

Department of Econometrics and Business Statistics

<http://business.monash.edu/econometrics-and-business-statistics/research/publications>

A Linear Estimator for Factor-Augmented Fixed-T Panels with Endogenous Regressors

Arturas Juodis and Vasilis Sarafidis

February 2020

Working Paper 05/20

A Linear Estimator for Factor-Augmented Fixed-T Panels with Endogenous Regressors

Artūras Juodis^a, Vasilis Sarafidis^{b,*}

^a*Faculty of Economics and Business, University of Groningen.*

^b*Department of Econometrics and Business Statistics, Monash University.*

Abstract

A novel method-of-moments approach is proposed for the estimation of factor-augmented panel data models with endogenous regressors when T is fixed. The underlying methodology involves approximating the unobserved common factors using observed factor proxies. The resulting moment conditions are linear in the parameters. The proposed approach addresses several issues which arise with existing nonlinear estimators that are available in fixed T panels, such as local minima-related problems, a sensitivity to particular normalisation schemes, and a potential lack of global identification. We apply our approach to a large panel of households and estimate the price elasticity of urban water demand. A simulation study confirms that our approach performs well in finite samples.

Keywords: Panel Data, Common Factors, Fixed T Consistency, Moment Conditions, Urban Water Management.

JEL: C13, C15, C23.

1. Introduction

The common factor approach has attracted considerable interest within panel data analysis because it offers a wide scope for controlling for unobservables, including situations where there is cross-sectional dependence, see e.g. Sarafidis and Wansbeek (2012).

Holtz-Eakin et al. (1988), Ahn et al. (2013), Robertson and Sarafidis (2015), Robertson et al. (2018), and Juodis (2018), among others, have proposed various estimators for panels with endogenous covariates and “fixed T ”, where T denotes the number of time

*Corresponding author. Address: 900 Dandenong Road, Caulfield East, Victoria 3145, Australia. E-mail: vasilis.sarafidis@monash.edu

series observations. A common feature of these approaches is that identification relies on nonlinear moment conditions. For panels with “large T ”, popular (least-squares type) methods include those developed by Pesaran (2006) and Bai (2009), known in the literature as Common Correlated Effects “CCE” and Principal Components “PC”, respectively.

Given the computational simplicity inherent in CCE and PC, there exists a large number of empirical papers across several fields employing these estimators when T is large. In contrast, it is fair to say that the literature on fixed- T panels remains largely unnoticed by empirical researchers, despite its aforementioned volume. There are several factors that might help to explain this observation. Firstly, there is usually no underlying theory to guide the selection of good starting values for a potentially large number of nuisance parameters, such as the unobserved common factors. This can potentially lead to local minima-related problems, frequently arising in estimation of factor models using iterative algorithms, see e.g. Jiang et al. (2017). Secondly, identification of the parameters of interest typically requires imposing certain normalisation restrictions. However, as was discussed in Kruiniger (2008) and Juodis and Sarafidis (2018), the choice of the normalisation scheme can be crucial for the properties of estimators that rely on some form of quasi-differencing, depending on the underlying data generating process for the unknown factors. For example, the approach proposed by Holtz-Eakin et al. (1988) requires that all factors take non-zero values in all time periods. Finally, as shown by Hayakawa (2016), the nonlinear moment conditions proposed in the literature do not always satisfy the global identification assumption, which is a necessary condition for consistency of GMM estimation.

The present paper develops a novel GMM approach; the main idea is to replace the unobserved factors with proxies constructed from observables. We put forward two distinct methods for constructing factor proxies. The first one involves the use of a multiple weighting scheme applied on a single observable, which can be either an external variable or a regressor. The second method employs a single weighting scheme applied to multiple variables. We show that these two methods can also be combined. In both cases, the underlying assumption is that the variables employed to construct the factor proxies are driven by the common shocks that are relevant for the main variable of interest. The resulting method of moments estimator has a closed form solution, and avoids the aforementioned issues associated with nonlinear estimators. Under suitable regularity conditions (discussed in the paper), the proposed estimator is consistent and

asymptotically mixed-normal.

In response to issues associated with nonlinear method of moments estimators, West-erlund et al. (2019) recently advocated the use of pooled CCE in fixed- T panels. How-ever, the computational simplicity of CCE comes with a price, as the method requires all regressors to be strictly exogenous and to exhibit a common factor structure. The latter restriction prohibits, for example, nonlinear partial effects, such as in regressions with quadratic terms. Recently, De Vos and Everaert (2019) extended CCE in fixed- T panels to the case of a lagged dependent variable, assuming that all other covariates are strictly exogenous. However, the proposed procedure is no longer linear. In con-trast, the approach presented in this paper does not need bias correction of any sort, and accommodates regressors with different degrees of exogeneity, without imposing the requirement that *all* covariates (and/or instruments) have as many factors as the main variable on interest. At the same time, the resulting GMM estimator possesses the appealing linearity property of the CCE estimator.

We use our approach to estimate the price elasticity of residential water usage demand. This topic is of large interest, not only among economists but also across international environmental agencies, regulators, water utilities and the general public. We find that urban water demand is more elastic in the long-run than in the short-run, which may be attributed to habit formation and technological constraints of water appliance efficiency. This result casts doubt on the potential effectiveness of scarcity pricing to balance demand and supply of water in periods of transitory droughts.

The remainder of this paper is organised as follows. Section 2 introduces a linear panel model with common shocks. Section 3 develops the proposed approach. Section 4 reports a Monte Carlo study to assess the finite sample performance of the estima-tor. Section 5 presents the empirical application. Finally, Section 6 concludes. A Supplementary Appendix to this paper discusses extensions to unbalanced panels and models with observed factors. In addition, it provides further finite sample evidence and contains proofs of our theoretical results.

Notation: The generic constants δ and M are used to denote a small and a large positive real number, respectively. For a generic matrix \mathbf{A} , $\text{vec}(\mathbf{A})$ denotes the vertical column stacking operator, and $\text{Col}(\mathbf{A})$ denotes the column space of \mathbf{A} . Moreover, \otimes denotes the Kronecker product. $\mathbf{l}_{i,s;q}$, $s \leq q$, is defined as $\mathbf{l}_{i,s;q} = (l_{i,s}, \dots, l_{i,q})'$. Finally, all random variables are defined on a common probability space (Ω, \mathcal{A}, P) .

2. Model

We consider the following panel data model with a multi-factor error structure:

$$y_{i,t} = \mathbf{x}'_{i,t}\boldsymbol{\beta} + \boldsymbol{\lambda}'_i\mathbf{f}_t + \varepsilon_{i,t}; \quad i = 1, \dots, N; \quad t = 1, \dots, T, \quad (1)$$

where $\mathbf{x}_{i,t} = (x_{i,t}^{(1)}, \dots, x_{i,t}^{(K)})'$ denotes the vector of explanatory variables and $\boldsymbol{\beta} = (\beta_1, \dots, \beta_K)'$ is the vector of the parameters of interest. \mathbf{f}_t denotes the $[L \times 1]$ vector of unobserved common factors, $\boldsymbol{\lambda}_i$ denotes the associated factor loadings for individual i , and $\varepsilon_{i,t}$ is the remaining error term.

The multi-factor error structure is appealing because it allows for multiple sources of *multiplicative* unobserved heterogeneity, as opposed to the one-way (or two-way) error components structure, which represents *additive* heterogeneity. For example, in a partial adjustment model of factor input prices, the factor component may capture common shocks that hit all producers, albeit with different intensities. In the estimation of production functions, the factor component may absorb different sources of technical inefficiency, which vary over time in an arbitrary way. In an empirical model of household water usage demand, the factor component may capture nonlinear effects of household size (typically unobserved) that depend on time-varying weather conditions.

Stacking the observations over time for each i , the model can be rewritten in vector form as:

$$\mathbf{y}_i = \mathbf{X}_i\boldsymbol{\beta} + \mathbf{F}\boldsymbol{\lambda}_i + \boldsymbol{\varepsilon}_i; \quad i = 1, \dots, N, \quad (2)$$

where $\mathbf{y}_i = (y_{i,1}, \dots, y_{i,T})'$, $\mathbf{X}_i = (\mathbf{x}_{i,1}, \dots, \mathbf{x}_{i,T})'$, $\boldsymbol{\varepsilon}_i = (\varepsilon_{i,1}, \dots, \varepsilon_{i,T})'$, while $\mathbf{F} = (\mathbf{f}_1, \dots, \mathbf{f}_T)'$ is of dimension $[T \times L]$.

Denote by \mathcal{F} the σ -field generated by all common shocks driving the individual-specific variables in the system. As such \mathcal{F} contains all factors \mathbf{F} , but we also allow variables $\{(\mathbf{x}_{i,t}, \boldsymbol{\lambda}_i)\}_{t=1}^T$ to be a function of other common shocks (not necessarily of linear factor structure), resulting in additional sources of dependence across cross-sectional units. For example, one can allow $x_{i,t} = b(\psi_i, g_t, \zeta_{i,t})$, where $b(\cdot)$ is a linear/nonlinear function in all arguments.

Assumption 2.1. *The DGP for all i and t satisfies the following restrictions:*

- (a) $(\mathbf{X}_i, \boldsymbol{\varepsilon}_i, \boldsymbol{\lambda}_i)$ are identically distributed and independent across i , conditional on \mathcal{F} .
- (b) Each time-varying element $p_{i,t}^{(\cdot)}$ in $\mathbf{p}_{i,t} = (p_{i,t}^{(1)}, \dots, p_{i,t}^{(K+1)})' \equiv (\mathbf{x}'_{i,t}, \varepsilon_{i,t})'$ satisfies

$$\mathbb{E} \left[\left| p_{i,t}^{(\cdot)} \right|^{4+\delta} \right] < \infty \text{ for all } t.$$

- (c) Each time-invariant element $\lambda_i^{(\cdot)}$ in $\boldsymbol{\lambda}_i = (\lambda_i^{(1)}, \dots, \lambda_i^{(L)})'$ satisfies $E \left[\left| \lambda_i^{(\cdot)} \right|^{4+\delta} \right] < \infty$.
- (d) $E_{\mathcal{F}}[\varepsilon_{i,t} | \boldsymbol{\lambda}_i, \mathbf{x}_{i,1:\tau_1(t)}^{(1)}, \dots, \mathbf{x}_{i,1:\tau_K(t)}^{(K)}] = 0 \forall t$, for some positive integers $\tau_1(t), \dots, \tau_K(t)$.

Besides the fact that \mathbf{F} is assumed to be random, the above assumptions are standard in the literature, see e.g. Assumption BA.1 in Ahn et al. (2013). Given the conditional independence assumption, all stochastic convergence modes in this paper are conditional on \mathcal{F} . We will emphasize this technicality further in Section 3, when we discuss the asymptotic distribution of the proposed estimator. These assumptions are general enough to allow for conditional heteroskedasticity in both dimensions, e.g. $\varepsilon_{i,t} = \sigma_i \xi_t \eta_{i,t}$, where $\eta_{i,t}$ is i.i.d. over i and t with unit variance, ξ_t is a sequence of constants, while σ_i is an i.i.d. sequence over i . Subject to some additional summability restrictions, the conditional i.i.d. restriction can be further relaxed to conditional independence with heterogenous population moments, without affecting the consistency of the estimator. The Supplementary Appendix provides one such example. For instance, σ_i and $\boldsymbol{\lambda}_i$ could be treated as a sequence of fixed constants, as was advocated in Hsiao et al. (2002). However, such setup would require additional technical restrictions in order for the limits to be well defined, see e.g. p. 996 in Gagliardini et al. (2016) for a related discussion.

Assumption 2.1(d) characterises the exogeneity properties of the covariates. In particular, covariates that satisfy $\tau_k(t) = T$ ($\tau_k(t) = t$) are strictly (weakly) exogenous with respect to the idiosyncratic error component, and endogenous otherwise. The estimator proposed in this paper allows for strictly/weakly exogenous regressors, such as lagged dependent variables and endogenous regressors.

Let \mathbf{z}_i be a $[d \times 1]$ vector containing all internal instruments that are available by Assumption 2.1(d), as well as external instruments satisfying the corresponding assumption. Also, let $\mathbf{S} = \text{diag}(\mathbf{S}_1, \dots, \mathbf{S}_T)$ denote a block-diagonal matrix with a typical block-diagonal entry equal to \mathbf{S}_t , where \mathbf{S}_t is a $[\zeta_t \times d]$ selection matrix of zeros and ones that picks ζ_t valid instruments at time t from \mathbf{z}_i .

