimates  $\hat{\theta}^c$ , across various  $n$  and  $T$ , each with 1000 Monte Carlo replications, and where  $R = R^0 = 3$ .

<sup>11</sup>For simplicity results are reported here only for idiosyncratic errors generated according to a standard normal. Similar results can be obtained under alternative error distributions and additional simulation results are available in Appendix J in the Supplementary Material.

<sup>12</sup>Assumptions 1–8 are verified for this design in Appendix I of the Supplementary Material.

Table 2: Bias of bias corrected estimates of nonzero parameters ( $R = R^0$ )

$n$	$T$	$\rho_1$	$\rho_2$	$\rho_4$	$\delta_1$	$\delta_3$	$\delta_5$	$\delta_{11}$	$\delta_{13}$	$\delta_{15}$	$\phi_1$	$\phi_3$
25	25	0.0002	-0.0004	-	0.0008	-0.0014	-	-0.0027	0.0031	-	-0.0004	0.0004
	50	0.0001	-0.0002	-	-0.0002	-0.0006	-	-0.0016	0.0026	-	-0.0002	0.0002
	100	0.0001	-0.0002	-	0.0001	-0.0005	0.0005	-0.0014	0.0017	-	-0.0001	0.0001
50	25	0.0001	0	-0.0001	0.0005	0.0005	-	-0.0007	0.0005	-	-0.0002	0.0002
	50	0.0002	-0.0003	0	0.0002	-0.0006	-	-0.0005	0.0013	-	-0.0001	0.0001
	100	0	-0.0001	0	-0.0001	-0.0003	0.0003	0	0.0005	-	-0.0002	0.0002
100	25	-0.0001	-0.0002	0.0002	0.0003	-0.0011	-	-0.0004	0.0022	-0.0006	-0.0003	0.0003
	50	0	0	0	0.0003	-0.0002	-	0.0001	0.0005	-0.0006	-0.0003	0.0002
	100	0.0001	-0.0001	0	0	-0.0001	0.0002	-0.0004	0.0007	-0.0001	-0.0002	0.0002

Table 2 shows that the biases are generally decreasing with both  $n$  and  $T$  and tend to be larger for the parameters  $\delta_1, \delta_3$  and  $\delta_5$ , as well as the exogenous spillovers  $\delta_{11}, \delta_{13}$  and  $\delta_{15}$ . This is unsurprising since the covariates  $\mathbf{X}_{\kappa}^*$  are directly correlated with the loadings and the factors by design. The biases of the  $\rho_q$  parameters are lower since these implicitly use the instrument  $\mathbf{G}_q \mathbf{X}_t \boldsymbol{\beta}^0$ , which may not itself be strongly correlated with the factors and the loadings. The same is true of the coefficients  $\phi_1$  and  $\phi_3$ , since the lags  $\mathbf{Y}_{-1}$  and interactions  $\mathbf{W}_q \mathbf{Y}_{-1}$  are less directly correlated with the factors and the loadings. These biases can be favourably compared with Table 6 in Appendix J in the Supplementary Material, which presents biases of the PQMLE without controlling for interactive effects, where there are large biases which persist with  $n$  and  $T$ .

Table 3: Coverage of nonzero parameter estimates ( $R = R^0$ )

$n$	$T$	$\rho_1$	$\rho_2$	$\rho_4$	$\delta_1$	$\delta_3$	$\delta_5$	$\delta_{11}$	$\delta_{13}$	$\delta_{15}$	$\phi_1$	$\phi_3$
25	25	0.901	0.902	-	0.885	0.908	-	0.891	0.897	-	0.904	0.907
	50	0.906	0.922	-	0.921	0.924	-	0.922	0.928	-	0.916	0.922
	100	0.930	0.919	-	0.926	0.929	0.920	0.917	0.930	-	0.929	0.915
50	25	0.920	0.932	0.931	0.924	0.927	-	0.927	0.927	-	0.913	0.920
	50	0.939	0.935	0.931	0.936	0.926	-	0.932	0.917	-	0.926	0.930
	100	0.946	0.942	0.922	0.932	0.934	0.932	0.945	0.921	-	0.921	0.928
100	25	0.929	0.929	0.923	0.930	0.921	-	0.926	0.916	0.932	0.934	0.931
	50	0.937	0.935	0.947	0.941	0.926	-	0.920	0.939	0.939	0.931	0.934
	100	0.947	0.930	0.942	0.950	0.946	0.948	0.941	0.957	0.942	0.922	0.921

Table 3 presents coverage probabilities of Wald confidence intervals based on Theorem 1 and with a nominal coverage of 95%. These generally improve with  $n$  and  $T$ , though due to the complexity of the design it is unsurprising that they do not do so monotonically. Table 4 shows the percentage of true zero parameters correctly estimated as such, with the



procedure performing well and achieving near 100% accuracy across all  $n$  and  $T$ .

Table 4: Percentage of true zeros ( $R = R^0$ )

$n$	$T$	$\rho_3$	$\rho_5$	$\delta_2$	$\delta_4$	$\delta_{12}$	$\delta_{14}$	$\phi_2$	$\phi_4$	$\phi_5$
25	25	99.9	-	100	-	99.9	-	99.9	-	-
	50	99.8	-	100	100	100	-	99.8	-	-
	100	99.8	-	100	100	100	-	99.9	-	-
50	25	100	-	100	-	100	100	100	100	-
	50	99.9	-	100	100	100	100	99.9	99.9	-
	100	99.6	-	100	100	100	100	99.6	99.6	-
100	25	99.9	99.9	99.9	-	99.9	99.9	99.9	99.9	99.9
	50	99.8	99.8	100	100	100	100	99.8	99.8	99.8
	100	99.7	99.8	100	100	100	100	99.7	99.7	99.7

The results reported in Tables 2–4 are computed with the correct number of factors inputted ( $R = R^0 = 3$ ), however, in practice, the true number of factors will not be known. To address this it was suggested in Section 4 to first perform penalised estimation of the model using an upper bound on the number of factors ( $R = R_{\max}$ ) and then to construct a pure factor model and use the information criterion described in Section 4.2 to detect the true number of factors. After this the model can be re-estimated inputting the detected number of factors to obtain the final estimates. In order to assess the effectiveness of this strategy, additional estimations are performed using an upper bound on the number of factors  $R_{\max} = 6 > R^0$ .<sup>13</sup> A pure factor model is then constructed using these estimates and the information criterion (24) computed. Table 5 presents the number of times, as a percentage, that the true number of factors is found to minimise the information criterion. Three variants of this criterion are used (IC1, IC2 and IC3) which differ only in their choice of penalty function  $\varrho_f$ .<sup>14</sup> As sample size increases, the performance of all three variants improves, though there is significant variability between the three criteria.<sup>15</sup>

<sup>13</sup>Table 15 in Appendix J provides additional results with  $R_{\max} = 10$ ; the results are very similar.

<sup>14</sup>The functions used in IC1, IC2 and IC3 are, respectively,  $\log(\min\{n, T\})/\min\{n, T\}$ ,  $((n + T)/(nT)) \log(\min\{n, T\})$  and  $((n + T)/(nT)) \log((nT)/(n + T))$ . For both  $\varrho_\rho$  and  $\varrho_\beta$  in IC\*,  $\log(\min\{n, T\})/\min\{n, T\}$  is used.

<sup>15</sup>The penalty function IC1 is smaller in magnitude than IC2 and IC3 across all samples sizes. Moreover, unlike IC2 and IC3, IC1 only decreases when  $\min\{n, T\}$  decreases. The overall result of this is under-penalisation for a larger  $R$  and poor performance in smaller samples when  $n = T$ .

Table 5: True number of factors is selected % ( $R = R_{\max} = 6$ )

$T$	25			50			100		
$n$	IC1	IC2	IC3	IC1	IC2	IC3	IC1	IC2	IC3
25	0	96.5	79.2	46.1	99.6	99.4	99.9	99.9	99.9
50	43.8	99.1	98.8	7.3	100	100	100	100	100
100	99.7	99.8	99.8	100	100	100	99.9	100	100

To gauge the likely impact of the factors not being known, estimation results with the number factors misspecified are provided in Appendix J in the Supplementary Material. These results illustrate cases in which the correct number of factors  $R^0$  remains fixed at 3, and yet  $R = 1$ ,  $R = 6$  and  $R = 10$  are inputted in estimation. In line with the result in Proposition 1, when the number of factors is underestimated ( $R = 1$ ) large biases persist, while the estimator remains consistent with the number of factors overestimated ( $R = 6$ ), even significantly so ( $R = 10$ ), though overestimation can result in considerable inefficiency.

## 5.2 Application

As an empirical demonstration, the method is applied to study the determinants of economic growth, using a panel data set where several countries are observed over multiple time periods. It is natural to suppose that economic growth might be influenced by unobserved shocks, as well as observable regressors, and in this spirit Lu and Su (2016) estimate a model of economic growth controlling for unobserved factors. In that paper, the authors focus, in particular, on applying shrinkage methods to determine an unknown number of factors. Extending their work to include interaction is well motivated, since one might reasonably expect the growth rates of different countries to be interrelated. Yet in such cases it can be difficult to specify weights matrices a priori. Indeed Durlauf et al. (2009) remark: *“Spatial methods may yet have an important role to play in growth econometrics. However, when these methods are adapted from the spatial statistics literature, they raise the problem of identifying the appropriate notion of space .... countries are perhaps best thought of as occupying some general socio-economic-political space defined by a range of factors; spatial methods then require a means to identify their locations”*. The model studied in this paper may provide insight into growth rate determination, where uncertainty in specifying cross-national interactions provides an example of the type of uncertainty which the present methodology seeks to address.