Under Assumption 2.1, the following set of $\zeta \equiv \sum_{t=1}^T \zeta_t$ population moment conditions is valid by construction:

$$E_{\mathcal{F}}[\mathbf{Z}'_i(\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta}_0 - \mathbf{F}\boldsymbol{\lambda}_i)] = \mathbf{S} \left(\text{vec} \left(E_{\mathcal{F}} \left[\mathbf{z}_i(\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta}_0)' \right] - \mathbf{G}_{z,\lambda}\mathbf{F}' \right) \right) = \mathbf{0}_{\zeta}, \quad (3)$$

where $\mathbf{Z}'_i \equiv \mathbf{S}(\mathbf{I}_T \otimes \mathbf{z}_i)$ and $\mathbf{G}_{z,\lambda} \equiv E_{\mathcal{F}}(\mathbf{z}_i\boldsymbol{\lambda}'_i)$ is a $[d \times L]$ unknown population matrix that absorbs the unobserved covariances between instruments and factor loadings.

The moment conditions in Eq. (3) give rise to the estimators proposed by Robertson and Sarafidis (2015) (with $\mathbf{G}_{z,\lambda}$ estimated) and Ahn et al. (2013) (with $\mathbf{G}_{z,\lambda}$ quasi-differenced). Either way, the resulting moment conditions are nonlinear and, hence, potentially subject to the issues discussed in Section 1. In what follows, we put forward a strategy that circumvents the nonlinearity of the moment conditions in Eq. (3). For future reference, notice that the last term in Eq. (3) can be rewritten as $\mathbf{S}(\text{vec}(\mathbf{G}_{z,\lambda}\mathbf{F}')) = \mathbf{S}(\mathbf{F} \otimes \mathbf{I}_d)\mathbf{g}_{z,\lambda}$, where $\mathbf{g}_{z,\lambda} = \text{vec}(\mathbf{G}_{z,\lambda})$.

3. A New Approach for Dealing with Unobserved Factors in Fixed-T Panels

Let \mathbf{F}_e denote a $[T \times L_e]$ dimensional matrix with $L_e \geq L$, such that $\mathbf{F} \in \text{Col}(\mathbf{F}_e)$. Furthermore, let $\widehat{\mathbf{F}}_e$ be a consistent estimator of the column space of \mathbf{F}_e . In Section 3.1 we derive the asymptotic properties of the proposed GMM estimator based on $\widehat{\mathbf{F}}_e$, assuming that the model is identified from Eq. (3). In Sections 3.2 and 3.3 we discuss different methods for constructing $\widehat{\mathbf{F}}_e$, depending on the model at hand. Finally, in Sections 3.4 and 3.5 we analyse identification, and we put forward two alternative procedures for implementing our approach in practice.

3.1. The Estimator

Assumption 3.1 below ensures that $\widehat{\mathbf{F}}_e$ is an appropriate plug-in estimator of \mathbf{F}_e . Our setup is sufficiently general in that it allows for two important cases, namely: (i) the number of estimated factors is larger than the true number of factors that enter into the error term of the equation for y ; (ii) the number of factor proxies in $\widehat{\mathbf{F}}_e$ that is required to identify \mathbf{F} is larger than L . These cases are illustrated below in Examples 1 (or 3), and 2, respectively.

Assumption 3.1. *Factor proxies are asymptotically linear such that $\sqrt{N}(\widehat{\mathbf{F}}_e - \mathbf{F}_e\mathbf{A}_N) = (N^{-1/2} \sum_{i=1}^N \boldsymbol{\Psi}_i) + o_P(1)$ and $\boldsymbol{\Psi}_i = (\mathbf{I}_T \otimes \boldsymbol{\psi}_i')\mathbf{B}'_N$. Here \mathbf{A}_N is an $[L_e \times L_e]$ rotation matrix, $\boldsymbol{\psi}_i$ is a $[q \times 1]$ vector, and \mathbf{B}_N is a $[L_e \times Tq]$ selection matrix. Furthermore,*

- (a) $\boldsymbol{\psi}_i$ is identically distributed and independent across i , conditional on \mathcal{F} . Moreover, $\psi_i^{(\cdot)}$ in $\boldsymbol{\psi}_i = (\psi_i^{(1)}, \dots, \psi_i^{(q)})'$ satisfies $\mathbb{E} \left[\left| \psi_i^{(\cdot)} \right|^{4+\delta} \right] < \infty$, with $\mathbb{E}_{\mathcal{F}}[\boldsymbol{\psi}_i] = \mathbf{0}_q$.
- (b) \mathbf{A}_N and \mathbf{B}_N are such that $\mathbf{A}_N \xrightarrow{p} \mathbf{A}$ and $\mathbf{B}_N \xrightarrow{p} \mathbf{B}$. Here \mathbf{A}, \mathbf{B} are \mathcal{F} -measurable.
- (c) For any value of N , including $N \rightarrow \infty$: $\text{rk}(\mathbf{A}_N) = L_e$ a.s.
- (d) $\mathbf{F} \in \text{Col}(\mathbf{F}_e)$ and $\text{rk}(\mathbf{F}_e) = L_e$ a.s.

Assumption 3.1 is fairly intuitive. In particular, parts (a)-(b) are employed so as to enable the application of a standard central limit theorem. Notice that, similarly to \mathbf{Z}'_i , $\mathbf{\Psi}'_i$ has a Kronecker product form. Thus, \mathbf{B}_N is a selection matrix which, for each point t (i.e. for each $\hat{\mathbf{f}}_{t,e}$), selects those elements of $\boldsymbol{\psi}_i$ that are of first-order (\sqrt{N}) importance. Therefore, the length of $\boldsymbol{\psi}_i$ need not be equal to TL_e , the number of elements in $\hat{\mathbf{F}}_e$. Section 3.3 provides more details. Parts (c)-(d) ensure that the factor proxies in $\hat{\mathbf{F}}_e$ asymptotically identify the L_e -dimensional column space of \mathbf{F}_e a.s. Finally, if part (d) is violated such that $\mathbf{F} \notin \text{Col}(\mathbf{F}_e)$, then in the limit $\hat{\mathbf{F}}_e$ will not approximate \mathbf{F} in the model equation.

Let $\boldsymbol{\beta}_0$ denote the true value of $\boldsymbol{\beta}$, and $\mathbf{g}_0 \equiv \text{vec} \left(\mathbf{G}_{z,\lambda_e} (\mathbf{A}_N^{-1})' \right)$, which is of dimension $[dL_e \times 1]$. That is, we define \mathbf{G}_{z,λ_e} with respect to the *extended* $[L_e \times 1]$ vector of factor loadings $\boldsymbol{\lambda}_{i,e}$. Here we define $\boldsymbol{\lambda}_{i,e} \equiv \mathbf{R}\boldsymbol{\lambda}_i$, where \mathbf{R} is the selection matrix of the form $\mathbf{F} = \mathbf{F}_e\mathbf{R}$ (the existence of \mathbf{R} is guaranteed by part (d)). Moreover, denote by $\boldsymbol{\theta} = (\boldsymbol{\beta}', \mathbf{g}')'$ the full parameter vector, and the corresponding true parameter vector by $\boldsymbol{\theta}_0$. It is worth mentioning that \mathbf{G}_{z,λ_e} and \mathbf{A}_N cannot be separately identified due to the usual rotation problem in factor models.

Using the plug-in principle and replacing \mathbf{F} by $\hat{\mathbf{F}}_e$ in Eq. (3), we define the following set of ζ estimating equations for $\boldsymbol{\theta}$:

$$\bar{\boldsymbol{\mu}}_N(\boldsymbol{\theta}) \equiv \frac{1}{N} \sum_{i=1}^N \mathbf{Z}'_i (\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta}) - \mathbf{S} \left(\hat{\mathbf{F}}_e \otimes \mathbf{I}_d \right) \mathbf{g}. \quad (4)$$

The GMM estimator is defined as the minimizer of the following objection function:

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \bar{\boldsymbol{\mu}}_N(\boldsymbol{\theta})' \boldsymbol{\Omega}_N \bar{\boldsymbol{\mu}}_N(\boldsymbol{\theta}), \quad (5)$$

where $\boldsymbol{\Omega}_N$ is some pre-specified positive definite matrix. Notice that $\bar{\boldsymbol{\mu}}_N(\boldsymbol{\theta})$ is linear in $\boldsymbol{\theta}$, and thus the estimator $\hat{\boldsymbol{\theta}}$ has a closed form solution.

The asymptotic distribution of the GMM estimator is determined primarily by the leading term in Eq. (4). In particular, $\bar{\boldsymbol{\mu}}_N(\boldsymbol{\theta})$ can be expanded as follows:

$$\bar{\boldsymbol{\mu}}_N(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N \boldsymbol{\mu}_i(\boldsymbol{\theta}) + o_P(N^{-1/2}), \quad (6)$$

where

$$\boldsymbol{\mu}_i(\boldsymbol{\theta}) = \mathbf{Z}'_i (\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta}) - \mathbf{S} \left((\mathbf{F}_e\mathbf{A}_N + \mathbf{\Psi}_i) \otimes \mathbf{I}_d \right) \mathbf{g}. \quad (7)$$

The following assumption imposes appropriate regularity conditions on $\bar{\boldsymbol{\mu}}_N(\boldsymbol{\theta})$.

Assumption 3.2.

- (a) $\mathbf{\Gamma}_\beta \equiv \text{plim}_{N \rightarrow \infty} - [\partial \bar{\boldsymbol{\mu}}_N(\boldsymbol{\theta}) / \partial \boldsymbol{\beta}'] = E_{\mathcal{F}}[\mathbf{Z}'_i \mathbf{X}_i]$ is \mathcal{F} -measurable and has full column rank a.s.
- (b) $\mathbf{\Gamma}_g \equiv \text{plim}_{N \rightarrow \infty} - [\partial \bar{\boldsymbol{\mu}}_N(\boldsymbol{\theta}) / \partial \mathbf{g}'] = \mathbf{S}(\mathbf{F}_e \mathbf{A} \otimes \mathbf{I}_d)$ is \mathcal{F} -measurable and has full column rank a.s.
- (c) The full-parameter Jacobian matrix $\mathbf{\Gamma} \equiv (\mathbf{\Gamma}_\beta, \mathbf{\Gamma}_g)$ has full column rank a.s.
- (d) $\boldsymbol{\Delta} \equiv \text{plim}_{N \rightarrow \infty} N^{-1} \sum_{i=1}^N \boldsymbol{\mu}_i(\boldsymbol{\theta}_0) \boldsymbol{\mu}_i(\boldsymbol{\theta}_0)'$ is \mathcal{F} -measurable and has full column rank a.s.

It is worth noting that Assumptions 2.1 and 3.1 are sufficient to ensure convergence in probability of the matrices defined in Assumption 3.2. Thus, the only non-trivial restrictions imposed in Assumption 3.2 are the rank restrictions. Part (a) is a standard identification condition in IV estimation and requires the instruments to be correlated with the regressors. Part (b) requires that $\mathbf{F}_e \mathbf{A}_N$ has full column rank, which is already implied by Assumption 3.1. Violations of this restriction are examined in Section 3.4. Part (c) assumes that $\mathbf{\Gamma}_\beta$ and $\mathbf{\Gamma}_g$ are linearly independent. That is, $E_{\mathcal{F}}[\mathbf{Z}'_i \mathbf{X}_i]$ cannot lie in the column space of \mathbf{F}_e . As an example, part (c) excludes situations where \mathbf{Z}_i and/or \mathbf{X}_i have degenerate idiosyncratic components with variance that is local-to-zero. Essentially, part (c) is the GMM analogue of the generalized non-collinearity condition of least-squares based factor estimates, as per Bai (2009) and Moon and Weidner (2015). Lastly, part (d) is also a standard condition and ensures that point-identified inference is asymptotically valid.

Remark 1. The Jacobian matrix $\mathbf{\Gamma}$ has ζ rows and $K + dL_e$ columns. Therefore, Assumption 3.2 requires that $\zeta \geq K + dL_e$. To illustrate the meaning of this requirement, suppose that all elements of \mathbf{z}_i are strictly exogenous, such that the largest possible set of (internal) instruments is given by $\zeta = dT$. A necessary condition for identification is that $d(T - L_e) \geq K$, which means that the number of factor proxies, L_e , should be strictly smaller than the number of time periods, T . Similar conclusions apply when some of the elements in \mathbf{z}_i are weakly exogenous or endogenous, except in this case \mathbf{g} is not identifiable without additional normalizations, see Section 3.3 in Juodis and Sarafidis (2018) for details.

The following theorem summarises the properties of the proposed estimator.