The data are obtained from Lu and Su (2016), with additional data on income classifications from the World Bank. The outcome  $y_{it}$  is the growth rate (Grth) in real GDP

per capita for one of a cross-section of 108 countries observed between the years 1970–2005. The same 9 primitive exogenous covariates are used as in Lu and Su (2016), which include variables such as life expectancy, population growth, and consumption, investment and government expenditure shares. A series of weights matrices are specified based on grouping countries according to four World Bank classifications: high income ( $\mathbf{W}_1$ ), upper-middle income ( $\mathbf{W}_2$ ), lower-middle income ( $\mathbf{W}_3$ ) and low income ( $\mathbf{W}_4$ ) economies, and reflect the more general notion of a socio-economic space remarked upon on by Durlauf et al. (2009). Each of these weights matrices are constructed by setting the  $(i, j)$ -th element to 1 if country  $i$  and  $j$  share the same income classification, and setting it equal to zero otherwise, before then row normalising each of the matrices.

Table 6: Estimation results without interaction.

$R$		Young	Fert	Life	Popu	Invpri	Con	Gov	Inv	Open	Lag1	IC1	IC2	IC3
0	estimate	0	0	0	−0.462	0	0	0	0.099	0	0.161	3.662	3.662	3.662
	t-stat	0	0	0	−8.030	0	0	0	17.394	0	10.386			
1	estimate	0	0	0	−0.474	0	0	−0.051	0.118	0	0.137	3.508	3.541	3.531
	t-stat	0	0	0	−7.317	0	0	−4.224	18.504	0	8.855			
2	estimate	0	0.444	0	−0.489	0	0	−0.238	0.228	0	0	3.449	3.515 <sup>†</sup>	3.494
	t-stat	0	4.804	0	−5.186	0	0	−9.424	19.112	0	0			
3	estimate	0	0	0	−0.061	0	0	−0.170	0.228	0	0	3.420 <sup>†</sup>	3.519	3.487 <sup>†</sup>
	t-stat	0	0	0	−0.690	0	0	−8.644	19.821	0	0			
6	estimate	0	0.165	0	−0.432	0	0	−0.174	0.217	0	0	3.437	3.636	3.572
	t-stat	0	2.131	0	−4.393	0	0	−7.779	19.524	0	0			

Table 6 reports bias corrected estimates  $\hat{\theta}^c$  in the absence of interaction.<sup>16</sup> Three variants (IC1, IC2 and IC3) of the information criterion given in (24) are computed using estimates generated inputting  $R = R_{\max} = 6$ .<sup>17</sup> In two out of three cases, the information criteria suggest that the number of factors  $R$  is 3, matching the number suggested in Lu and Su (2016). The estimates corresponding to  $R = 3$  can be compared to the results for the AgLasso (which selects  $R = 3$ ) given in Table 7 of Lu and Su (2016). In this case coefficient estimates and t-statistics are similar.

<sup>16</sup>Note that, in the absence of interaction, the quasi-maximum likelihood estimator reduces to the usual principal component least squares estimator (e.g., Bai, 2009).

<sup>17</sup>These variants are the same as those used in simulations.

Table 7: Estimation results with endogenous interaction and temporal lags.

$R$		$\mathbf{W}_1 \times \text{Grth}$	$\mathbf{W}_2 \times \text{Grth}$	$\mathbf{W}_3 \times \text{Grth}$	$\mathbf{W}_4 \times \text{Grth}$	Young	Fert	Life	Popu	Invpri	Con	Gov	Inv	Open
0	estimate	0.210	0.150	0	0.258	0	0	0	-0.492	0	0	0	0.090	0
	t-stat	3.167	1.297	0	3.795	0	0	0	-8.460	0	0	0	15.115	0
1	estimate	0.295	0.289	-0.192	0.345	0	-0.070	0	-0.443	0	0	-0.050	0.111	0
	t-stat	3.688	4.060	-1.475	5.141	0	-1.239	0	-4.977	0	0	-4.511	16.372	0
2	estimate	0.100	0	-0.325	0.207	0	0.355	0	-0.477	0	0	-0.237	0.218	0
	t-stat	1.323	0	-2.383	2.808	0	3.958	0	-5.107	0	0	-9.493	17.823	0
3	estimate	0.195	0	-0.305	0.227	0	-0.001	0	-0.095	0	0	-0.188	0.215	0
	t-stat	2.603	0	-2.277	3.099	0	-0.016	0	-0.953	0	0	-8.129	18.055	0
6	estimate	0	0	-0.202	0	0.093	-0.946	0	-0.570	0	0	-0.225	0.220	0
	t-stat	0	0	-2.224	0	6.179	-4.760	0	-5.703	0	0	-8.536	16.954	0

Table 7 Continued: Estimation results with endogenous interaction and temporal lags.

$R$		Lag1	$\mathbf{W}_1 \times \text{Lag1}$	$\mathbf{W}_2 \times \text{Lag1}$	$\mathbf{W}_3 \times \text{Lag1}$	$\mathbf{W}_4 \times \text{Lag1}$	IC1	IC2	IC3
0	estimate	0.159	0	0	0	0	3.723	3.723	3.723
	t-stat	10.279	0	0	0	0			
1	estimate	0.129	0.172	0	0.400	0	3.496	3.529	3.519
	t-stat	8.145	1.730	0	2.695	0			
2	estimate	0.031	0	0	0.177	0	3.442	3.508 <sup>†</sup>	3.487
	t-stat	1.965	0	0	1.137	0			
3	estimate	0.033	0	0	0.233	0	3.417 <sup>†</sup>	3.516	3.484 <sup>†</sup>
	t-stat	2.070	0	0	1.572	0			
6	estimate	0	0	0	0	0	3.433	3.632	3.568
	t-stat	0	0	0	0	0			

Table 7 reports estimation results once endogenous interaction and dynamic interaction is added. Government spending and investments shares in particular remain highly significant. However there is also evidence to suggest that there are significant endogenous spillovers, especially between high income and low income countries. The results indicate that amongst these two groups of countries, growth rates are interrelated with a positive spillover. In addition, there is evidence to suggest the presence of dynamic spillovers, these being positive, between lower-middle income countries.

## 6 Conclusion

To conclude, this paper considers the estimation of a model of cross-section interaction, whose salient features are a potentially increasing number of weights matrices and a factor structure in the error term. A penalised quasi-maximum likelihood estimator is proposed, in order to perform inference on network spillovers of various kinds, and its asymptotic properties are studied. A small Monte Carlo study reports good finite sample performance, and an empirical application studying the determinants of economic growth finds positive

spillovers between the growth rates of high income and low income countries.

This work could be extended in several directions. For instance, one might consider possible endogeneity of the weights matrices as in Shi and Lee (2018) and Kuersteiner and Prucha (2020), or extend the use of weights matrices to the error term. Since they are observed, the possibility of time varying weights matrices might also be of interest. With some modified assumptions, the consistency result in Proposition 1 could be extended quite readily to this case, though additional work would be required to characterise the asymptotic distribution. Another prospect might be to consider higher dimensional settings, for example, one might consider an entirely unknown weights matrix, modelled in this framework as a series of weights matrices containing a single unitary element. However, identification in this setting would need to be carefully studied since including parameters which increase too quickly with  $n$ , alongside the factor loadings, may present complications. As a final thought, it might also be natural to allow the number of factors to increase with sample size. When the number of interacting cross-sectional units increases, and more units in a network are observed, it might be expected that additional latent structures in the error term would lead to an increase in the rank of the factor term.

## Appendix A. Proofs of Main Results

This appendix provides proofs of the main results before which a series of lemmas are stated. The proofs of these lemmas are given in the Supplementary Material. The following facts are used repeatedly (proofs can be found, for instance, in Moon and Weidner, 2017). Let  $\mathbf{A}$  and  $\mathbf{B}$  be two conformable matrices. Then  $\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_F \leq \sqrt{\text{rank}(\mathbf{A})}\|\mathbf{A}\|_2$ ,  $\|\mathbf{A}\|_2 \leq \sqrt{\|\mathbf{A}\|_1\|\mathbf{A}\|_\infty}$  and  $\|\mathbf{AB}\|_F \leq \|\mathbf{A}\|_F\|\mathbf{B}\|_2 \leq \|\mathbf{A}\|_F\|\mathbf{B}\|_F$ . Let the  $i$ -th row of an  $n \times m$  matrix  $\mathbf{B}$  be denoted  $(\mathbf{B})_{i\cdot}$ , and the  $j$ -th column be denoted  $(\mathbf{B})_{\cdot j}$ . Then  $(\sum_{j=1}^m \|\mathbf{B}_{\cdot j}\|_2^2)^{\frac{1}{2}} = (\sum_{i=1}^n \|\mathbf{B}_{i\cdot}\|_2^2)^{\frac{1}{2}} = \|\mathbf{B}\|_F$ . Finally, under Assumption 1.1,  $\|\varepsilon\|_2 = O_P(\sqrt{\min\{n, T\}})$  (see Latala, 2005).

**Estimated factors and loadings:** The maximiser of  $\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\Lambda})$  with respect to  $\boldsymbol{\Lambda}$  is not unique, since for any  $\boldsymbol{\Lambda}^* = \boldsymbol{\Lambda}\mathbf{H}$ , with  $\mathbf{H}$  being an  $R \times R$  invertible matrix,  $\mathbf{M}_{\boldsymbol{\Lambda}} = \mathbf{M}_{\boldsymbol{\Lambda}^*}$ . In order to achieve uniqueness of the estimators of  $\boldsymbol{\Lambda}$  and  $\mathbf{F}$ , the normalisations that  $\frac{1}{n}\boldsymbol{\Lambda}'\boldsymbol{\Lambda} = \mathbf{I}_R$  and  $\mathbf{F}'\mathbf{F}$  is a diagonal matrix are adopted, see for example Bai (2009).<sup>18</sup>

<sup>18</sup>It is straightforward to see that such a rotation exists. For example, by the singular value decomposition, decompose  $\boldsymbol{\Lambda}\mathbf{F}' = \mathbf{USV}'$ . Let  $\check{\mathbf{\Lambda}}$  be the  $R$  columns of  $\sqrt{n}\mathbf{U}$  associated with the nonzero singular values and  $\check{\mathbf{F}}'$  be the corresponding  $R$  rows of  $\mathbf{SV}'/\sqrt{n}$ . As the columns of  $\mathbf{U}$  and  $\mathbf{V}$  are orthogonal, and  $\mathbf{S}$  is diagonal, it follows that  $\check{\mathbf{\Lambda}}'\check{\mathbf{\Lambda}}/n = \mathbf{I}_R$ ,  $\check{\mathbf{F}}'\check{\mathbf{F}}$  is diagonal and  $\check{\mathbf{\Lambda}}\check{\mathbf{F}}' = \boldsymbol{\Lambda}\mathbf{F}'$ .