Theorem 1. *Suppose that Assumptions 2.1, 3.1 and 3.2 hold true. Then for $N \rightarrow \infty$,*

$$\sqrt{N} \left(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0 \right) \xrightarrow{d} \left[(\mathbf{\Gamma}' \boldsymbol{\Omega} \mathbf{\Gamma})^{-1} \mathbf{\Gamma}' \boldsymbol{\Omega} \right] \boldsymbol{\Delta}^{1/2} \boldsymbol{\pi} \quad (\mathcal{F} - \text{stably}), \quad (8)$$

where $\mathbf{\Omega}$ is some \mathcal{F} -measurable matrix such that $\text{plim}_{N \rightarrow \infty} \mathbf{\Omega}_N = \mathbf{\Omega}$, while $\mathbf{\Gamma}$, $\mathbf{\Omega}$, and $\mathbf{\Delta}$ are independent of $\boldsymbol{\pi} \sim N(\mathbf{0}_\zeta, \mathbf{I}_\zeta)$.

Proof. See the Supplementary Appendix. \square

Theorem 1 adopts the notion of \mathcal{C} -stable convergence, introduced by Kuersteiner and Prucha (2013), and characterises convergence as \mathcal{F} -stable. Hence, the GMM estimator $\widehat{\boldsymbol{\theta}}$ is consistent, and asymptotically mixed-normal. While this is an important distinction between the properties of the proposed estimator and those of Robertson and Sarafidis (2015) (who treat factors as fixed), it plays no role for inference procedures based on standardised statistics, as long as $\mathbf{\Gamma}$ and $\mathbf{\Delta}$ can be consistently estimated from their sample analogues. The result of Theorem 1 is general enough to establish general stable convergence, but we present the result as \mathcal{F} -stable because we wish to emphasize measurability of all random matrices with respect to \mathcal{F} .

As our focus lies on asymptotically linear estimators of \mathbf{F}_e , consistent estimation of $\mathbf{\Delta}$ requires a plug-in estimator of $\boldsymbol{\Psi}_i$. In particular, if such estimator $\widehat{\boldsymbol{\Psi}}_i$ is available, $\mathbf{\Delta}$ can be estimated consistently using the conventional formula

$$\widehat{\mathbf{\Delta}} = \frac{1}{N} \sum_{i=1}^N \widehat{\boldsymbol{\mu}}_i(\widehat{\boldsymbol{\theta}}) \widehat{\boldsymbol{\mu}}_i(\widehat{\boldsymbol{\theta}})', \quad (9)$$

where $\widehat{\boldsymbol{\mu}}_i(\boldsymbol{\theta})$ is the feasible plug-in estimate of $\boldsymbol{\mu}_i(\boldsymbol{\theta})$, i.e.

$$\widehat{\boldsymbol{\mu}}_i(\boldsymbol{\theta}) = \mathbf{Z}'_i(\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta}) - \mathbf{S} \left((\widehat{\mathbf{F}}_e + \widehat{\boldsymbol{\Psi}}_i) \otimes \mathbf{I}_d \right) \mathbf{g}. \quad (10)$$

Finally, the usual two-step GMM estimator can be obtained by setting $\mathbf{\Omega}_N = \widehat{\mathbf{\Delta}}^{-1}$.

3.2. Construction of Factor Proxies, $\widehat{\mathbf{F}}_e$

This section puts forward two specific methods for constructing $\widehat{\mathbf{F}}_e$. These methods are motivated by both our empirical application and common practice in the large- T panel data literature.

Method I. One variable and multiple weights. Suppose there exists a single variable $\mathbf{v}_i = (v_{i,1}, \dots, v_{i,T})'$ driven by \mathbf{F} , as well as possibly additional factors. That is, \mathbf{v}_i satisfies

$$\mathbf{v}_i = \mathbf{F}_e \boldsymbol{\gamma}_i + \mathbf{u}_i, \quad (11)$$

where \mathbf{F}_e is $[T \times L_e]$ with $L_e \geq L$, such that $\mathbf{F} \in \text{Col}(\mathbf{F}_e)$. \mathbf{v}_i can be either internal, i.e. one of the regressors, or external, in the spirit of Pesaran et al. (2013). The existence of such variable is quite plausible in panel data models (see e.g. Hansen and Liao

2018 and Karabiyik et al. 2019) because economic agents inhabit common economic environments and therefore many variables are often subject to common shocks, such as changes in technology and productivity, changes in preferences and tastes, and so on.

Let \mathbf{w}_i denote an $[L_e \times 1]$ vector of individual-specific weights, such that $\mathbb{E}_{\mathcal{F}}[\mathbf{u}_i \mathbf{w}'_i] = \mathbf{O}_{T \times L_e}$ and $\mathbb{E}_{\mathcal{F}}[\gamma_i \mathbf{w}'_i] = \mathbf{G}_{\gamma, w}$ a.s. In terms of the general notation employed in Assumption 3.1, this setup corresponds to setting $\mathbf{A} = \mathbf{G}_{\gamma, w}$ and $\boldsymbol{\Psi}_i = \mathbf{u}_i \mathbf{w}'_i + \mathbf{F}_e(\gamma_i \mathbf{w}'_i - \mathbf{G}_{\gamma, w})$. Thus, Assumption 3.1 translates into the requirement that $\text{rk}(\mathbf{G}_{\gamma, w}) = L_e$. In this case,

$$\widehat{\mathbf{F}}_e = \frac{1}{N} \sum_{i=1}^N \mathbf{v}_i \mathbf{w}'_i, \quad (12)$$

is a suitable estimator of \mathbf{F}_e . Furthermore, the corresponding plug-in estimator of $\widehat{\boldsymbol{\Psi}}_i$ is simply given by

$$\widehat{\boldsymbol{\Psi}}_i = \mathbf{v}_i \mathbf{w}'_i - \widehat{\mathbf{F}}_e. \quad (13)$$

It is apparent that Method I permits cases where the L_e factor proxies in $\widehat{\mathbf{F}}_e$ identify more factors than those entering into the error term of the equation for y . Such generality is appealing because \mathbf{v}_i may contain more factors than those that already drive $y_{i,t}$.

Example 1. Suppose that $L = 1$, such that the factor component in Eq. (2) reduces to $\mathbf{f}^{(1)} \lambda_i^{(1)}$, but \mathbf{v}_i is driven by two factors, i.e. $\mathbf{F}_e = (\mathbf{f}^{(1)}, \mathbf{f}^{(2)})$ and \mathbf{w}_i is a $[2 \times 1]$ vector. Then,

$$\widehat{\mathbf{F}}_e = \frac{1}{N} \sum_{i=1}^N \mathbf{v}_i \mathbf{w}'_i \xrightarrow{p} \mathbf{F}_e \mathbf{A}, \quad (14)$$

where $\mathbf{A} = \mathbb{E}_{\mathcal{F}}[\gamma_i \mathbf{w}'_i]$. Consistent estimability of $\boldsymbol{\beta}$ requires that $\text{rk}(\mathbf{A}) = 2$. Notice that if one makes use of one weight only, w_i , $\mathbf{A} = \mathbb{E}_{\mathcal{F}}[\gamma_i w_i]$ has at most rank 1. Thus, in this case $\widehat{\mathbf{F}}_e$ cannot estimate $\mathbf{f}^{(1)}$ consistently, unless $\mathbb{E}_{\mathcal{F}}[\gamma_i^{(2)} w_i] = 0$, where $\gamma_i^{(2)}$ denotes the bottom entry in $\boldsymbol{\gamma}_i$. Section 3.4 discusses in detail the situation where the full rank restriction on \mathbf{A} is violated.

The aforescribed method requires two distinct ingredients: \mathbf{v}_i and \mathbf{w}_i . There are several potential choices for \mathbf{w}_i . Consider initially the case $L_e = L = 1$ and suppose that the factor component in Eq. (11) reduces to $\mathbf{f}^{(1)} \gamma_i^{(1)}$. One simple choice is to set the value of the weight across all individuals equal to a fixed constant, i.e. $w_i = 1$. In this case, Eq. (12) becomes the cross-sectional average of \mathbf{v}_i , whereas $\mathbf{G}_{\gamma, w}$ reduces

to $\mu_\gamma \equiv E_{\mathcal{F}} [\gamma_i^{(1)}]$. Thus, Assumption 3.1 implies that $\mu_\gamma \neq 0$ and $E_{\mathcal{F}}[\mathbf{u}_i] = \mathbf{0}_T$. An alternative choice, and one which is especially appealing in autoregressive models, is to use random weights that are independent across i conditional on \mathcal{F} , such as lagged values of the observed data. For instance, in ARDL(p,q) models a natural choice is to set $w_i = y_{i,0}$ (or $w_i = x_{i,0}^{(k)}$), which requires $E_{\mathcal{F}}[\gamma_i^{(1)} y_{i,0}] \neq 0$ (or $E_{\mathcal{F}}[\gamma_i^{(1)} x_{i,0}^{(k)}] \neq 0$), as in the “correlated random effects” framework.

When $L_e > 1$, the use of a single weight would violate Assumption 3.1, as was pointed out in Example 1. Thus, a possible “automated” strategy for constructing the $[L_e \times 1]$ dimensional vector \mathbf{w}_i is to pick the first $L_e - 1$ observations of a single regressor, such that $\mathbf{w}_i = \left(1, \mathbf{x}_{i,1:L-1}^{(k)}\right)$ for some k . Alternatively, one can rely upon the initial condition of dependent and independent variables, i.e. $\mathbf{w}_i = \left(1, y_{i,0}, x_{i,0}^{(1)}, \dots, x_{i,0}^{(L_e-2)}\right)'$. Lastly, one can use powers of $y_{i,0}$ or $x_{i,0}^{(k)}$, such as $y_{i,0}^2$ and so on. This option does not require the distribution of $y_{i,t}$ to be symmetric because non-central moments are used. To illustrate the aforementioned strategies, let $K = 2$, $L = 3$, $x_{i,t}^{(1)} \equiv y_{i,t-1}$, and let $x_{i,t}^{(2)} \equiv x_{i,t}$ be treated as weakly exogenous. For this specification, one could set $\mathbf{w}_i = (1, y_{i,0}, y_{i,1})'$, $\mathbf{w}_i = (1, y_{i,0}, x_{i,0})'$, or $\mathbf{w}_i = (1, y_{i,0}, y_{i,0}^2)'$.

Since the parameters of interest can be estimated based on different weighting schemes, the approach proposed in this paper provides a more flexible way of dealing with unobserved factors compared to other methods that are available in the literature. The practical question of how to actually select among different values of \mathbf{w}_i (and/or different v 's) is discussed in Section 3.5.

Remark 2. Independently from this research, Gagliardini and Gouriéroux (2017) and Fan and Liao (2019) have recently advocated a similar construction of factor proxies, which involves pre-specified (potentially arbitrary) weights \mathbf{w}_i . Unlike our study, the prime focus of those studies lies on the asymptotic properties of factor estimates when both N and T are large.

Method II. Multiple variables, single weight. The second approach involves the construction of factor proxies based on a pre-specified *single* weighting scheme, w_i , and L_e distinct time-varying variables. As before, these variables can be either internal or external. In particular, let $\mathbf{V}_i = (\mathbf{v}_i^{(1)}, \dots, \mathbf{v}_i^{(L_e)})$ be a $[T \times L_e]$ matrix, such that

$$\mathbf{V}_i = \mathbf{F}_e \boldsymbol{\Upsilon}_i + \mathbf{U}_i, \tag{15}$$

where $\boldsymbol{\Upsilon}_i = (\boldsymbol{\gamma}_i^{(1)}, \dots, \boldsymbol{\gamma}_i^{(L_e)})$ and $\mathbf{U}_i = (\mathbf{u}_i^{(1)}, \dots, \mathbf{u}_i^{(L_e)})$. In this case, $\widehat{\mathbf{F}}_e$ is defined as:

$$\widehat{\mathbf{F}}_e = \frac{1}{N} \sum_{i=1}^N \mathbf{V}_i w_i = \frac{1}{N} \sum_{i=1}^N (\mathbf{F}_e \boldsymbol{\Upsilon}_i w_i + \mathbf{U}_i w_i). \quad (16)$$

Notice that the use of a single weight requires that such weight is valid for each column of \mathbf{V}_i , i.e. $\mathbb{E}_{\mathcal{F}}[\mathbf{U}_i w_i] = \mathbf{O}_{T \times L_e}$. However, the above formulation does not imply that each column in \mathbf{V}_i must be driven by the full set of columns in \mathbf{F}_e . For instance, let $L = L_e = 2$, $\mathbf{F} = \mathbf{F}_e = (\mathbf{f}^{(1)}, \mathbf{f}^{(2)})$ with $\mathbf{v}_i^{(1)} = \mathbf{f}^{(1)} \boldsymbol{\gamma}_i^{(1)} + \mathbf{u}_i^{(1)}$ and $\mathbf{v}_i^{(2)} = \mathbf{f}^{(2)} \boldsymbol{\gamma}_i^{(2)} + \mathbf{u}_i^{(2)}$. This structure can be represented by Eq. (15) using a diagonal $\boldsymbol{\Upsilon}_i$ matrix.

We note that the setup in Method II is sufficiently general in that it accommodates cases where the number of factor proxies required to approximate \mathbf{F} is larger than L itself. The following example illustrates this point.

Example 2. Consider an autoregressive, AR(1), model with $L = 1$:

$$y_{i,t} = \alpha y_{i,t-1} + \lambda_i f_t + \varepsilon_{i,t}; \quad t = 1, \dots, T, \quad (17)$$

or, in vector form,

$$\mathbf{y}_i = \alpha \mathbf{y}_{i,-1} + \mathbf{f} \lambda_i + \boldsymbol{\varepsilon}_i, \quad (18)$$

where $\mathbf{y}_{i,-1} = (y_{i,0}, \dots, y_{i,T-1})'$, $\mathbf{f} = (f_1, \dots, f_T)'$ and \mathbf{y}_i is defined below Eq. (2). As it is shown in Everaert and De Groot (2016) and Juodis et al. (2020), the CCE-style factor proxies, $\widehat{\mathbf{F}}_e = (\overline{\mathbf{y}}, \overline{\mathbf{y}}_{-1})$, which arise in our context by setting $\mathbf{V}_i = (\mathbf{y}_i, \mathbf{y}_{i,-1})$ and $w_i = 1$, satisfy Assumption 3.1 provided that $\mathbb{E}_{\mathcal{F}}[\lambda_i] \neq 0$. In this setup, none of the two factor proxies in $\widehat{\mathbf{F}}_e$ alone is able to estimate \mathbf{f} consistently; instead only a linear combination of both columns can be estimated consistently. In particular, if we define $\mathbf{F}_e \equiv \mathbb{E}_{\mathcal{F}}[\mathbf{V}_i]$, then $\mathbf{F}_e = (\mathbf{f}_e^{(1)}, \mathbf{f}_e^{(2)})$ with $\mathbf{f}_e^{(1)} = \alpha \mathbf{f}_e^{(2)} + \mathbb{E}_{\mathcal{F}}[\lambda_i] \mathbf{f}$, while $\mathbf{f}_e^{(2)}$ is a cumulative function of the lags of \mathbf{f} and the initial condition, i.e. $f_{1,e}^{(2)} = \mathbb{E}_{\mathcal{F}}[y_{i,0}]$, and $f_{t,e}^{(2)} = \mathbb{E}_{\mathcal{F}}[\lambda_i] f_{t-1} + \alpha f_{t-1,e}^{(2)}$ for $t = 2, \dots, T$. Notice that in this case $\mathbf{R} = (1, -\alpha)' / \mathbb{E}_{\mathcal{F}}[\lambda_i]$, and the extended vector of factor loadings $\boldsymbol{\lambda}_{i,e}$ is defined accordingly. As has been discussed, alternative choices for $w_i = 1$ do exist. For instance, when $\mathbb{E}_{\mathcal{F}}[\lambda_i] = 0$, a plausible strategy is to set $w_i = y_{i,0}$ and drop the first observation from estimation of the AR(1) model in Eq. (17).

Remark 3. Method II is motivated by, although it is not nested within, the seminal CCE approach of Pesaran (2006). In particular, CCE considers deterministic weights (mostly $w_i = 1$), and assumes that *all* K regressors, \mathbf{X}_i , are strictly exogenous with

respect to $\boldsymbol{\varepsilon}_i$. In contrast, here w_i can be stochastic, as Example 2 points out. Moreover, not all regressors need to be strictly exogenous. Finally, the restriction $E_{\mathcal{F}}[\lambda_i] \neq 0$ is testable within our framework, see Section 3.4. In the Supplementary Appendix we study a stylised setup in which both the proposed linear GMM approach and CCE are consistent for T fixed; specifically, we focus on a model with one strictly exogenous regressor, which is driven by the same (single) factor that enters into the error term of the equation for y , such that $E_{\mathcal{F}}[\gamma_i] \neq 0$. We demonstrate that under spherical error components, the CCE estimator and the linear GMM estimator based on $w_i = 1$ have the same asymptotic variance.

Remark 4. Depending on the model at hand, one can also combine multiple time-varying variables, \mathbf{V}_i , with multiple time-invariant weights, \mathbf{w}_i , in order to construct $\widehat{\mathbf{F}}_e$. As an example, let $\mathbf{V}_i = (\mathbf{v}_i^{(1)}, \mathbf{v}_i^{(2)})$ and $\mathbf{w}_i = (w_i^{(1)}, w_i^{(2)})'$, in which case potential factor proxies can be constructed by using all possible combinations of \mathbf{V}_i and \mathbf{w}_i , i.e.

$$\widehat{\mathbf{F}}_e = \frac{1}{N} \sum_{i=1}^N \mathbf{V}_i \otimes \mathbf{w}_i'. \quad (19)$$

Alternatively, one can consider only a subset of available proxies. For instance,

$$\widehat{\mathbf{F}}_e = \frac{1}{N} \sum_{i=1}^N (\mathbf{v}_i^{(1)} w_i^{(1)}, \mathbf{v}_i^{(2)} w_i^{(2)}). \quad (20)$$

This flexibility provides a further advantage of the factor-proxies formulation used in this paper over that of Pesaran (2006). The practical question of how to actually select among different values of \mathbf{w}_i (and/or different values of \mathbf{v}_i) is discussed in Section 3.5.

Irrespective of the method that is considered to construct factor proxies, notice that $\widehat{\mathbf{F}}_e$ can be expressed as

$$\widehat{\mathbf{F}}_e = \mathbf{F}_e \mathbf{A} + \overline{\boldsymbol{\Psi}}; \quad \overline{\boldsymbol{\Psi}} = \frac{1}{N} \sum_{i=1}^N \boldsymbol{\Psi}_i, \quad (21)$$

where here $\boldsymbol{\Psi}_i$ is defined as

$$\boldsymbol{\Psi}_i = (\mathbf{F}_e (\mathbf{A}_i - \mathbf{A}) + \mathbf{E}_i), \quad (22)$$

for some \mathbf{A}_i . Thus, $\widehat{\mathbf{F}}_e$ is a linear estimator of the column space of \mathbf{F}_e . In what follows we put forward an estimator that is only asymptotically linear, but not linear for fixed N and T .

3.3. Regularized Factor Proxies

This section analyses the important case where the model is “fundamentally identified”, but the practitioner includes more weights, or more variables than necessary in the approximation of \mathbf{F}_e . In the context of Assumption 3.1, this implies that part (c) is violated but (d) is not.

To formalise this idea, let $R \geq L_e$ be the total number of factor proxies and $\widehat{\mathbf{F}}_R$ denote the corresponding factor proxies with dimension $[T \times R]$. Following Eq. (21), irrespective of the method considered to construct factor proxies, $\widehat{\mathbf{F}}_R$ can be decomposed as

$$\widehat{\mathbf{F}}_R = \mathbf{F}_e \mathbf{A}_R + \overline{\boldsymbol{\Psi}}^{(R)}; \quad \overline{\boldsymbol{\Psi}}^{(R)} = \frac{1}{N} \sum_{i=1}^N \boldsymbol{\Psi}_i^{(R)}, \quad (23)$$

where $\overline{\boldsymbol{\Psi}} = (\overline{\boldsymbol{\psi}}_1^{(R)}, \dots, \overline{\boldsymbol{\psi}}_T^{(R)})'$ is $[T \times R]$, and \mathbf{A}_R is $[L_e \times R]$ with rank L_e . As a result, $\text{rk}(\text{E}_{\mathcal{F}}[\widehat{\mathbf{F}}_R]) = L_e$ and $\boldsymbol{\Gamma}_g$, defined in Assumption 3.2, is of reduced rank.

Example 3. Consider the following model in vector form:

$$\mathbf{y}_i = \beta_1 \mathbf{x}_i^{(1)} + \beta_2 \mathbf{x}_i^{(2)} + \mathbf{f}^{(1)} \lambda_i^{(1)} + \boldsymbol{\varepsilon}_i, \quad (24)$$

where $\mathbf{f}^{(1)}$ is $[T \times 1]$. Thus, $K = 2$ and $L = 1$. Let

$$\mathbf{x}_i^{(1)} = \mathbf{f}^{(1)} \pi_i + \mathbf{u}_i^{(1)}, \quad (25)$$

and $L_e = 2$, such that

$$\mathbf{x}_i^{(2)} = \mathbf{F}_e \boldsymbol{\xi}_i + \mathbf{u}_i^{(2)}, \quad (26)$$

where $\mathbf{F}_e = (\mathbf{f}^{(1)}, \mathbf{f}^{(2)})$ is of order $[T \times 2]$, and $\boldsymbol{\xi}_i$ is a vector of order $[2 \times 1]$. One possible strategy for proceeding is to proxy \mathbf{F}_e using the full set of observables, by setting $\mathbf{V}_i = (\mathbf{x}_i^{(1)}, \mathbf{x}_i^{(2)}, \mathbf{y}_i)$ and $w_i = 1$. This is standard practice in CCE estimation, with the aim of avoiding the need to estimate the number of factors in the model. In terms of the notation used in Eq. (15), we have $\boldsymbol{\Upsilon}_i = (\gamma_i^{(1)}, \gamma_i^{(2)}, \gamma_i^{(3)})$, which is a $[2 \times 3]$ matrix with $\gamma_i^{(1)} = (\pi_i, 0)'$, $\gamma_i^{(2)} = \boldsymbol{\xi}_i$ and $\gamma_i^{(3)} = \boldsymbol{\lambda}_i + \beta_1 \gamma_i^{(1)} + \beta_2 \boldsymbol{\xi}_i$, where $\boldsymbol{\lambda}_i = (\lambda_i^{(1)}, 0)'$. In this case, $\text{E}_{\mathcal{F}}[\boldsymbol{\Upsilon}_i w_i]$ does not have full column rank. Essentially, too many variables are used for approximating \mathbf{F}_e . In this section we put forward a method that overcomes this problem within our GMM approach.

Under these circumstances, it is straightforward to show that the GMM estimator considered thus far remains consistent. However, it turns out that the asymptotic distribution of $\widehat{\boldsymbol{\beta}}$ can be highly non-standard, due to a degeneracy of the \mathbf{g} estimates.

For this reason, it is essential to use factor proxies that are non-degenerate, that is, none of the columns of $\widehat{\mathbf{F}}_e$ are asymptotically collinear.

In what follows, we put forward a regularization approach for constructing $\widehat{\mathbf{F}}_e$ such that parts (a)-(c) of Assumption 3.1 are satisfied. Our regularization method uses the singular value decomposition of $\widehat{\mathbf{F}}_R$, or, equivalently, the principal components of $\widehat{\mathbf{F}}_R \widehat{\mathbf{F}}_R'$.

To begin with, let \mathbf{V}_N be a $[L_e \times L_e]$ diagonal matrix containing the L_e largest eigenvalues of $T^{-1} \widehat{\mathbf{F}}_R \widehat{\mathbf{F}}_R'$ in descending order. The following assumption ensures that the asymptotic distribution of the proposed regularized estimator is well defined, see e.g. Bai (2003) for a similar condition:

Assumption 3.3. *The eigenvalues of the $[L_e \times L_e]$ matrix $(\mathbf{A}_R \mathbf{A}_R') (\mathbf{F}_e' \mathbf{F}_e)$ are distinct a.s.*

Let $\widetilde{\mathbf{F}} = \sqrt{T} \widehat{\mathbf{U}}_{L_e}$ denote the scaled regularized, principal components (PC) estimator for \mathbf{F}_e , where $\widehat{\mathbf{U}}_{L_e}$ denotes the associated eigenvectors (left singular vectors) corresponding to the L_e largest eigenvalues of $T^{-1} \widehat{\mathbf{F}}_R \widehat{\mathbf{F}}_R'$. The fixed T consistency of the PC factor estimator $\widetilde{\mathbf{F}}$, which is of dimension $[T \times L_e]$, follows intuitively from the results in Connor and Korajczyk (1986) and Bai (2003). As it is emphasized in Bai (2003), a necessary and sufficient condition for fixed T consistency is that $\widehat{\mathbf{F}}_R \widehat{\mathbf{F}}_R' \xrightarrow{p} \mathbf{F}_e \mathbf{M} \mathbf{F}_e' + \kappa \mathbf{I}_T$, for some matrix \mathbf{M} and some scalar $\kappa \geq 0$. In our case, this necessary condition is satisfied trivially with $\kappa = 0$. However, a new proof for the asymptotic distribution of $\widetilde{\mathbf{F}} = (\widetilde{\mathbf{f}}_1, \dots, \widetilde{\mathbf{f}}_T)'$ is required when T is fixed. This is achieved by Theorem 2 below.