Under these normalisations, define

$$\hat{\mathbf{\Lambda}}(\boldsymbol{\theta}) := \arg \min_{\mathbf{\Lambda}: \frac{1}{n} \mathbf{\Lambda}' \mathbf{\Lambda} = \mathbf{I}_R} \left\{ \frac{1}{nT} \sum_{t=1}^T \mathbf{e}_t' \mathbf{M}_{\mathbf{\Lambda}} \mathbf{e}_t \right\} = \arg \max_{\mathbf{\Lambda}: \frac{1}{n} \mathbf{\Lambda}' \mathbf{\Lambda} = \mathbf{I}_R} \left\{ \frac{1}{n} \text{tr} \left( \mathbf{\Lambda}' \frac{1}{nT} \sum_{t=1}^T \mathbf{e}_t \mathbf{e}_t' \mathbf{\Lambda} \right) \right\}. \quad (\text{A.1})$$

It can be shown that the columns of  $\hat{\mathbf{\Lambda}}(\boldsymbol{\theta})$  are equal to  $R$  orthonormal eigenvectors of the matrix  $\frac{1}{nT} \sum_{t=1}^T \mathbf{e}_t' \mathbf{e}_t$  associated with the  $R$  largest eigenvalues and are unique, up to a column-wise sign change. Hereafter  $\hat{\mathbf{\Lambda}} := \hat{\mathbf{\Lambda}}(\hat{\boldsymbol{\theta}})$ .

**Additional notation:** For matrices  $\mathbf{B}$  and  $\mathbf{B}^*$ ,  $\mathbf{B} = \mathbf{B}^* + O_P(a_{nT})$  means that  $\|\mathbf{B} - \mathbf{B}^*\|_2 = O_P(a_{nT})$ . Similarly  $\mathbf{B} = \mathbf{B}^* + o_P(a_{nT})$  means that  $\|\mathbf{B} - \mathbf{B}^*\|_2 = o_P(a_{nT})$ . The elements of the matrices  $\boldsymbol{\mathcal{X}}_{\kappa}^*$ ,  $\boldsymbol{\mathcal{X}}_k$ ,  $\boldsymbol{\mathcal{Z}}_p$ ,  $\boldsymbol{\varepsilon}$ ,  $\mathbf{\Lambda}$  and  $\mathbf{F}$  are respectively denoted  $x_{\kappa it}^*$ ,  $x_{kit}$ ,  $z_{pit}$ ,  $\varepsilon_{it}$ ,  $\lambda_{ir}$  and  $f_{tr}$ . For any other  $n \times m$  matrix  $\mathbf{B}$ , the  $(i, j)$ -th element is denoted  $(\mathbf{B})_{ij}$ . For brevity,  $\hat{\sigma}^2 := \hat{\sigma}^2(\hat{\boldsymbol{\theta}}, \hat{\mathbf{\Lambda}})$ ,  $\mathbf{Z}_t^*$  denotes the  $n \times P$  matrix  $(\mathbf{W}_1 \mathbf{y}_t, \dots, \mathbf{W}_Q \mathbf{y}_t, \mathbf{X}_t)$ . Finally, the  $l$ -th raw moment of some random variable  $s$  is denoted  $\mathcal{M}_s^l$ .

**Lemma A.1.** *For any positive definite matrix  $\mathbf{B}$ ,  $\det(\mathbf{B})^{\frac{1}{n}} \leq \frac{1}{n} \text{tr}(\mathbf{B})$ , with equality if and only if  $\mathbf{B} = c \mathbf{I}_n$  for some  $c > 0$ .*

**Lemma A.2.** *Under Assumptions 1–2,*

- (i)  $\mathbf{S}(\boldsymbol{\rho}) \mathbf{S}^{-1} = \mathbf{I}_n + \sum_{q=1}^Q (\rho_q^0 - \rho_q) \mathbf{G}_q$ ;
- (ii)  $\|\boldsymbol{\mathcal{Z}}_p\|_2 \leq \|\boldsymbol{\mathcal{Z}}_p\|_F = O_P(\sqrt{nT})$  for  $p = 1, \dots, P$ ;
- (iii)  $\|\mathbf{\Lambda}^0\|_2 \leq \|\mathbf{\Lambda}^0\|_F = O_P(\sqrt{n})$ ,  $\|\mathbf{F}^0\|_2 \leq \|\mathbf{F}^0\|_F = O_P(\sqrt{T})$ ;
- (iv)  $(\sum_{p=1}^P \|\boldsymbol{\mathcal{Z}}_p\|_2^2)^{\frac{1}{2}}, (\sum_{t=1}^T \|\mathbf{Z}_t\|_2^2)^{\frac{1}{2}} = O_P(\sqrt{PnT})$ ;
- (v)  $\mathbb{E}[\sum_{p=1}^P (\text{tr}(\boldsymbol{\mathcal{Z}}_p' \mathbf{S}(\boldsymbol{\rho}) \mathbf{S}^{-1} \boldsymbol{\varepsilon})^2) = O(PnT)$ ;
- (vi)  $\|\boldsymbol{\varepsilon}\|_F = O_P(\sqrt{nT})$ ;
- (vii)  $(\sum_{t=1}^T \|\mathbf{X}_t \boldsymbol{\beta}^0\|_2^2)^{\frac{1}{2}} = O_P(\sqrt{nT})$ ;
- (viii)  $\|\mathbf{S}(\boldsymbol{\rho}) \mathbf{S}^{-1} - \mathbf{I}_n\|_2 = O_P(\sqrt{Q} \|\boldsymbol{\theta}^0 - \boldsymbol{\theta}\|_2)$ .

**Lemma A.3.** *Under Assumptions 1–4,*

- (i)  $(\frac{1}{nT} \sum_{t=1}^T \|\mathbf{Z}_t(\boldsymbol{\theta}^0 - \boldsymbol{\theta})\|_2^2)^{\frac{1}{2}} = O_P(\|\boldsymbol{\theta}^0 - \boldsymbol{\theta}\|_2)$ ;
- (ii)  $\hat{\sigma}^{-2}(\hat{\boldsymbol{\theta}}, \mathbf{\Lambda}) = O_P(1)$ .

**Lemma A.4.** *Under Assumptions 1–6,*

$$\begin{aligned} D\sqrt{nT}(\hat{\boldsymbol{\theta}}_{(1)} - \boldsymbol{\theta}_{(1)}^0) &= \frac{1}{\sigma_0^2} \frac{1}{\sqrt{nT}} \mathbf{Z}'_{(1)} (\mathbf{M}_{\mathbf{F}^0} \otimes \mathbf{M}_{\Lambda^0}) \text{vec}(\boldsymbol{\varepsilon}) \\ &+ \frac{1}{\sigma_0^2} \frac{1}{\sqrt{nT}} \begin{pmatrix} \text{tr}((\mathbf{G}_1^* \boldsymbol{\varepsilon})' \mathbf{M}_{\Lambda^0} \boldsymbol{\varepsilon} \mathbf{M}_{\mathbf{F}^0}) \\ \vdots \\ \text{tr}((\mathbf{G}_{Q^0}^* \boldsymbol{\varepsilon})' \mathbf{M}_{\Lambda^0} \boldsymbol{\varepsilon} \mathbf{M}_{\mathbf{F}^0}) \\ \mathbf{0}_{K^0 \times 1} \end{pmatrix} + o_P(1), \end{aligned}$$

where the matrix  $\mathbf{D}$  is defined in equation (16) and the matrices  $\mathbf{G}_q^*$  with  $q = 1, \dots, Q^0$  are those associated with nonzero coefficients.

**Lemma A.5.** *Under Assumptions 1–6,*

$$\begin{aligned} \text{(i)} \quad & \|\mathbf{D}^{-1} - \hat{\mathbf{D}}^{-1}\|_2 = O_P((Q^0)^{1.5} P^0 \|\boldsymbol{\theta}^0 - \hat{\boldsymbol{\theta}}\|_2) + O_P\left(\frac{Q^0 P^0}{\sqrt{\min\{n, T\}}}\right); \\ \text{(ii)} \quad & \mathbb{E} \left[ \sum_{q=1}^{Q^0} (\text{tr}((\mathbf{G}_q^* \boldsymbol{\varepsilon})' \mathbf{P}_{\Lambda^0} \boldsymbol{\varepsilon}) - \sigma_0^2 T \text{tr}(\mathbf{P}_{\Lambda^0} \mathbf{G}_q^*))^2 \right] = O(Q^0 T); \\ \text{(iii)} \quad & \mathbb{E} \left[ \sum_{q=1}^{Q^0} (\text{tr}((\mathbf{G}_q^* \boldsymbol{\varepsilon})' \mathbf{P}_{\Lambda^0} \boldsymbol{\varepsilon} \mathbf{P}_{\mathbf{F}^0}) - \sigma_0^2 R^0 \text{tr}(\mathbf{P}_{\Lambda^0} \mathbf{G}_q^*))^2 \right] = O(Q^0); \\ \text{(iv)} \quad & \mathbb{E} \left[ \sum_{q=1}^{Q^0} (\text{tr}((\mathbf{G}_q^* \boldsymbol{\varepsilon})' \boldsymbol{\varepsilon} \mathbf{P}_{\mathbf{F}^0}) - \sigma_0^2 R^0 \text{tr}(\mathbf{G}_q^*))^2 \right] = O(Q^0 n); \\ \text{(v)} \quad & \frac{1}{\sigma_0^2} \frac{1}{\sqrt{nT}} \begin{pmatrix} \text{tr}((\mathbf{Z}_1 - \bar{\mathbf{Z}}_1)' (\mathbf{P}_{\Lambda^0} \boldsymbol{\varepsilon} + \mathbf{M}_{\Lambda^0} \boldsymbol{\varepsilon} \mathbf{P}_{\mathbf{F}^0})) \\ \vdots \\ \text{tr}((\mathbf{Z}_{P^0} - \bar{\mathbf{Z}}_{P^0})' (\mathbf{P}_{\Lambda^0} \boldsymbol{\varepsilon} + \mathbf{M}_{\Lambda^0} \boldsymbol{\varepsilon} \mathbf{P}_{\mathbf{F}^0})) \end{pmatrix} = \begin{pmatrix} \mathbf{b}^{(2)} \\ \mathbf{0}_{K^0 \times 1} \\ \mathbf{b}^{(3)} \end{pmatrix} + o_P(1), \end{aligned}$$