Theorem 2. *Suppose that Assumption 3.3 holds true, and $\widehat{\mathbf{F}}_R$ satisfies Eq. (21) together with Assumption 3.1, except that $\text{rk}(\mathbf{A}_R) = L_e < R$. Then the regularised factor estimator $\widetilde{\mathbf{F}}$ satisfies Assumption 3.1, with*

$$\mathbf{A}_N = (\mathbf{A}_R \mathbf{A}_R') \left(\frac{\widetilde{\mathbf{F}}' \mathbf{F}_e}{T} \right) \mathbf{V}_N^{-1}. \quad (27)$$

Furthermore, each row $\Psi_i^{(t)}$ in Ψ_i is of the following form:

$$\Psi_i^{(t)} = \mathbf{V}_N^{-1} \frac{1}{T} \left(\sum_{s=1}^T \widetilde{\mathbf{f}}_s \left(\mathbf{f}'_{s,e} \mathbf{A}_R \psi_{i,t}^{(R)} + \mathbf{f}'_{t,e} \mathbf{A}_R \psi_{i,s}^{(R)} \right) \right); \quad t = 1, \dots, T. \quad (28)$$

Proof. See the Supplementary Appendix. \square

It is straightforward to see from expression Eq. (28) that the leading variance term of $\widetilde{\mathbf{F}}$ is affine in each $\{\mathbf{A}_R \psi_{i,t}\}_{t=1}^T$. This implies that the inclusion of *uninformative* factor

proxies, i.e. proxies that do not identify \mathbf{F}_e , has no impact on the first-order asymptotic properties of $\tilde{\mathbf{F}}$. Essentially, the PC estimator of factors performs estimation of the factor space and factor proxy selection at the same time.

Since $\tilde{\mathbf{F}}$ satisfies Assumption 3.1, it can be used as a plug-in estimator in the GMM objective function. Consistent estimation of the variance-covariance matrix $\mathbf{\Delta}$, requires replacing unknown quantities in Ψ_i with their plug-in counterparts. In particular, the $\mathbf{f}'_{s,e}\mathbf{A}_R$ terms in Eq. (28) can be consistently estimated by $(\hat{\mathbf{f}}_s^{(R)})'$. Furthermore, depending on the method used to construct the factor proxies, the plug-in counterparts of $\psi_{i,t}^{(R)}$ are given by either

$$\hat{\psi}_{i,t}^{(R)} = \mathbf{w}_i \mathbf{v}_{i,t} - \hat{\mathbf{f}}_t^{(R)}, \quad (29)$$

or

$$\hat{\psi}_{i,t}^{(R)} = \mathbf{v}_{i,t} \mathbf{w}_i - \hat{\mathbf{f}}_t^{(R)}, \quad (30)$$

where $\mathbf{v}_{i,t} = (\mathbf{v}_{i,t}^{(1)}, \dots, \mathbf{v}_{i,t}^{(L_e)})'$.

Remark 5. The above strategy for selection of factor proxies is not optimal from the point of view of obtaining a GMM estimator for β_0 with minimal asymptotic variance. Specifically, the PC estimator is a weighted average of all individual factor proxies, with corresponding weights being determined outside the GMM objective function. One could proxy \mathbf{F}_e by combining factor proxies optimally. However, such an approach has a major drawback, in that the resulting moment conditions are nonlinear. To illustrate this, consider a linear combination of $\mathbf{w}_i = (w_i^{(1)}, w_i^{(2)})'$ that takes the form $\tilde{w}_i = a w_i^{(1)} + (1-a) w_i^{(2)}$. It is clear that in the absence of knowledge of a , the resulting moment conditions involve products of unknown parameters, namely \mathbf{A}_R and a . Hence, the appealing linearity of the proposed approach no longer holds, and instead one may use e.g. the FIVU and FIVR estimators of Robertson and Sarafidis (2015), which are asymptotically more efficient because they involve joint estimation of β and \mathbf{F} using the full set of moment conditions.

3.4. Identification

Assumption 3.1 plays a major role in characterizing the large sample properties of the proposed GMM estimator. In this section, we discuss several departures from Assumption 3.1, as well as diagnostic checks that can be used to detect these departures. To save space, we focus on factor proxies constructed using multiple \mathbf{w}_i , as in Eq. (12).

At first, consider the case where $\mathbf{F} \notin \text{Col}(\mathbf{F}_e)$ but otherwise all remaining parts of Assumption 3.1 are satisfied. For instance, suppose that some of the factors in \mathbf{y}_i are

entirely different from those that drive \mathbf{v}_i , or alternatively $L_e < L$. In neither case can the factor proxies approximate the column space of \mathbf{F} asymptotically. As a result, it is straightforward to show that the GMM estimator for β_0 is inconsistent. That is,

$$\widehat{\beta} - \beta_0 \xrightarrow{p} \left(\Gamma'_\beta \Gamma_\beta - \Gamma'_\beta \Gamma_g (\Gamma'_g \Gamma_g)^{-1} \Gamma'_g \Gamma_\beta \right)^{-1} \left(\Gamma'_\beta - \Gamma'_\beta \Gamma_g (\Gamma'_g \Gamma_g)^{-1} \Gamma'_g \right) \mathbf{S} \text{vec}(\mathbf{G}_{z,\lambda} \mathbf{F}'), \quad (31)$$

where Γ_β and Γ_g denote the two constituent blocks of the Jacobian matrix of the moment conditions, as defined in Assumption 3.2. Such identification failure can be detected using the usual overidentifying restrictions test, commonly referred to as J-statistic.

Next, consider the case where $\mathbf{G}_{\gamma,w}$ is not of full rank, such that $\mathbf{F} \notin \text{Col}(\mathbf{F}_e \mathbf{G}_{\gamma,w})$. As an example, let $L_e = L = 1$ and $w_i = 1$, such that $\text{E}_{\mathcal{F}}(\gamma_i w_i) = \mu_\gamma$. The full rank condition on $\mathbf{G}_{\gamma,w}$ is violated when $\mu_\gamma = 0$. More generally, suppose that $\text{rk}(\mathbf{G}_{\gamma,w}) = Q < L_e \leq L$ a.s. Let $\mathbf{G}_{\gamma,w} = \mathbf{C} \mathbf{D}'$, where both \mathbf{C} and \mathbf{D} are $[L_e \times Q]$ matrices of rank Q a.s. Furthermore, let \mathbf{D}_\perp denote the orthogonal complement of \mathbf{D} , i.e. \mathbf{D}_\perp satisfies $\text{rk}(\mathbf{D}_\perp) = L_e - Q$ and $\mathbf{D}' \mathbf{D}_\perp = \mathbf{O}_{Q \times (L_e - Q)}$. Theorem 3 summarizes the asymptotic distribution of $\widehat{\beta}$ when $\mathbf{G}_{\gamma,w}$ is rank-deficient.

Theorem 3. *Suppose that Assumption 2.1 is satisfied, and consider factor proxies with $\mathbf{A}_N = \mathbf{G}_{\gamma,w}$ with $\text{rk}(\mathbf{G}_{\gamma,w}) = Q < L_e \leq L$. Then for $\Omega_N = \mathbf{I}$, as $N \rightarrow \infty$:*

$$\widehat{\beta} - \beta_0 \xrightarrow{d} \left(\Gamma'_\beta \Gamma_\beta - \Gamma'_\beta \Xi (\Xi' \Xi)^{-1} \Xi' \Gamma_\beta \right)^{-1} \left(\Gamma'_\beta - \Gamma'_\beta \Xi (\Xi' \Xi)^{-1} \Xi' \right) \mathbf{S} \text{vec}(\mathbf{G}_{z,\lambda} \mathbf{F}'), \quad (32)$$

\mathcal{F} -stably, where $\Xi = \mathbf{S}((\Xi_d, \Xi_s) \otimes \mathbf{I}_d)$, $\Xi_d = \mathbf{F}_e \mathbf{C} (\mathbf{D}' \mathbf{D})$, and Ξ_s is such that $\text{vec}(\Xi_s) \sim \Delta_F^{1/2} \boldsymbol{\pi}_F$ with $\Delta_F = \text{E}_{\mathcal{F}}[\text{vec}(\mathbf{v}_i \mathbf{w}'_i \mathbf{D}_\perp) \text{vec}(\mathbf{v}_i \mathbf{w}'_i \mathbf{D}_\perp)']$, and $\boldsymbol{\pi}_F \sim N(\mathbf{0}_{(L_e - Q)T}, \mathbf{I}_{(L_e - Q)T})$.

Proof. See the Supplementary Appendix. □

It is clear that the GMM estimator converges to a random limit because the Ξ_s matrix is stochastic. This result is in line with existing weak instruments results for IV/GMM estimators, see e.g. Staiger and Stock (1997). However, a major difference between Theorem 3 and existing literature is that in the usual weak-IV setup, the limit remains random even after conditioning on the (random) Jacobian matrix. In contrast, in the present case, the right-hand side of Eq. (32) is a non-zero constant vector, conditional on \mathcal{F} and $\boldsymbol{\pi}_F$. As a result, failure of identification associated with $\text{rk}(\mathbf{G}_{\gamma,w}) = Q < L_e \leq L$ a.s. implies model mis-specification. This can be detected using the usual J-statistic again.

Remark 6. The stochastic nature of the identification failure in Theorem 3 can be easily avoided using regularized factor proxies as in Section 3.3, obtained using a consistent estimate of L_e .

3.5. Implementation

We take it as given in the section that a set of potential factor proxies has been collected into a $[T \times R]$ matrix $\widehat{\mathbf{F}}_R$, using either Method I, Method II, or a combination of both. For example, under Method II one could construct $\widehat{\mathbf{F}}_R$ based on all available observables, i.e. the dependent variable, the regressors, as well as possible external variables. Unfortunately, as we discussed in Section 3.3, it turns out that including more variables (or, equivalently for Method I, more weights) than necessary in the approximation of \mathbf{F}_e , can render the asymptotic distribution of the GMM estimator highly non-standard. To circumvent this potential issue, in what follows we put forward two distinct methods that practitioners may use to implement the plug-in principle embedded in our approach. Firstly, a “regularization” method, then a “best-subset selection” method. Both methods are illustrated in the empirical section.

Regularization. This approach builds upon the regularized factor proxies analysed in Section 3.3 and consists of the following steps:

- (1) Obtain a consistent estimate of the underlying number of factors in $\widehat{\mathbf{F}}_R$, given by \widehat{L}_e , using either a sequential pivotal rank testing as proposed by Kleibergen and Paap (2006), or the eigenvalue ratio (ER) and the growth ratio (GR) statistics as in Ahn and Horenstein (2013).
- (2) Use the regularized estimator $\tilde{\mathbf{F}}$ as the plug-in estimator for the GMM objective function to estimate $\widehat{\boldsymbol{\theta}}$. If the model is not rejected by the J-statistic, no further steps are required.

Assuming that \widehat{L}_e is not large relative to \widehat{L} , the only disadvantage of this approach is ultimately some loss in terms of efficiency due to the fact that for $L_e > L$ one estimates the extended $\mathbf{G}_{z, \lambda_e}$, as e.g. in Example 2.

Remark 7. Let $r_{\max} = \min(T, R) - 1$. The estimator for L_e based on the ER-statistic is defined as

$$\widehat{L}_e = \arg \max_{r \in \{1, \dots, r_{\max}\}} ER(r); \quad ER(r) = \frac{\lambda_r(T^{-1} \widehat{\mathbf{F}}_R \widehat{\mathbf{F}}_R')}{\lambda_{r+1}(T^{-1} \widehat{\mathbf{F}}_R \widehat{\mathbf{F}}_R')}, \quad (33)$$

where $\lambda_r(T^{-1} \widehat{\mathbf{F}}_R \widehat{\mathbf{F}}_R')$ is the r^{th} largest eigenvalue of $T^{-1} \widehat{\mathbf{F}}_R \widehat{\mathbf{F}}_R'$. This procedure yields a consistent estimate of L_e because for $r < L_e$ $ER(r)$ remains bounded a.s., whereas

for $r = L_e$ $ER(r) \rightarrow \infty$ as $\lambda_{r+1}(T^{-1}\widehat{\mathbf{F}}_R\widehat{\mathbf{F}}_R') \xrightarrow{p} 0$ by the continuity of the eigenvalues. However, as it currently stands, the ER procedure excludes the possibility that $\widehat{L}_e = R < T$, i.e. the dimension of the column space of $\widehat{\mathbf{F}}_R$ is exactly R . We fix this shortcoming by borrowing the mock-eigenvalue idea of Ahn and Horenstein (2013) and suggesting that a single, redundant column in $\widehat{\mathbf{F}}_R$ always be included. For example, one can easily construct a redundant factor proxy as $N^{-1} \sum_{i=1}^N \mathbf{v}_i w_i^+$, where w_i^+ is randomly drawn from any zero mean distribution (e.g. the Rademacher $\{-1; 1\}$ distribution) and \mathbf{v}_i is either defined in Eq. (11), or it is one of the columns of \mathbf{V}_i in Eq. (15). In this way, we extend the definition of R to $R + 1$ manually, thus avoiding the boundary problem of the original ER-statistic. The GR-statistic is discussed in greater detail in the Supplementary Appendix.

Best-subset selection. The best-subset selection method is a model selection approach that is motivated by the machine learning literature, see e.g. Section 3.3. in Hastie et al. (2017). In the present context, the method aims to determine the combination of factor proxies that yields the smallest BIC value. In particular, let R be the number of factor proxies at hand, $\widehat{\mathbf{F}}_R$, and $L_{\max} (\geq L_e)$ be the maximum number of unobserved factors considered in estimation. In practice, L_{\max} could be set as the largest possible value of L that is feasible to allow in estimation. Furthermore, let $\mathcal{B}(L_{\max})$ denote the set of different combinations of columns in $\widehat{\mathbf{F}}_R$ of sizes $P = 1, \dots, L_{\max}$. To illustrate, consider Example 3 with $L_{\max} = 2$. Therein, $R = 3$ since $\widehat{\mathbf{F}}_R = N^{-1} \sum_{i=1}^N (\mathbf{x}_i^{(1)}, \mathbf{x}_i^{(2)}, \mathbf{y}_i) w_i \equiv (\widehat{\mathbf{f}}_R^{(1)}, \widehat{\mathbf{f}}_R^{(2)}, \widehat{\mathbf{f}}_R^{(3)})$. For $P = 1$, one can use either one of $\widehat{\mathbf{f}}_R^{(1)}$, $\widehat{\mathbf{f}}_R^{(2)}$ and $\widehat{\mathbf{f}}_R^{(3)}$. For $P = 2$, possible combinations of factor proxies are $\{\widehat{\mathbf{f}}_R^{(1)}, \widehat{\mathbf{f}}_R^{(2)}\}$, $\{\widehat{\mathbf{f}}_R^{(1)}, \widehat{\mathbf{f}}_R^{(3)}\}$, and $\{\widehat{\mathbf{f}}_R^{(2)}, \widehat{\mathbf{f}}_R^{(3)}\}$. In total, the cardinality of the set $\mathcal{B}(L_{\max})$ is at most $|\mathcal{B}(L_{\max})| = \binom{3}{1} + \binom{3}{2} = 6$.

Within our framework, a consistent estimate of the number of factors and the selected combination of factor proxies can be determined using a Schwartz/Bayesian model information criterion (BIC), as proposed originally by Cragg and Donald (1997) and Ahn et al. (2013). This is formalized right below.

Proposition 1. Let $p = 1, \dots, |\mathcal{B}(L_{\max})|$ and $Q_{N,p}(\widehat{\boldsymbol{\theta}}(\boldsymbol{\Omega}_N) | \widehat{\mathbf{F}}^{(p)})$ be the value of the objective function evaluated at $\widehat{\boldsymbol{\theta}}$ given $\boldsymbol{\Omega}_N$ and some $\widehat{\mathbf{F}}^{(p)} \in \mathcal{B}(L_{\max})$:

$$Q_{N,p}(\widehat{\boldsymbol{\theta}}(\boldsymbol{\Omega}_N) | \widehat{\mathbf{F}}^{(p)}) = \bar{\boldsymbol{\mu}}_N(\widehat{\boldsymbol{\theta}})' \boldsymbol{\Omega}_N \bar{\boldsymbol{\mu}}_N(\widehat{\boldsymbol{\theta}}).$$

Consider the following BIC:

$$S_{N,p} = N \times Q_{N,p}(\widehat{\boldsymbol{\theta}}(\boldsymbol{\Omega}_N) | \widehat{\mathbf{F}}^{(p)}) - \ln(N) \times h(p), \quad (34)$$

where $h(p) = \rho \times \left(\zeta - \dim \left(\widehat{\boldsymbol{\theta}}|_p(P) \right) \right) = \mathcal{O}(1)$, a strictly decreasing function of P with $0 < \rho < \infty$. Under the set of our assumptions, we have

$$\widehat{P} = \arg \min_{\widehat{\boldsymbol{F}}^{(p)} \in \mathcal{B}(L_{\max})} S_{N,p} \xrightarrow{p} L \text{ as } N \rightarrow \infty.$$

Proof. The proof follows directly from Cragg and Donald (1997). The only difference is that for every value of L several weights are potentially available, i.e. one needs to consider minimum BIC for each value of $P = 1, \dots, L_{\max}$ first. \square

Proposition 1 implies that practitioners can estimate models for all choices of $\mathcal{B}(L_{\max})$ and pick \widehat{P} (together with the associated combination of factor proxies) as the value of L that corresponds to the smallest BIC value. The above result holds as long as for $P = L$ there exists an element in $\mathcal{B}(L_{\max})$ that ensures that the GMM estimator is consistent. For instance, in the context of Example 3 where $L = 1$, one such element in $\mathcal{B}(L_{\max})$ is given by $N^{-1} \sum_i^N \boldsymbol{x}_i^{(2)} = \widehat{\boldsymbol{f}}_R^{(2)}$. Notice that in this case, the GMM estimator based on best-subset selection is more efficient than the GMM estimator based on regularization because asymptotically the former employs 1 factor proxy and the latter employs 2 factor proxies. On the other hand, if no such element exists (as e.g. in Example 2), but there exists at least one $\widehat{\boldsymbol{F}}^{(p)}$ that satisfies Assumption 3.1 with $L_e > L$, then BIC will consistently estimate L_e instead. Finally, if no such $\widehat{\boldsymbol{F}}^{(p)}$ exists, then BIC will not consistently estimate L or L_e due to lack of identification, as discussed in Section 3.4. Note that the selected model will be rejected with high probability by the J-statistic in this case.

4. Finite Sample Evidence

4.1. Setup

We consider the following dynamic model with one or two factors:

$$y_{i,t} = \alpha y_{i,t-1} + \beta x_{i,t} + \sum_{r=1}^2 \lambda_{r,i}^y f_{r,t} + \varepsilon_{i,t}^y; \quad x_{i,t} = \delta y_{i,t-1} + \alpha_x x_{i,t-1} + \lambda_{1,i}^x f_{1,t} + \varepsilon_{i,t}^x,$$

for $t > 0$, while for $t = 0$ we set:

$$y_{i,0} = \sum_{r=1}^2 \lambda_{r,i} f_{r,0} + \varepsilon_{i,0}^y; \quad x_{i,0} = \lambda_{1,i}^x f_{1,0} + \varepsilon_{i,0}^x.$$

Additional covariates to be used in the construction of factor proxies are generated as

$$\boldsymbol{v}_{i,t}^{(1)} = \lambda_{1,i}^{v1} f_{1,t} + \varepsilon_{i,t}^{v1}; \quad \boldsymbol{v}_{i,t}^{(2)} = \sum_{r=1}^2 \lambda_{r,i}^{v2} f_{r,t} + \varepsilon_{i,t}^{v2}.$$

The factor loadings for the first factor are normally distributed with mean equal to μ_λ and unit variance, such that

$$\lambda_{1,i}^\psi = \mu_\lambda + \rho(\lambda_{1,i}^y - \mu_\lambda) + \sqrt{1 - \rho^2}v_{1,i}^\psi; \quad v_{1,i}^\psi \sim \mathcal{N}(0, 1),$$

where $\psi \in \{x, v1, v2\}$, and $\lambda_{1,i}^y \sim \mathcal{N}(\mu_\lambda, 1)$. ρ denotes the correlation coefficient between the factor loadings of the $y_{i,t}$ and $x_{i,t}$, and $y_{i,t}$ and $v_{i,t}^{(1)}, v_{i,t}^{(2)}$ processes.

The properties of the factor loadings that correspond to $f_{2,t}$ depend on the setup we consider. In particular, in the case where $L = L_e = 1$ we simply set $\lambda_{2,i}^y = \lambda_{2,i}^{v2} = 0$ for all i . In the case where $L = L_e = 2$ the corresponding factor loadings are drawn independently of other factors loadings, i.e.

$$\lambda_{2,i}^y \sim \mathcal{N}(\mu_\lambda, 1); \quad \lambda_{2,i}^{v2} \sim \mathcal{N}(\mu_\lambda^{v2}, 1). \quad (35)$$

Such a setup facilitates the interpretability of the simulation results, without over-parameterising what is already a large set of nuisance parameters. In what follows, we fix $\mu_\lambda^{v2} = 1$. Finally, all factors are drawn as $f_{r,t} \sim \mathcal{N}(0, 1)$.

The idiosyncratic errors are generated as

$$\varepsilon_{i,t}^y \sim \mathcal{N}(0, 1); \quad \varepsilon_{i,t}^x \sim \mathcal{N}(0, \sigma_x^2); \quad \varepsilon_{i,t}^{v1} \sim \mathcal{N}(0, 1); \quad \varepsilon_{i,t}^{v2} \sim \mathcal{N}(0, 1); \quad t \geq 0.$$

In all designs the value of σ_x^2 is fixed to ensure that the signal-to-noise ratio of the model

$$\text{SNR} \equiv \frac{1}{T} \sum_{t=1}^T \frac{\text{var}(y_{i,t} | \lambda_{1,i}^y, \lambda_{2,i}^y, \lambda_{1,i}^x, \lambda_{2,i}^x, \{f_{1,s}\}_{s=0}^t, \{f_{2,s}\}_{s=0}^t)}{\text{var}(\varepsilon_{i,t}^y)} - 1,$$

equals 5. The chosen SNR value lies within the range of values considered in the literature, e.g. Bun and Kiviet (2006) specifies $\text{SNR} \in \{3; 9\}$.

We consider $N \in \{200; 800\}$, $T \in \{4; 8\}$, $\alpha \in \{0.4; 0.8\}$, and we set $\beta = 1 - \alpha$, such that the ‘‘long-run’’ parameter equals 1. The values of the remaining parameters are as follows: $\delta \in \{0; 0.3\}$, $\mu_\lambda = 1$, $\rho = 0.6$, and $\alpha_x = 0.6$. The number of replications performed equals 2,000 for each design and the factors are drawn in each replication.

Remark 8. The Supplementary Appendix to the present paper provides further finite sample results that correspond to alternative designs. In particular, among other setups, we examine the case where \mathbf{v}_i contains more factors than those that already drive y (i.e. $L_e > L$); moreover, we examine the effect of lack of identification by setting $\mu_\lambda = 0$, and we also specify a model with an additional covariate. In addition, we report results with respect to the nonlinear GMM estimators of Ahn et al. (2013) and Robertson and Sarafidis (2015) for the one-factor design.

4.2. Results

We investigate the finite sample properties of the following four estimators: “F1” denotes the GMM estimator that uses $\mathbf{v}_i^{(1)}$ and $w_i = 1$; “F2” is the GMM estimator that uses $(\mathbf{v}_i^{(1)}, \mathbf{v}_i^{(2)})$ and $w_i = 1$; “Fr” denotes the regularized GMM estimator that uses $(\mathbf{v}_i^{(1)}, \mathbf{v}_i^{(2)})$ and $\mathbf{w}_i = (1, y_{i,0})'$, as described in the regularization approach outlined in Section 3.5. In order to isolate the effect of regularization on the construction of factor proxies, we take L_e as given. Results corresponding to Fr based on an estimate of L_e are reported in the Supplementary Appendix of the paper; finally, “Fbic” denotes the estimator that employs different combinations of factor proxies based on the choices above, and picks up the proxies corresponding to the minimum BIC value. Thus, Fbic represents the best-subset selection approach outlined in Section 3.5. The value of ρ in the BIC function is set equal to 0.75, following BIC1 in Ahn et al. (2013), and Zhu et al. (2019).