where the matrices  $\mathbf{G}_q^*$  with  $q = 1, \dots, Q^0$ , and the variables  $\mathbf{Z}_p - \bar{\mathbf{Z}}_p$  with  $p = 1, \dots, P^0$  are those associated with nonzero coefficients.

**Lemma A.6.** *Under Assumptions 1–7,  $\frac{1}{\sqrt{nT}} \frac{1}{\sigma_0^2} (\mathbf{S}(\mathbf{D} + \mathbf{V})\mathbf{S}')^{-\frac{1}{2}} \mathbf{S}\mathbf{c} \xrightarrow{d} \mathcal{N}(\mathbf{0}_{L \times 1}, \mathbf{I}_L)$ , where  $\mathbf{c} := \mathbf{Z}'_{(1)} \text{vec}(\boldsymbol{\varepsilon}) + (\text{tr}(\boldsymbol{\varepsilon}' \mathbf{G}_1^* \boldsymbol{\varepsilon}), \dots, \text{tr}(\boldsymbol{\varepsilon}' \mathbf{G}_{Q^0}^* \boldsymbol{\varepsilon}), \mathbf{0}_{1 \times K^0})'$ , the matrices  $\mathbf{S}, \mathbf{D}$  and  $\mathbf{V}$  are defined in Assumptions 7.1 and 7.2, and the matrices  $\mathbf{G}_q^*$  with  $q = 1, \dots, Q^0$  are those associated with nonzero coefficients.*

**Proof of Proposition 1.** Here only a sketch of the proof is provided. A more detailed version can be found in Appendix D of the Supplementary Material.

### Consistency of the QMLE $\tilde{\boldsymbol{\theta}}$

First, consider the average concentrated quasi-likelihood

$$\mathcal{L}(\boldsymbol{\theta}) := \sup_{\Lambda \in \mathbb{R}^{n \times R}} \left\{ \frac{1}{n} \log(\det(\mathbf{S}(\boldsymbol{\rho}))) - \frac{1}{2} \log(\hat{\sigma}^2(\boldsymbol{\theta}, \Lambda)) \right\}. \quad (\text{A.2})$$

A lower bound for this, denoted  $\underline{\mathcal{L}}(\boldsymbol{\theta}^0)$ , can be established by substituting in the true DGP, and using Assumptions 1.1 and 1.2,

$$\begin{aligned}\underline{\mathcal{L}}(\boldsymbol{\theta}^0) &:= \frac{1}{n} \log(\det(\mathbf{S})) - \frac{1}{2} \log \left( \sigma_0^2 + O_P \left( \frac{1}{\min\{n, T\}} \right) \right) \\ &= \frac{1}{n} \log(\det(\mathbf{S})) - \frac{1}{2} \log (\sigma_0^2 + O_P(a_{nT}^2)) \leq \mathcal{L}(\boldsymbol{\theta}^0).\end{aligned}\tag{A.3}$$

Second, using Lemmas A.2(i), A.2(iv), A.2(v), A.3(i) and Assumption 4.2, an upper bound for  $\mathcal{L}(\boldsymbol{\theta})$ , denoted  $\bar{\mathcal{L}}(\boldsymbol{\theta})$ , can also be established,

$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}) &\leq \frac{1}{n} \log(\det(\mathbf{S}(\boldsymbol{\rho}))) - \frac{1}{2} \log \left( c_1 \|\boldsymbol{\theta}^0 - \boldsymbol{\theta}\|_2^2 + O_P \left( \frac{1}{\min\{n, T\}} \right) \right) \\ &\quad + \frac{\sigma_0^2}{n} \text{tr}((\mathbf{S}(\boldsymbol{\rho})\mathbf{S}^{-1})' \mathbf{S}(\boldsymbol{\rho})\mathbf{S}^{-1}) + O_P \left( \frac{1}{\sqrt{nT}} \right) + \|\boldsymbol{\theta}^0 - \boldsymbol{\theta}\|_2 O_P \left( \sqrt{\frac{P}{nT}} \right) \\ &= \frac{1}{n} \log(\det(\mathbf{S}(\hat{\boldsymbol{\rho}}))) - \frac{1}{2} \log \left( c_1 \|\boldsymbol{\theta} - \boldsymbol{\theta}^0\|_2^2 + O_P(a_{nT}) \|\boldsymbol{\theta} - \boldsymbol{\theta}^0\|_2 + O_P(a_{nT}^2) \right) \\ &\quad + \frac{\sigma_0^2}{n} \text{tr}((\mathbf{S}(\hat{\boldsymbol{\rho}})\mathbf{S}^{-1})' \mathbf{S}(\hat{\boldsymbol{\rho}})\mathbf{S}^{-1}) \\ &=: \bar{\mathcal{L}}(\boldsymbol{\theta}).\end{aligned}\tag{A.4}$$

Now, since  $\tilde{\boldsymbol{\theta}}$  is a global maximiser,  $\mathcal{L}(\boldsymbol{\theta}^0) \leq \mathcal{L}(\tilde{\boldsymbol{\theta}})$  and therefore  $\underline{\mathcal{L}}(\boldsymbol{\theta}^0) \leq \bar{\mathcal{L}}(\tilde{\boldsymbol{\theta}})$ . Using the expressions for these bounds derived in (A.3) and (A.4) gives

$$\begin{aligned}&\frac{1}{n} \log(\det(\mathbf{S})) - \frac{1}{2} \log (\sigma_0^2 + O_P(a_{nT}^2)) \\ &\leq \frac{1}{n} \log(\det(\mathbf{S}(\hat{\boldsymbol{\rho}}))) - \frac{1}{2} \log \left( c_1 \|\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^0\|_2^2 + O_P(a_{nT}) \|\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^0\|_2 + O_P(a_{nT}^2) \right) \\ &\quad + \frac{\sigma_0^2}{n} \text{tr}((\mathbf{S}(\hat{\boldsymbol{\rho}})\mathbf{S}^{-1})' \mathbf{S}(\hat{\boldsymbol{\rho}})\mathbf{S}^{-1}).\end{aligned}\tag{A.5}$$

Multiplying both sides of (A.5) by  $-2$ , exponentiating, and then noticing that, by Lemma A.1,  $\sigma_0^2 \det((\mathbf{S}(\hat{\boldsymbol{\rho}})\mathbf{S}^{-1})' \mathbf{S}(\hat{\boldsymbol{\rho}})\mathbf{S}^{-1})^{\frac{1}{n}} \leq \frac{\sigma_0^2}{n} \text{tr}((\mathbf{S}(\hat{\boldsymbol{\rho}})\mathbf{S}^{-1})' \mathbf{S}(\hat{\boldsymbol{\rho}})\mathbf{S}^{-1})$ , results in

$$0 \geq c_1 \|\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^0\|_2^2 + O_P(a_{nT}) \|\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^0\|_2 + O_P(a_{nT}^2).\tag{A.6}$$

Completing the square,  $0 \geq (\sqrt{c_1} \|\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^0\|_2 + O_P(a_{nT}))^2 + O_P(a_{nT}^2)$ , whereby it follows that  $\|\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^0\|_2 = O_P(a_{nT})$ .

### Consistency of the PQMLE $\hat{\boldsymbol{\theta}}$

Since  $\hat{\boldsymbol{\theta}}$  is the maximiser of the penalised quasi-likelihood function,  $\mathcal{Q}(\boldsymbol{\theta}^0) \leq \mathcal{Q}(\hat{\boldsymbol{\theta}})$ . Thus,

$$\mathcal{Q}(\boldsymbol{\theta}^0) = \mathcal{L}(\boldsymbol{\theta}^0) - \left( \gamma_\rho \sum_{q=1}^Q \omega_q |\rho_q^0| + \gamma_\beta \sum_{k=1}^K \omega_{Q+k} |\beta_k^0| \right)$$



$$\begin{aligned}
&\leq \mathcal{Q}(\hat{\boldsymbol{\theta}}) \\
&= \mathcal{L}(\hat{\boldsymbol{\theta}}) - \left( \gamma_\rho \sum_{q=1}^Q \omega_q |\hat{\rho}_q| + \gamma_\beta \sum_{k=1}^K \omega_{Q+k} |\hat{\beta}_k| \right) \\
&\leq \mathcal{L}(\hat{\boldsymbol{\theta}}).
\end{aligned} \tag{A.7}$$