All estimators use the maximum number of moment conditions with respect to lagged values of $y_{i,t}$ and $x_{i,t}$, although for comparison purposes $x_{i,t}$ is always treated as weakly exogenous regardless of the value of δ . Notice that, for reasons explained in the Supplementary Appendix, setting $w_i = 1$ leads to slightly better results compared to $w_i = y_{i,0}$ in the present design. Therefore for $L = L_e = 1$, F1 can be viewed as an “oracle” estimator because it makes use of the true number of factors, as well as the “right” choice of \mathbf{v}_i and w_i . Likewise, F2 can be viewed as an oracle estimator for $L = L_e = 2$. Therefore, the performances of F1 and F2 may serve as good benchmarks for the performance of Fr and Fbic.

Tables A1 and A2 report results in terms of bias, RMSE, standard deviation and empirical size of the t -test for $L = L_e = 1$ and $L = L_e = 2$, respectively. Nominal size is set equal to 5%. We make use of corrected standard errors, which are computed based on Windmeijer (2005). This is important because, as is well known in the dynamic panel data literature, the two-step GMM estimator may exhibit substantial size distortions, especially when the number of moment conditions is large relative to N .

As can be seen from Table A1, all four estimators show negligible bias under all combinations of N , T , α and δ . In addition, the RMSE values of all estimators are small and fall roughly at the rate of \sqrt{N} , as predicted by our theory. Fr performs very similar to F1 in terms of RMSE, which indicates that the proposed regularization approach works very well. On the other hand, RMSE is slightly higher for Fbic (especially when $N = 200$), which mostly reflects the fact that L is treated as unknown. Finally, the

empirical size of all estimators is close to the nominal value for both α and β in most cases, unless N is relatively small and T is relatively large, in which case there appear to be upward size distortions. This outcome implies that using too many moment conditions when the cross-sectional dimension is small can result in size distortions, despite the standard error correction. However, it is worth emphasizing that in practice this problem can be mitigated substantially by using only a subset of the moment conditions available. Since this has already been demonstrated by Juodis and Sarafidis (2018) using simulated data for the factor-augmented model, we do not explore this possibility here.

Similar conclusions can be drawn from Table A2 where $L = L_e = 2$, focusing upon the performance of Fr and Fbic vis-a-vis F2. It is worth noting that F1 is not consistent in this case because it estimates one factor only. Thus, it is not surprising that in most cases F1 exhibits large RMSE and substantial size distortions.

Table A3 reports results on the overidentifying restrictions (J) test statistic (nominal size equals 5%) for F1, F2 and Fr. In addition, this table reports selection frequencies of \hat{L} based on BIC, and \hat{L}_e using the ER-statistic, as outlined in Section 3.5. Notice that since $L = L_e$, $P = L$ and therefore we do not distinguish between \hat{P} and \hat{L} . As we can see, for $L = L_e = 1$, the size of the J-statistic corresponding to both F1 and Fr is close to the 5% level in most circumstances, unless N is relatively small and T is relatively large, in which case there is a small downward size distortion, especially for F2. As has been mentioned, such small-sample distortions can be mitigated in practice by using a smaller number of moment conditions. Similar conclusions hold in regards to the performance of the J-statistic corresponding to F2 and Fr for $L = L_e = 2$. Note that since F1 is not consistent in the two-factor case, the results of the J-statistic corresponding to F1 reflect power, which appears to be high under all circumstances.

In regards to model selection, for $L = L_e = 1$ all methods appear to perform well and select the true number of factors with high frequency. As expected, model selection becomes less straightforward for $L = L_e = 2$, especially when both N and T are small. Partially, this may be attributed to the fact that for small values of T the two factors can be highly collinear in some simulated samples, even if they are drawn independently. However, it is clear that the frequency of selecting $\hat{L} = 2$ and $\hat{L}_e = 2$ rises substantially with larger values of either N or T . Finally, the corresponding rows of ER-statistic in Table A3 need not sum to 1, due to rounding.

5. Application: Estimation of the Price Elasticity of Urban Water Demand

A large number of studies have focused on the estimation of the price elasticity of water usage demand, see e.g. House-Peters and Chang (2011) and Araral and Wang (2013) for excellent surveys. The vast majority of the literature assumes that the effect of weather is linear. However, as Maidment and Miaou (1986) and Gato et al. (2007) pointed out, water usage is most likely to respond to changes in weather conditions in a nonlinear fashion. We address this concern in what follows by allowing for nonlinear effects of weather conditions, depending on household/property-specific unobserved characteristics, such as household size, garden size, and others. This setup is represented by a common factor structure.

5.1. Data and methodology

We make use of publicly-available multi-household level data from New South Wales, Australia, provided by the Sydney Water Corporation (SWC), see also Abrams et al. (2012). SWC is the largest water utility in Australia, serving more than 4 million people, while its area of operations covers around 12,700 km^2 .

Our sample contains 4,500 multi-household units, each one being observed over a period of 5 years, 2004-2008. Each unit represents an average of four to six households, in order to preserve privacy. Additional descriptive statistics of the data are reported in the Supplementary Appendix to this paper. The model that we consider for studying the price elasticity of water demand is as follows:

$$y_{i,t} = \alpha y_{i,t-1} + \beta_1 price_{i,t} + \beta_2 rain_{i,t} + \beta_3 temp_{i,t} + \epsilon_{i,t}; \quad \epsilon_{i,t} = \boldsymbol{\lambda}'_i \mathbf{f}_t + \varepsilon_{i,t}, \quad (36)$$

for $i = 1, \dots, N (= 4,500)$ and $t = 1, \dots, T (= 4)$, where $y_{i,t}$ denotes the natural logarithm of the average daily water consumption for household i at year t , expressed in thousands of litres of water (kL); $price_{i,t}$ is the average real price paid per kilolitre of water used by household i at time t ; and $rain_{i,t}$ and $temp_{i,t}$ denote the average amounts of daily rainfall (mm) and temperature (degrees Celsius) during year t . The dynamic specification accommodates partial adjustment mechanisms in water demand. This is due to both habit formation in water consumption and energy efficiency constraints associated with the existing stock of durable goods within households. The log-linear functional form implies that price elasticity depends on the level of price itself. That is, the higher the level of the price, the more sensitive consumers become to changes in price. This is consistent with utility theory, see e.g. Al-Qunaibet and Johnston (1985).

The weather variables are unit-specific because they are determined by the physical proximity of each property to a total of thirteen weather stations that exist across Sydney, operated by the Australian Bureau of Meteorology. This reflects the fact that weather patterns can vary substantially across NSW and, more specifically, the coast generally has more rainfall and cooler conditions than many areas located inland.

The common factor structure allows for nonlinear unobserved heterogeneity across individual units, due to differences in the number of people living in a household, pool ownership, garden size and structure, and so on. As an example, consider two properties: one with two household members, no swimming pool and a small garden, and another one with five members, a large garden and a pool. The difference in average yearly consumption between the two properties is expected to be proportionately larger under extreme weather conditions than under normal conditions. To put this differently, the change in water consumption following an extreme weather event is likely to be smaller for the former property than the latter.

The formulation above requires that household size (and other property-specific characteristics) remains constant over the 5-year period of our analysis. While it is unreasonable to expect such condition to be fulfilled for all households in the sample, recall that each individual unit $i = 1, \dots, N$ represents an average of four to six households, which is due to the data aggregation implemented by SWC in the original dataset. As a result, changes in household size over time are likely to be smoothed out to a large extent. If this is not true, our estimator will not be consistent. This implies that violation of such restriction can indeed be detected in practice using the usual J-statistic.

We implement our methodology by assuming that the unobserved factor component is approximated by an external variable, the average daily soil moisture index (smi) observed for unit i at period t . smi is computed based on a combination of precipitation, temperature and soil moisture, and is used widely as an index accounting for extreme weather and soil drought intensity, see e.g. Hunt et al. (2009). Hence, in terms of the notation used in Section 3.1, we set $\mathbf{v}_{i,t} = smi_{i,t}$.

We estimate the model by fitting $L = \{0, 1, 2\}$ factors. We employ three weights, namely $\mathbf{w}_i = (1, y_{i,0}, y_{i,0}^2)'$. In terms of the notation introduced in Section 3, we have $L_{\max} = 2$, $R = 3$ and the set $\mathcal{B}(L_{\max})$ contains all 6 possible permutations of $\widehat{\mathbf{F}}_R = \frac{1}{N} \sum_{i=1}^N \mathbf{v}_i \mathbf{w}_i'$, where \mathbf{v}_i is constructed by stacking $smi_{i,t}$ for $t = 1, \dots, T$ in a $[T \times 1]$ vector. The implementation of our approach based on regularization follows closely

the steps described in Section 3.5, with L_e estimated from $T^{-1}\widehat{\mathbf{F}}_R\widehat{\mathbf{F}}_R'$ using the ER-statistic outlined therein. The implementation based on best-subset selection estimates the model using six different combinations of weights, i.e. those in $\mathcal{B}(L_{\max})$. The estimated number of factors, \widehat{P} , as well as the associated selected weight combination, are determined using the BIC criterion documented in Section 3.5. The value of ρ in the BIC function is set equal to 0.75, as in the Monte Carlo setup.

5.2. Results

Table 1 reports results for ten different specifications. In specific, $M0$ denotes the estimator that imposes $\epsilon_{i,t} = \varepsilon_{i,t}$, i.e. it assumes unobserved heterogeneity away. M_{DIF} and M_{SYS} denote the first-differenced and system GMM estimators proposed by Arellano and Bond (1991) and Arellano and Bover (1995), respectively. Both estimators allow for an additive two-way error components structure and are obtained using the `xtabond2` command in Stata 15, see Roodman (2009). $M_{\widetilde{F}}$ denotes the GMM estimator that uses regularized factor proxies, $\widetilde{\mathbf{F}}$, based on $\widehat{L}_e = 1$ (the value of which is obtained from the ER-statistic). In addition, $M1_c$, $M1_{y_0}$ and $M1_{y_0^2}$ denote the GMM estimators that impose $L_e = 1$ and use $\mathbf{w}_i = 1$, $\mathbf{w}_i = y_{i,0}$ and $\mathbf{w}_i = y_{i,0}^2$, respectively. Finally, $M2_{c,y_0}$ imposes $L_e = 2$ with weights given by $\mathbf{w}_i = (1, y_{i,0})'$. The same holds for $M2_{c,y_0^2}$ and $M2_{y_0,y_0^2}$, except that $\mathbf{w}_i = (1, y_{i,0}^2)'$ and $\mathbf{w}_i = (y_{i0}, y_{i0}^2)'$, respectively.

In all models the price variable is treated as endogenous and is instrumented by appropriate lagged values of the same variable. This is because during the period of the analysis a two-tier pricing scheme was in place in NSW, such that consumers paid a higher price when their consumption levels exceeded a certain threshold. On the other hand, all weather variables are treated as exogenous with respect to $\varepsilon_{i,t}$.

All estimators that allow for a common factor structure, make use of $\zeta = 40$ moment conditions, whereas M_{DIF} and M_{SYS} make use of 17 and 21 moment conditions, respectively. The difference in the number of moment conditions between the GMM estimators that impose an additive structure and those that allow for a genuine factor model is mainly due to their treatment of the exogenous weather variables. In the former case, standard practice involves taking first differences to remove unobserved (linear) heterogeneity and then using the exogenous weather variables as Anderson-Hsiao type instruments. In the latter case, present and lagged values of the weather variables are included as instruments in each time period, in order to allow for possible arbitrary correlations between nonlinear heterogeneity and weather conditions.

Results are reported in terms of the estimated coefficients and standard errors (in

parentheses), where $\widehat{\beta}_1/(1 - \widehat{\alpha})$ corresponds to the long-run price coefficient, the standard error of which has been obtained using the Delta method. Furthermore, Table 1 reports the J-statistic and its p-value (in square brackets), the number of moment conditions and of parameters estimated for each model, and finally the value of the *BIC*.