Consider the penalty term. Under Assumption 3.1,

$$\gamma_\rho \sum_{q=1}^Q \omega_q |\rho_q^0| + \gamma_\beta \sum_{k=1}^K \omega_{Q+k} |\beta_k^0| \leq c_2 \max\{\gamma_\rho, \gamma_\beta\} P^0 \left( \left| \frac{\theta_{\underline{p}}^\dagger}{\theta_{\underline{p}}^0} \right| \right)^{-\zeta} |\theta_{\underline{p}}^0|^{-\zeta}, \tag{A.8}$$

where  $\underline{p} := \arg \min_{1 \leq p \leq P, \theta_p^0 \neq 0} |\theta_p^\dagger|$ . Since the initial estimate  $\boldsymbol{\theta}^\dagger$  satisfies  $\|\boldsymbol{\theta}^\dagger - \boldsymbol{\theta}^0\|_2 = O_P(c_{nT}) = o_P(1)$ , it follows that  $|\theta_{\underline{p}}^\dagger / \theta_{\underline{p}}^0 - 1| \leq \frac{1}{|\theta_{\underline{p}}^0|} \|\boldsymbol{\theta}^\dagger - \boldsymbol{\theta}^0\|_2 = o_P(1)$  which implies  $\theta_{\underline{p}}^\dagger / \theta_{\underline{p}}^0 = O_P(1)$ . Hence,

$$\gamma_\rho \sum_{q=1}^Q \omega_q |\rho_q^0| + \gamma_\beta \sum_{k=1}^K \omega_{Q+k} |\beta_k^0| = \max\{\gamma_\rho, \gamma_\beta\} O_P(P^0) = O_P(a_{nT}^2), \tag{A.9}$$

under Assumption 3.2. Next, using (A.9), and applying the lower and upper bounds derived in (A.3) and (A.4) to (A.7) gives

$$\begin{aligned}
&\frac{1}{n} \log(\det(\mathbf{S})) - \frac{1}{2} \log(\sigma_0^2 + O_P(a_{nT}^2)) + O_P(a_{nT}^2) \leq \frac{1}{n} \log(\det(\mathbf{S}(\hat{\boldsymbol{\rho}}))) \\
&- \frac{1}{2} \log \left( c_1 \|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^0\|_2^2 + O_P(a_{nT}) \|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^0\|_2 + O_P(a_{nT}^2) + \frac{\sigma_0^2}{n} \text{tr}((\mathbf{S}(\hat{\boldsymbol{\rho}}) \mathbf{S}^{-1})' \mathbf{S}(\hat{\boldsymbol{\rho}}) \mathbf{S}^{-1}) \right).
\end{aligned} \tag{A.10}$$

After rearranging and simplifying this becomes

$$\begin{aligned}
&\log \left( \sigma_0^2 \det((\mathbf{S}(\hat{\boldsymbol{\rho}}) \mathbf{S}^{-1})' \mathbf{S}(\hat{\boldsymbol{\rho}}) \mathbf{S}^{-1})^{\frac{1}{n}} + O_P(a_{nT}^2) \right) + O_P(a_{nT}^2) \\
&\geq \log \left( c_1 \|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^0\|_2^2 + O_P(a_{nT}) \|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^0\|_2 + O_P(a_{nT}^2) + \frac{\sigma_0^2}{n} \text{tr}((\mathbf{S}(\hat{\boldsymbol{\rho}}) \mathbf{S}^{-1})' \mathbf{S}(\hat{\boldsymbol{\rho}}) \mathbf{S}^{-1}) \right).
\end{aligned} \tag{A.11}$$

Exponentiating, using Lemma A.1, and the fact that by Assumption 4.4  $O_P(a_{nT}^2) = o_P(1)$  gives the result

$$0 \geq c_1 \|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^0\|_2^2 + O_P(a_{nT}) \|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^0\|_2 + O_P(a_{nT}^2), \tag{A.12}$$

whereby completing the square yields  $\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^0\|_2 = O_P(a_{nT})$ .  $\square$

**Proof of Proposition 2.** Since the PQMLE  $\hat{\theta}$  is consistent for  $\theta^0$  by Proposition 1, and by Assumption 2.1  $\theta^0$  is in the interior of  $\Theta$ ,  $\hat{\theta}$  must also be in the interior of  $\Theta$  w.p.a.1 as  $n, T \rightarrow \infty$ . Thus, w.p.a.1,  $\hat{\theta}$  must solve the first order condition

$$\frac{\partial \mathcal{Q}(\theta, \Lambda)}{\partial \theta} = \frac{\partial \mathcal{L}(\theta, \Lambda)}{\partial \theta} - \frac{\partial \varrho(\theta, \gamma, \zeta)}{\partial \theta} = \mathbf{0}_{P \times 1}, \quad (\text{A.13})$$

where

$$\frac{\partial \mathcal{L}(\theta, \Lambda)}{\partial \theta} = \begin{pmatrix} -\frac{1}{n} \text{tr}(\mathbf{G}_1(\rho)) + \frac{1}{\hat{\sigma}^2(\theta, \Lambda)} \frac{1}{nT} \sum_{t=1}^T (\mathbf{W}_1 \mathbf{y}_t)' \mathbf{M}_\Lambda(\mathbf{S}(\rho) \mathbf{y}_t - \mathbf{X}_t \beta) \\ \vdots \\ -\frac{1}{n} \text{tr}(\mathbf{G}_Q(\rho)) + \frac{1}{\hat{\sigma}^2(\theta, \Lambda)} \frac{1}{nT} \sum_{t=1}^T (\mathbf{W}_Q \mathbf{y}_t)' \mathbf{M}_\Lambda(\mathbf{S}(\rho) \mathbf{y}_t - \mathbf{X}_t \beta) \\ \frac{1}{\hat{\sigma}^2(\theta, \Lambda)} \frac{1}{nT} \sum_{t=1}^T \mathbf{x}'_{1t} \mathbf{M}_\Lambda(\mathbf{S}(\rho) \mathbf{y}_t - \mathbf{X}_t \beta) \\ \vdots \\ \frac{1}{\hat{\sigma}^2(\theta, \Lambda)} \frac{1}{nT} \sum_{t=1}^T \mathbf{x}'_{Kt} \mathbf{M}_\Lambda(\mathbf{S}(\rho) \mathbf{y}_t - \mathbf{X}_t \beta) \end{pmatrix}. \quad (\text{A.14})$$

In the following it is shown that, as  $n, T \rightarrow \infty$ , this first order condition cannot hold unless the estimators of those  $\theta_p$  which have a true value of zero also take a value of exactly zero w.p.a.1. To reach a contradiction, suppose that there is some  $p$ , call this  $p^*$ , for which  $\theta_p^0 = 0$  yet  $\Pr(\hat{\theta}_p = 0)$  does not go to 1 as  $n, T \rightarrow \infty$ . It is first shown that  $\frac{\partial \mathcal{L}(\hat{\theta}, \Lambda)}{\partial \theta_{p^*}}|_{\theta=\hat{\theta}} = O_P(1)$ , i.e., the first order condition evaluated at  $\hat{\theta}$  is not explosive in probability. Since  $\theta_{p^*}$  could be some  $\rho_q$  or  $\beta_k$ , both cases are examined in turn. Consider first the case where  $\theta_{p^*}$  is some  $\rho_q$ . Substituting in the true data generating process, the element of  $\frac{\partial \mathcal{L}(\theta, \Lambda)}{\partial \theta}|_{\theta=\hat{\theta}}$  relating to  $\rho_q$  is equal to

$$\begin{aligned} & -\frac{1}{n} \text{tr}(\mathbf{G}_q(\hat{\rho})) + \frac{1}{\hat{\sigma}^2(\theta, \Lambda)} \frac{1}{nT} \sum_{t=1}^T (\mathbf{W}_q \mathbf{y}_t)' \mathbf{M}_\Lambda(\mathbf{S}(\hat{\rho}) \mathbf{y}_t - \mathbf{X}_t \hat{\beta}) \\ = & -\frac{1}{n} \text{tr}(\mathbf{G}_q(\hat{\rho})) + \frac{1}{\hat{\sigma}^2(\hat{\theta}, \Lambda)} \frac{1}{nT} \sum_{t=1}^T (\mathbf{G}_q \mathbf{X}_t \beta^0)' \mathbf{M}_\Lambda \mathbf{Z}_t (\theta^0 - \hat{\theta}) \\ & + \frac{1}{\hat{\sigma}^2(\hat{\theta}, \Lambda)} \frac{1}{nT} \sum_{t=1}^T (\mathbf{G}_q \mathbf{X}_t \beta^0)' \mathbf{M}_\Lambda \mathbf{S}(\hat{\rho}) \mathbf{S}^{-1} \Lambda^0 \mathbf{f}_t^0 + \frac{1}{\hat{\sigma}^2(\hat{\theta}, \Lambda)} \frac{1}{nT} \sum_{t=1}^T (\mathbf{G}_q \mathbf{X}_t \beta^0)' \mathbf{M}_\Lambda \mathbf{S}(\hat{\rho}) \mathbf{S}^{-1} \varepsilon_t \\ & + \frac{1}{\hat{\sigma}^2(\hat{\theta}, \Lambda)} \frac{1}{nT} \sum_{t=1}^T (\mathbf{G}_q \Lambda^0 \mathbf{f}_t^0)' \mathbf{M}_\Lambda \mathbf{Z}_t (\theta^0 - \hat{\theta}) + \frac{1}{\hat{\sigma}^2(\hat{\theta}, \Lambda)} \frac{1}{nT} \sum_{t=1}^T (\mathbf{G}_q \Lambda^0 \mathbf{f}_t^0)' \mathbf{M}_\Lambda \mathbf{S}(\hat{\rho}) \mathbf{S}^{-1} \Lambda^0 \mathbf{f}_t^0 \\ & + \frac{1}{\hat{\sigma}^2(\hat{\theta}, \Lambda)} \frac{1}{nT} \sum_{t=1}^T (\mathbf{G}_q \Lambda^0 \mathbf{f}_t^0)' \mathbf{M}_\Lambda \mathbf{S}(\hat{\rho}) \mathbf{S}^{-1} \varepsilon_t + \frac{1}{\hat{\sigma}^2(\hat{\theta}, \Lambda)} \frac{1}{nT} \sum_{t=1}^T (\mathbf{G}_q \varepsilon_t)' \mathbf{M}_\Lambda \mathbf{Z}_t (\theta^0 - \hat{\theta}) \\ & + \frac{1}{\hat{\sigma}^2(\hat{\theta}, \Lambda)} \frac{1}{nT} \sum_{t=1}^T (\mathbf{G}_q \varepsilon_t)' \mathbf{M}_\Lambda \mathbf{S}(\hat{\rho}) \mathbf{S}^{-1} \Lambda^0 \mathbf{f}_t^0 + \frac{1}{\hat{\sigma}^2(\hat{\theta}, \Lambda)} \frac{1}{nT} \sum_{t=1}^T (\mathbf{G}_q \varepsilon_t)' \mathbf{M}_\Lambda \mathbf{S}(\hat{\rho}) \mathbf{S}^{-1} \varepsilon_t \end{aligned}$$

$$=: k_1 + \dots + k_{10}. \quad (\text{A.15})$$

Since  $\mathbf{G}_q(\boldsymbol{\rho})$  is UB, terms  $k_5, \dots, k_{10}$  are  $O_P(1)$  by the same arguments as for their counterparts in the proof of Lemma A.3(ii) (terms  $l_2, \dots, l_6$ ; see Supplementary Material), and using the result in that lemma (whereby  $1/\hat{\sigma}^2(\hat{\boldsymbol{\theta}}, \boldsymbol{\Lambda}) = O_P(1)$ ). Since the rank of  $\mathbf{G}_q(\boldsymbol{\rho})$  can be no more than  $n$ , using  $|\text{tr}(\mathbf{B})| \leq \text{rank}(\mathbf{B})\|\mathbf{B}\|_2$  for some square matrix  $\mathbf{B}$  (Moon and Weidner, 2017, Lemma S.4.1(v)), and that  $\mathbf{S}^{-1}(\boldsymbol{\rho})$  and  $\mathbf{W}_q$  are UB, one has

$$|k_1| = \frac{1}{n} |\text{tr}(\mathbf{G}_q(\hat{\boldsymbol{\rho}}))| \leq \|\mathbf{G}_q(\hat{\boldsymbol{\rho}})\|_2 \leq \|\mathbf{S}^{-1}(\hat{\boldsymbol{\rho}})\|_2 \|\mathbf{W}_q\|_2 = O_P(1). \quad (\text{A.16})$$

Using Lemmas A.2(vii), A.3(i) and A.3(ii), as well as Proposition 1, yields

$$\begin{aligned} |k_2| &\leq \frac{1}{\sqrt{nT}} \frac{1}{\hat{\sigma}^2(\hat{\boldsymbol{\theta}}, \boldsymbol{\Lambda})} \|\mathbf{G}_q\|_2 \|\mathbf{M}_\Lambda\|_2 \left( \sum_{t=1}^T \|\mathbf{X}_t \boldsymbol{\beta}^0\|_2^2 \right)^{\frac{1}{2}} \left( \frac{1}{nT} \sum_{t=1}^T \|\mathbf{Z}_t(\boldsymbol{\theta}^0 - \hat{\boldsymbol{\theta}})\|_2^2 \right)^{\frac{1}{2}} \\ &= \frac{1}{\sqrt{nT}} O_P(\sqrt{nT}) O_P(a_{nT}) = O_P(1). \end{aligned} \quad (\text{A.17})$$

The remaining terms,  $k_3$  and  $k_4$ , can be shown to be  $O_P(1)$  similarly, using Lemmas A.2(iii), A.2(vi) A.2(vii) and A.3(ii). Next consider the case where  $\theta_{p^*}$  is some  $\beta_k$ . The element of  $\frac{\partial \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\Lambda})}{\partial \boldsymbol{\theta}}|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}}$  corresponding to  $\beta_k$  is

$$\begin{aligned} &\frac{1}{\hat{\sigma}^2(\hat{\boldsymbol{\theta}}, \boldsymbol{\Lambda})} \frac{1}{nT} \sum_{t=1}^T \mathbf{x}'_{kt} \mathbf{M}_\Lambda (\mathbf{S}(\hat{\boldsymbol{\rho}}) \mathbf{S}^{-1}(\mathbf{X}_t \boldsymbol{\beta}^0 + \boldsymbol{\Lambda}^0 \mathbf{f}_t^0 + \boldsymbol{\varepsilon}_t) - \mathbf{X}_t \hat{\boldsymbol{\beta}}) \\ &= \frac{1}{\hat{\sigma}^2(\hat{\boldsymbol{\theta}}, \boldsymbol{\Lambda})} \frac{1}{nT} \sum_{t=1}^T \mathbf{x}'_{kt} \mathbf{M}_\Lambda \mathbf{Z}_t(\boldsymbol{\theta}^0 - \boldsymbol{\theta}) + \frac{1}{\hat{\sigma}^2(\hat{\boldsymbol{\theta}}, \boldsymbol{\Lambda})} \frac{1}{nT} \sum_{t=1}^T \mathbf{x}'_{kt} \mathbf{M}_\Lambda \mathbf{S}(\hat{\boldsymbol{\rho}}) \mathbf{S}^{-1} \boldsymbol{\Lambda}^0 \mathbf{f}_t^0 \\ &\quad + \frac{1}{\hat{\sigma}^2(\hat{\boldsymbol{\theta}}, \boldsymbol{\Lambda})} \frac{1}{nT} \sum_{t=1}^T \mathbf{x}'_{kt} \mathbf{M}_\Lambda \mathbf{S}(\hat{\boldsymbol{\rho}}) \mathbf{S}^{-1} \boldsymbol{\varepsilon}_t \\ &=: k_{11} + k_{12} + k_{13}. \end{aligned}$$

Using Lemmas A.2(ii), A.2(iii), A.2(vi), A.3(i) and A.3(ii), one has

$$\begin{aligned} |k_{11}| &\leq \frac{1}{nT} \frac{1}{\hat{\sigma}^2(\hat{\boldsymbol{\theta}}, \boldsymbol{\Lambda})} \|\mathbf{M}_\Lambda\|_2 \left( \sum_{t=1}^T \|\mathbf{x}_{kt}\|_2^2 \right)^{\frac{1}{2}} \left( \sum_{t=1}^T \|\mathbf{Z}_t(\boldsymbol{\theta}^0 - \hat{\boldsymbol{\theta}})\|_2^2 \right)^{\frac{1}{2}} \\ &= \frac{1}{\sqrt{nT}} \frac{1}{\hat{\sigma}^2(\hat{\boldsymbol{\theta}}, \boldsymbol{\Lambda})} \|\mathbf{M}_\Lambda\|_2 \|\boldsymbol{\chi}_k\|_F \left( \frac{1}{nT} \sum_{t=1}^T \|\mathbf{Z}_t(\boldsymbol{\theta}^0 - \hat{\boldsymbol{\theta}})\|_2^2 \right)^{\frac{1}{2}} \\ &= \frac{1}{\sqrt{nT}} O_P(\sqrt{nT}) O_P(a_{nT}) = O_P(1), \end{aligned} \quad (\text{A.18})$$

$$\begin{aligned}
|k_{12}| &\leq \frac{1}{nT} \frac{1}{\hat{\sigma}^2(\hat{\boldsymbol{\theta}}, \boldsymbol{\Lambda})} \|\mathbf{M}_{\boldsymbol{\Lambda}}\|_2 \|\mathbf{S}(\hat{\boldsymbol{\rho}}) \mathbf{S}^{-1}\|_2 \|\boldsymbol{\Lambda}^0\|_2 \left( \sum_{t=1}^T \|\mathbf{x}_{kt}\|_2^2 \right)^{\frac{1}{2}} \left( \sum_{t=1}^T \|\mathbf{f}_t^0\|_2^2 \right)^{\frac{1}{2}} \\
&= \frac{1}{nT} \frac{1}{\hat{\sigma}^2(\hat{\boldsymbol{\theta}}, \boldsymbol{\Lambda})} \|\mathbf{M}_{\boldsymbol{\Lambda}}\|_2 \|\mathbf{S}(\hat{\boldsymbol{\rho}}) \mathbf{S}^{-1}\|_2 \|\boldsymbol{\Lambda}^0\|_2 \|\mathcal{X}_k\|_F \|\mathbf{F}^0\|_F \\
&= \frac{1}{nT} O_P(\sqrt{n}) O_P(\sqrt{T}) O_P(\sqrt{nT}) = O_P(1),
\end{aligned} \tag{A.19}$$

and

$$\begin{aligned}
|k_{13}| &\leq \frac{1}{nT} \frac{1}{\hat{\sigma}^2(\hat{\boldsymbol{\theta}}, \boldsymbol{\Lambda})} \|\mathbf{M}_{\boldsymbol{\Lambda}}\|_2 \|\mathbf{S}(\hat{\boldsymbol{\rho}}) \mathbf{S}^{-1}\|_2 \left( \sum_{t=1}^T \|\mathbf{x}_{kt}\|_2^2 \right)^{\frac{1}{2}} \left( \sum_{t=1}^T \|\boldsymbol{\varepsilon}_t\|_2^2 \right)^{\frac{1}{2}} \\
&= \frac{1}{nT} \frac{1}{\hat{\sigma}^2(\hat{\boldsymbol{\theta}}, \boldsymbol{\Lambda})} \|\mathbf{M}_{\boldsymbol{\Lambda}}\|_2 \|\mathbf{S}(\hat{\boldsymbol{\rho}}) \mathbf{S}^{-1}\|_2 \|\mathcal{X}_k\|_F \|\boldsymbol{\varepsilon}\|_F \\
&= \frac{1}{nT} O_P(\sqrt{nT}) O_P(\sqrt{nT}) = O_P(1).
\end{aligned} \tag{A.20}$$

Combining the previous results gives  $\frac{\partial \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\Lambda})}{\partial \theta_{p^*}}|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}} = O_P(1)$ . Turning now to the derivative of the penalty term, evaluated at  $\hat{\boldsymbol{\theta}}$ ,

$$\frac{\partial \varrho(\boldsymbol{\theta}, \boldsymbol{\gamma}, \zeta)}{\partial \theta_{p^*}} \Big|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}} = -\gamma^* \frac{1}{|\theta_{p^*}^\dagger|^\zeta} \frac{\hat{\theta}_{p^*}}{|\hat{\theta}_{p^*}|}, \tag{A.21}$$

where  $\gamma^* \in \{\gamma_\rho, \gamma_\beta\}$  denotes the penalty parameter associated with  $\theta_{p^*}$ . By Assumption 5,  $\min\{\gamma_\rho, \gamma_\beta\} |\theta_{p^*}^\dagger|^{-\zeta}$  is explosive in probability because  $\theta_{p^*}^0 = 0$  and so  $|\theta_{p^*}^\dagger| = |\theta_{p^*}^\dagger - \theta_{p^*}^0| \leq \|\boldsymbol{\theta}^\dagger - \boldsymbol{\theta}^0\|_2 = o_P(1)$  by Assumption 3.3. As such, as  $n, T \rightarrow \infty$ , the first order condition cannot be satisfied since  $\frac{\partial \mathcal{L}(\hat{\boldsymbol{\theta}}, \boldsymbol{\Lambda})}{\partial \theta_{p^*}}|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}} = O_P(1)$  and yet the derivative of the penalty term diverges. This contradicts  $\hat{\boldsymbol{\theta}}$  being a maximiser of the objective function. Therefore, instead, it must be that  $\hat{\theta}_{p^*} = 0$  w.p.a.1 as  $n, T \rightarrow \infty$  for the first order condition (A.13) to be satisfied.  $\square$

**Proof of Theorem 1.** Starting with the expression obtained in Lemma A.4,

$$\begin{aligned}
D\sqrt{nT}(\hat{\boldsymbol{\theta}}_{(1)} - \boldsymbol{\theta}_{(1)}^0) &= \frac{1}{\sigma_0^2} \frac{1}{\sqrt{nT}} \mathbf{Z}'_{(1)} (\mathbf{M}_{\mathbf{F}^0} \otimes \mathbf{M}_{\boldsymbol{\Lambda}^0}) \text{vec}(\boldsymbol{\varepsilon}) \\
&\quad + \frac{1}{\sigma_0^2} \frac{1}{\sqrt{nT}} \begin{pmatrix} \text{tr}((\mathbf{G}_1^* \boldsymbol{\varepsilon})' \mathbf{M}_{\boldsymbol{\Lambda}^0} \boldsymbol{\varepsilon} \mathbf{M}_{\mathbf{F}^0}) \\ \vdots \\ \text{tr}((\mathbf{G}_{Q_0}^* \boldsymbol{\varepsilon})' \mathbf{M}_{\boldsymbol{\Lambda}^0} \boldsymbol{\varepsilon} \mathbf{M}_{\mathbf{F}^0}) \\ \mathbf{0}_{K^0 \times 1} \end{pmatrix} + o_P(1)
\end{aligned}$$

$$\begin{aligned}
&= \frac{1}{\sigma_0^2} \frac{1}{\sqrt{nT}} \mathbf{c} - \frac{1}{\sigma_0^2} \frac{1}{\sqrt{nT}} \begin{pmatrix} \text{tr}((\mathbf{G}_1^* \boldsymbol{\varepsilon})'(P_{\Lambda^0} \boldsymbol{\varepsilon} + M_{\Lambda^0} \boldsymbol{\varepsilon} P_{F^0})) \\ \vdots \\ \text{tr}((\mathbf{G}_{Q^0}^* \boldsymbol{\varepsilon})'(P_{\Lambda^0} \boldsymbol{\varepsilon} + M_{\Lambda^0} \boldsymbol{\varepsilon} P_{F^0})) \\ \mathbf{0}_{K^0 \times 1} \end{pmatrix} \\
&\quad - \frac{1}{\sigma_0^2} \frac{1}{\sqrt{nT}} \begin{pmatrix} \text{tr}((\mathbf{Z}_1 - \bar{\mathbf{Z}}_1)'(P_{\Lambda^0} \boldsymbol{\varepsilon} + M_{\Lambda^0} \boldsymbol{\varepsilon} P_{F^0})) \\ \vdots \\ \text{tr}((\mathbf{Z}_{P^0} - \bar{\mathbf{Z}}_{P^0})'(P_{\Lambda^0} \boldsymbol{\varepsilon} + M_{\Lambda^0} \boldsymbol{\varepsilon} P_{F^0})) \end{pmatrix} + o_P(1),
\end{aligned} \tag{A.22}$$

where  $\mathbf{c} := \mathbf{Z}'_{(1)} \text{vec}(\boldsymbol{\varepsilon}) + (\text{tr}(\boldsymbol{\varepsilon}' \mathbf{G}_1^* \boldsymbol{\varepsilon}), \dots, \text{tr}(\boldsymbol{\varepsilon}' \mathbf{G}_{Q^0}^* \boldsymbol{\varepsilon}), \mathbf{0}_{1 \times K^0})'$ , recalling the definition of  $\mathbf{Z}_p$  and  $\mathbf{Z}$  given just prior to the statement of Assumption 7. By expanding the second term on the right-hand side of (A.22) and applying Lemmas A.5(ii), A.5(iii), A.5(iv), and also applying Lemma A.5(v) to the third term on the right, one obtains

$$D\sqrt{nT}(\hat{\boldsymbol{\theta}}_{(1)} - \boldsymbol{\theta}_{(1)}^0) = \frac{1}{\sigma_0^2} \frac{1}{\sqrt{nT}} \mathbf{c} + \mathbf{b} + o_P(1). \tag{A.23}$$

Rearranging and premultiplying by  $(\mathbf{S}(D + V)\mathbf{S}')^{-\frac{1}{2}}$  gives

$$\sqrt{nT}(\mathbf{S}(D + V)\mathbf{S}')^{-\frac{1}{2}} \mathbf{S} D(\hat{\boldsymbol{\theta}}_{(1)} - D^{-1} \mathbf{b} - \boldsymbol{\theta}_{(1)}^0) = (\mathbf{S}(D + V)\mathbf{S}')^{-\frac{1}{2}} \mathbf{S} \frac{1}{\sqrt{nT}} \frac{1}{\sigma_0^2} \mathbf{c} + o_P(1). \tag{A.24}$$

Finally, using Lemma A.6 and Assumption 7.2,  $(\mathbf{S}(D + V)\mathbf{S}')^{-\frac{1}{2}} \mathbf{S} \frac{1}{\sqrt{nT}} \frac{1}{\sigma_0^2} \mathbf{c} \xrightarrow{d} \mathcal{N}(\mathbf{0}_{L \times 1}, \mathbf{I}_L)$ , which yields the result.  $\square$

**Proof of Proposition 3.** In order to prove the result, it suffices to show that  $\|D^{-1} \mathbf{b} - \hat{D}^{-1} \hat{\mathbf{b}}\|_2 = o_P(1)$ . Observe that

$$\|D^{-1} \mathbf{b} - \hat{D}^{-1} \hat{\mathbf{b}}\|_2 \leq \|D^{-1} - \hat{D}^{-1}\|_2 \|\hat{\mathbf{b}}\|_2 + \|D^{-1}\|_2 \|\mathbf{b} - \hat{\mathbf{b}}\|_2. \tag{A.25}$$

It is straightforward to establish that  $\|D^{-1} - \hat{D}^{-1}\|_2 \|\hat{\mathbf{b}}\|_2 = o_P(1)$  using Lemma A.5(i) and the fact that, under Assumptions 1–6,  $\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^0\|_2 = O_P\left(\sqrt{\frac{P}{nT}}\right)$ , which follows from (F.17) in the proof of Lemma A.4 in the Supplementary Material. For the second term in (A.25),  $\|\mathbf{b} - \hat{\mathbf{b}}\|_2 = o_P(1)$  can be shown using Lemmas A.5(ii)–A.5(v), and the following two results. To simplify notation, assume that  $P = P^0$ ,  $Q = Q^0$  and  $\phi_1^0$  is nonzero. Then,

$$\|\mathbf{G}_q^* - \mathbf{G}_q^*(\hat{\rho})\|_2 = \|\mathbf{G}_q - \mathbf{G}_q(\hat{\rho}) - \frac{1}{n} \text{tr}(\mathbf{G}_q) \mathbf{I}_n + \frac{1}{n} \text{tr}(\mathbf{G}_q(\hat{\rho})) \mathbf{I}_n\|_2$$

$$\begin{aligned}
&\leq \|\mathbf{G}_q - \mathbf{G}_q(\hat{\boldsymbol{\rho}})\|_2 + \frac{1}{n} |\text{tr}(\mathbf{G}_q(\hat{\boldsymbol{\rho}}) - \mathbf{G}_q)| \\
&\leq 2\|\mathbf{G}_q - \mathbf{G}_q(\hat{\boldsymbol{\rho}})\|_2 \\
&= 2\|\mathbf{G}_q(\hat{\boldsymbol{\rho}})(\mathbf{S}(\hat{\boldsymbol{\rho}})\mathbf{S}^{-1} - \mathbf{I}_n)\|_2 \\
&\leq 2\|\mathbf{G}_q(\hat{\boldsymbol{\rho}})\|_2 \|\mathbf{S}(\hat{\boldsymbol{\rho}})\mathbf{S}^{-1} - \mathbf{I}_n\|_2 \\
&= O_P(\sqrt{Q}\|\boldsymbol{\theta}^0 - \hat{\boldsymbol{\theta}}\|_2)
\end{aligned} \tag{A.26}$$

using Lemma A.2(viii). Second,

$$\begin{aligned}
\|\mathbf{A} - \mathbf{A}(\hat{\boldsymbol{\rho}}, \hat{\boldsymbol{\phi}})\|_2 &= \|\mathbf{S}^{-1}(\phi_1^0 \mathbf{I}_n + \sum_{q=1}^Q \phi_{q+1}^0 \mathbf{W}_q) - \mathbf{S}^{-1}(\hat{\boldsymbol{\rho}})(\hat{\phi}_1 \mathbf{I}_n + \sum_{q=1}^Q \hat{\phi}_{q+1} \mathbf{W}_q)\|_2 \\
&\leq \|\mathbf{S}^{-1}(\phi_1^0 \mathbf{I}_n + \sum_{q=1}^Q \phi_{q+1}^0 \mathbf{W}_q) - \mathbf{S}^{-1}(\hat{\boldsymbol{\rho}})(\phi_1^0 \mathbf{I}_n + \sum_{q=1}^Q \phi_{q+1}^0 \mathbf{W}_q)\|_2 \\
&\quad + \|\mathbf{S}^{-1}(\hat{\boldsymbol{\rho}})(\phi_1^0 \mathbf{I}_n + \sum_{q=1}^Q \phi_{q+1}^0 \mathbf{W}_q) - \mathbf{S}^{-1}(\hat{\boldsymbol{\rho}})(\hat{\phi}_1 \mathbf{I}_n + \sum_{q=1}^Q \hat{\phi}_{q+1} \mathbf{W}_q)\|_2 \\
&\leq \|\mathbf{S}^{-1}\|_2 \|(\mathbf{I}_n - \mathbf{S}\mathbf{S}^{-1}(\hat{\boldsymbol{\rho}}))(\phi_1^0 \mathbf{I}_n + \sum_{q=1}^Q \phi_{q+1}^0 \mathbf{W}_q)\|_2 \\
&\quad + \|\mathbf{S}^{-1}(\hat{\boldsymbol{\rho}})\|_2 \|(\phi_1^0 - \hat{\phi}_1) \mathbf{I}_n + \sum_{q=1}^Q (\phi_{q+1}^0 - \hat{\phi}_{q+1}) \mathbf{W}_q\|_2 \\
&\leq \|\mathbf{S}^{-1}\|_2 \|\mathbf{I}_n - \mathbf{S}\mathbf{S}^{-1}(\hat{\boldsymbol{\rho}})\|_2 \|\mathbf{S}\|_2 \|\mathbf{A}\|_2 \\
&\quad + \|\mathbf{S}^{-1}(\hat{\boldsymbol{\rho}})\|_2 (|\phi_1^0 - \hat{\phi}_1| + \sum_{q=1}^Q |\phi_{q+1}^0 - \hat{\phi}_{q+1}| \|\mathbf{W}_q\|_2) \\
&\leq \|\mathbf{S}^{-1}\|_2 \|\mathbf{I}_n - \mathbf{S}\mathbf{S}^{-1}(\hat{\boldsymbol{\rho}})\|_2 \|\mathbf{S}\|_2 \|\mathbf{A}\|_2 + \|\mathbf{S}^{-1}(\hat{\boldsymbol{\rho}})\|_2 \|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^0\|_2 \\
&\quad + \|\mathbf{S}^{-1}(\hat{\boldsymbol{\rho}})\|_2 \|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^0\|_2 \sqrt{Q} \max_{1 \leq q \leq Q} \|\mathbf{W}_q\|_2 \\
&= O_P(\sqrt{Q}\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^0\|_2),
\end{aligned} \tag{A.27}$$

where  $\|\mathbf{I}_n - \mathbf{S}\mathbf{S}^{-1}(\hat{\boldsymbol{\rho}})\|_2 = \|\mathbf{S}\mathbf{S}^{-1}(\hat{\boldsymbol{\rho}})\mathbf{S}(\hat{\boldsymbol{\rho}})\mathbf{S}^{-1} - \mathbf{S}\mathbf{S}^{-1}(\hat{\boldsymbol{\rho}})\|_2 \leq \|\mathbf{S}\mathbf{S}^{-1}(\hat{\boldsymbol{\rho}})\|_2 \|\mathbf{S}(\hat{\boldsymbol{\rho}})\mathbf{S}^{-1} - \mathbf{I}_n\|_2 = O_P(\sqrt{Q}\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^0\|_2)$  using Lemma A.2(viii) and Assumption 2.3. The result then follows.  $\square$

**Proof of Proposition 4.** The proof largely follows the same structure as the proof of Theorem 3.5 in Lu and Su (2016). Details can be found in Appendix D in the Supplementary Material.  $\square$

## References

- Bai, J., 2009. Panel data models with interactive fixed effects. *Econometrica* 77 (4), 1229–1279.
- Bai, J., Li, K., 2021. Dynamic spatial panel data models with common shocks. *forthcoming in Journal of Econometrics*.
- Bai, J., Liao, Y., 2017. Inferences in panel data with interactive effects using large covariance matrices. *Journal of Econometrics* 200 (1), 59–78.
- Bai, J., Ng, S., 2002. Determining the number of factors in approximate factor models. *Econometrica* 70 (1), 191–221.
- Blume, L. E., Brock, W. A., Durlauf, S. N., Jayaraman, R., 2015. Linear social interactions models. *Journal of Political Economy* 123 (2), 444–496.
- Bramoullé, Y., Djebbari, H., Fortin, B., 2009. Identification of peer effects through social networks. *Journal of Econometrics* 150 (1), 41–55.
- de Paula, Á., Rasul, I., Souza, P., 2020. Identifying network ties from panel data: Theory and an application to tax competition. Working paper, CeMMAP.
- Durlauf, S. N., Johnson, P. A., Temple, J. R., 2009. The methods of growth econometrics. In: Mills, T., Patterson, K. (Eds.), *Palgrave Handbook of Econometrics*. Vol. 2. Elsevier, Ch. 8, pp. 1119–1179.
- Fan, J., Peng, H., 2004. Nonconcave penalized likelihood with a diverging number of parameters. *The Annals of Statistics* 32 (3), 928–961.
- Gupta, A., Robinson, P. M., 2015. Inference on higher-order spatial autoregressive models with increasingly many parameters. *Journal of Econometrics* 186 (1), 19–31.
- Gupta, A., Robinson, P. M., 2018. Pseudo maximum likelihood estimation of spatial autoregressive models with increasing dimension. *Journal of Econometrics* 202 (1), 92–107.
- Horn, R. A., Johnson, C. R., 2012. *Matrix Analysis*, 2nd Edition. Cambridge University Press, New York, USA.
- Hsiao, C., 2018. Panel models with interactive effects. *Journal of Econometrics* 206 (2), 645–673.

- Kuersteiner, G. M., Prucha, I. R., 2020. Dynamic spatial panel models: Networks, common shocks, and sequential exogeneity. *Econometrica* 88 (5), 2109–2146.
- Lam, C., Souza, P. C., 2019. Estimation and selection of spatial weight matrix in a spatial lag model. *Journal of Business & Economic Statistics* 38 (3), 693–710.
- Latala, R., 2005. Some estimates of norms of random matrices. *Proceedings of the American Mathematical Society* 133 (5), 1273–1282.
- Lee, L.-F., 2007. Identification and estimation of econometric models with group interactions, contextual factors and fixed effects. *Journal of Econometrics* 140 (2), 333–374.
- Lee, L.-F., Liu, X., Lin, X., 2010. Specification and estimation of social interaction models with network structures. *The Econometrics Journal* 13 (2), 145–176.
- Lee, L.-F., Yu, J., 2010. Estimation of spatial autoregressive panel data models with fixed effects. *Journal of Econometrics* (1), 165–185.
- Lee, L.-F., Yu, J., 2014. Spatial panel data models. In: Baltagi, B. H. (Ed.), *The Oxford Handbook of Panel Data*. Oxford University Press, Ch. 12, pp. 363–401.
- Leeb, H., Pötscher, B. M., 2005. Model selection and inference: Facts and fiction. *Econometric Theory* 21 (1), 21–59.
- Lewbel, A., Qu, X., Tang, X., 2021. Social networks with unobserved links. Working paper.
- Liu, T., 2017. Model selection and adaptive lasso estimation of spatial models. PhD Thesis, The Ohio State University.
- Lu, X., Su, L., 2016. Shrinkage estimation of dynamic panel data models with interactive fixed effects. *Journal of Econometrics* 190 (1), 148–175.
- Moon, H. R., Weidner, M., 2015. Linear regression for panel with unknown number of factors as interactive fixed effects. *Econometrica* 83 (4), 1543–1579.
- Moon, H. R., Weidner, M., 2017. Dynamic linear panel regression models with interactive fixed effects. *Econometric Theory* 33 (1), 158–195.
- Newey, W. K., McFadden, D., 1994. Large sample estimation and hypothesis testing. In: Engle, R. F., McFadden, D. L. (Eds.), *Handbook of Econometrics*. Vol. 4. Elsevier, Ch. 36, pp. 2111–2245.



- Shi, W., Lee, L.-F., 2017. Spatial dynamic panel data models with interactive fixed effects. *Journal of Econometrics* 197 (2), 323–347.
- Shi, W., Lee, L.-F., 2018. A spatial panel data model with time varying endogenous weights matrices and common factors. *Regional Science and Urban Economics* 72, 6–34.
- Wang, Y., 2018. Panel data with high-dimensional factors: Inference on treatment effects with an application to sampled networks. Working paper.
- Yu, J., de Jong, R., Lee, L.-F., 2008. Quasi-maximum likelihood estimators for spatial dynamic panel data with fixed effects when both  $n$  and  $T$  are large. *Journal of Econometrics* 146 (1), 118–134.
- Zhang, X., Yu, J., 2018. Spatial weights matrix selection and model averaging for spatial autoregressive models. *Journal of Econometrics* 203 (1), 1–18.
- Zou, H., 2006. The adaptive lasso and its oracle properties. *Journal of the American Statistical Association* 101 (476), 1418–1429.