Table 1: Results

	$M0$	M_{DIF}	M_{SYS}	$M_{\widehat{F}}$	$M1_c$	$M1_{y_0}$	$M1_{y_0^2}$	$M2_{c,y_0}$	$M2_{c,y_0^2}$	$M2_{y_0,y_0^2}$
$\widehat{\alpha}$.942 (.004)	.504 (.045)	.771 (.018)	.405 (.047)	.414 (.048)	.405 (.048)	.393 (.047)	.354 (.082)	.382 (.082)	.395 (.087)
$\widehat{\beta}_1$.170 (.015)	-.032 (.361)	.261 (.419)	-.185 (.034)	-.178 (.035)	-.187 (.035)	-.192 (.034)	-.247 (.061)	-.225 (.054)	-.190 (.051)
$\widehat{\beta}_2$	-.051 (.003)	-.001 (.008)	-.003 (.009)	-.013 (.006)	-.013 (.006)	-.021 (.010)	-.013 (.006)	-.030 (.009)	-.011 (.010)	-.018 (.011)
$\widehat{\beta}_3$	-.008 (.001)	.034 (.012)	.029 (.013)	.050 (.010)	.050 (.010)	.084 (.019)	.050 (.010)	.032 (.013)	.060 (.015)	.051 (.016)
$\frac{\widehat{\beta}_1}{1-\widehat{\alpha}}$	2.93 (.075)	-.065 (.726)	1.14 (1.86)	-.310 (.078)	-.303 (.079)	-.314 (.079)	-.316 (.077)	-.383 (.137)	-.364 (.128)	-.314 (.113)
J test	156.3	27.6	49.2	28.8	28.7	28.9	28.9	13.6	16.1	15.6
p-val.	[.000]	[.002]	[.000]	[.092]	[.094]	[.092]	[.091]	[.096]	[.042]	[.050]
ζ	40	17	21	40	40	40	40	40	40	40
$dim(\boldsymbol{\theta})$	5	9	10	20	20	20	20	32	32	32
BIC	10.6	-17.8	-4.91	-54.48	-54.54	-54.46	-54.39	-19.69	-17.24	-17.74

It is clear that the p-value of the J-statistic is close to zero when we fit either zero factors or the two-way (additive) error components structure, which implies that the model is mis-specified. This is also reflected in the estimated price coefficient, which is largely statistically insignificant for both M_{DIF} and M_{SYS} , and for the latter it even has the wrong sign. We note that these results are not sensitive to the number of instruments employed. For example, using only the two most recent lags of the dependent variable as instruments, i.e. $y_{i,t-2}$ and $y_{i,t-3}$, along with the *collapse* option in Stata, the J-statistic for M_{SYS} roughly equals 27.6, and so the p-value remains close to zero.

On the other hand, fitting one or two genuine factors fails to reject the specified model at the 1% level of significance. This finding provides evidence that the factor structure is supported by the data compared to the additive error components model, and demonstrates the importance of controlling for nonlinear heterogeneity.

More specifically, recall that the ER-statistic indicates that $\widehat{L}_e = 1$. In addition, according to the BIC criterion, $\widehat{P} = 1$. Furthermore, $M1_c$, the estimator that uses a constant weight, $w_i = 1$, minimizes the overall BIC criterion. The results from adopting the regularization approach ($M_{\widehat{F}}$) and the best-subset selection approach ($M1_c$) are similar. While this may be partially attributed to $\widehat{L}_e = \widehat{P} = 1$, it is evident that the estimated price elasticity of demand appears to be robust across different factor proxies and different values of L . This is a desirable outcome.

In what follows we discuss further findings based on $M_{\widehat{F}}$. A unit (dollar) increase in the price of water is estimated to reduce water consumption by approximately 18.5% and 31.0% in the short- and long-run, respectively. Similarly, a unit increase in rain (temperature) is expected to reduce (increase) water consumption by approximately 1.13% (5.00%) in the short-run.

The price elasticity of demand is computed by multiplying the relevant price coefficients with a range of values for price. Table 2 presents elasticity estimates at four different values of price; namely, mean and median price, as well as the 10th and 90th percentiles. For example, at the median price of \$1.37 per kL, a 1% increase in the price of water lowers demand by about .25% in the short-run and .43% in the long-run.

Table 2: Point-wise predicted elasticities for M_{F_R} .

	10th perc.	mean	median	90th perc.
price	1.17	1.35	1.37	1.56
SR elasticity	-.216	-.249	-.253	-.288
LR elasticity	-.362	-.419	-.426	-.483

As expected, urban water demand appears to be much more elastic in the long-run than in the short-run. This may be attributed to habit formation and technological constraints of water appliance efficiency. Moreover, the value of $\widehat{\alpha} = .405$ implies that it takes about 2.5 years for 90% of the total price effect to be realized, all other things being constant. This outcome casts doubt on the potential effectiveness of scarcity pricing to balance demand and supply of water in periods of transitory droughts.

In comparison to other studies in the literature, the estimated price elasticity of de-

mand is in the low-to-middle range of results. For instance, the long-run price coefficient is not statistically different to the value obtained by Nauges and Thomas (2003) (see Table III in their paper) although theirs is derived from the constant-elasticity model using municipal-level data and includes average income but not weather conditions.

6. Concluding Remarks

This paper puts forward a novel method-of-moments approach for estimation of factor-augmented panel data models with endogenous regressors and T fixed. The underlying idea is to proxy the factors by observed variables, so that the resulting moment conditions are linear in the parameters. The proposed methodology addresses several issues that arise with existing nonlinear GMM estimators, such as local minima-related problems and a potential lack of global identification. At the same time, the proposed methodology retains the appealing features of the method of moments in that it accommodates weakly exogenous and endogenous regressors without the need for bias correction.

We note that although this paper has explicitly assumed that \mathbf{v}_i (or \mathbf{V}_i) has an additive factor structure with corresponding factor loadings $\boldsymbol{\gamma}_i(\boldsymbol{\mathcal{X}}_i)$, in practice, as it is clear from Eq. (21), the additive factor specification is sufficient but not necessary. In particular, since the proposed estimation method uses information from $E_{\mathcal{F}}[\widehat{\mathbf{F}}_e]$ only, it is sufficient that this expected value is of reduced rank structure, i.e. $E_{\mathcal{F}}[\widehat{\mathbf{F}}_e] = \mathbf{F}_e \mathbf{A}$. The deviations from the mean, i.e. $\overline{\boldsymbol{\Psi}} = N^{-1} \sum_{i=1}^N \boldsymbol{\Psi}_i$, can still contain unobserved common shocks so long as they satisfy the conditional independence in Assumption 2.1. Finally, \mathbf{F}_e and \mathbf{A} should be simply regarded as correspondingly the left and right singular vectors (up to a scaling) of $E_{\mathcal{F}}[\widehat{\mathbf{F}}_e]$. Thus, the $\mathbf{F}_e \mathbf{A}$ decomposition can be assumed completely without loss of generality, provided that L_e is defined accordingly.

We hope that the proposed methodology will enhance the application of multi-factor error structures in panels involving micro level data, and encourage empirical researchers to implement such approaches in practice. Furthermore, since the resulting method of moments estimator has a close form solution, our approach can be extended straightforwardly to several different set ups, such as multi-dimensional panels, spatial panels, pseudo panels, and threshold models, to mention a few. We leave these avenues for future research.

Acknowledgments

We would like to thank the editor, Jianqing Fan, and three anonymous referees for their constructive comments and suggestions, which have helped to improve our paper substantially. We are also grateful to Pavel Čížek, Mervyn Silvapulle, Rutger Teulings and Frank Windmeijer for helpful comments at an earlier stage of this project. Financial support from the NWO VENI grant number 451 – 17 – 002 is gratefully acknowledged by the first author. The second author acknowledges financial support from the Australian Research Council, under research grant number DP-170103135.

References

- ABRAMS, B., S. KUMARADEVAN, V. SARAFIDIS, AND F. SPANINKS (2012): “An Econometric Assessment of Pricing Sydneys Residential Water Use,” *Economic Record*, 88, 89–105.
- AHN, S. C. AND A. R. HORENSTEIN (2013): “Eigenvalue Ratio Test for the Number of Factors,” *Econometrica*, 81, 1203–1227.
- AHN, S. C., Y. H. LEE, AND P. SCHMIDT (2013): “Panel Data Models with Multiple Time-varying Individual Effects,” *Journal of Econometrics*, 174, 1–14.
- AL-QUNAIBET, M. H. AND R. S. JOHNSTON (1985): “Municipal Demand for Water in Kuwait: Methodological Issues and Empirical Results,” *Water Resources Research*, 21, 433–438.
- ARARAL, E. AND Y. WANG (2013): “Water Demand Management: Review of Literature and Comparison in South-East Asia,” *International Journal of Water Resources Development*, 29, 434–450.
- ARELLANO, M. AND S. BOND (1991): “Some Tests of Specification for Panel Data: Monte Carlo Evidence and an Application to Employment Equations,” *Review of Economic Studies*, 58, 277–297.
- ARELLANO, M. AND O. BOVER (1995): “Another Look at the Instrumental Variable Estimation of Error-components Models,” *Journal of Econometrics*, 68, 29–51.
- BAI, J. (2003): “Inferential Theory for Factor Models of Large Dimensions,” *Econometrica*, 71, 135–171.
- (2009): “Panel Data Models With Interactive Fixed Effects,” *Econometrica*, 77, 1229–1279.
- BUN, M. J. G. AND J. F. KIVIET (2006): “The Effects of Dynamic Feedbacks on LS

- and MM Estimator Accuracy in Panel Data Models,” *Journal of Econometrics*, 132, 409–444.
- CONNOR, G. AND R. A. KORAJCZYK (1986): “Performance Measurement with the Arbitrage Pricing Theory: A new Framework for Analysis,” *Journal of Financial Economics*, 15, 373 – 394.
- CRAGG, J. C. AND S. G. DONALD (1997): “Inferring the Rank of a Matrix,” *Journal of Econometrics*, 76, 223–250.
- DE VOS, I. AND G. EVERAERT (2019): “Bias-corrected Common Correlated Effects Pooled Estimation in Homogeneous Dynamic Panels,” *Journal of Business and Economic Statistics*, (forthcoming).
- EVERAERT, G. AND T. DE GROOTE (2016): “Common Correlated Effects Estimation of Dynamic Panels with Cross-Sectional Dependence,” *Econometric Reviews*, 35, 428–463.
- FAN, J. AND Y. LIAO (2019): “Learning Latent Factors from Diversified Projections and its Applications to Over-Estimated and Weak Factors,” *arXiv e-prints*, arXiv:1908.01252.
- GAGLIARDINI, P. AND C. GOURIÉROUX (2017): “Double Instrumental Variable Estimation of Interaction Models with Big Data,” *Journal of Econometrics*, 201, 176–197.
- GAGLIARDINI, P., E. OSSOLA, AND O. SCAILLET (2016): “Time-Varying Risk Premium in Large Cross-Sectional Equity Data Sets,” *Econometrica*, 84, 985–1046.
- GATO, S., N. JAYASURIYA, AND P. ROBERTS (2007): “Temperature and Rainfall Thresholds for Base Use Urban Water Demand Modelling,” *Journal of Hydrology*, 337, 364–376.
- HANSEN, C. AND Y. LIAO (2018): “The Factor-Lasso and K-Step Bootstrap Approach for Inference in High-Dimensional Economic Applications,” *Econometric Theory*, 1–45.
- HASTIE, T., R. TIBSHIRANI, AND J. FRIEDMAN (2017): *The Elements of Statistical Learning*, Springer, 2 ed.
- HAYAKAWA, K. (2016): “Identification Problem of GMM Estimators for Short Panel Data Models with Interactive Fixed Effects,” *Economics Letters*, 139, 22–26.
- HOLTZ-EAKIN, D., W. K. NEWEY, AND H. S. ROSEN (1988): “Estimating Vector Autoregressions with Panel Data,” *Econometrica*, 56, 1371–1395.
- HOUSE-PETERS, L. A. AND H. CHANG (2011): “Urban Water Demand Modeling: Review of Concepts, Methods, and Organizing Principles,” *Water Resources Research*,

- 47, W05401.
- HSIAO, C., M. H. PESARAN, AND A. K. TAHMISCIOGLU (2002): “Maximum Likelihood Estimation of Fixed Effects Dynamic Panel Data Models Covering Short Time Periods,” *Journal of Econometrics*, 109, 107–150.
- HUNT, E., K. HUBBARD, D. WILHITE, T. ARKEBAUER, AND A. DUTCHER (2009): “The Development and Evaluation of a Soil Moisture Index,” *International Journal of Climatology*, 29, 747–759.
- JIANG, B., Y. YANG, J. GAO, AND C. HSIAO (2017): “Recursive Estimation in Large Panel Data Models: Theory and Practice,” Mimeo.
- JUODIS, A. (2018): “Pseudo Panel Data Models with Cohort Interactive Effects,” *Journal of Business & Economic Statistics*, 36, 47–61.
- JUODIS, A., H. KARABIYIK, AND J. WESTERLUND (2020): “On the Robustness of the Pooled CCE Estimator,” *Journal of Econometrics*, (forthcoming).
- JUODIS, A. AND V. SARAFIDIS (2018): “Fixed T Dynamic Panel Data Estimators with Multifactor Errors,” *Econometric Reviews*, 37, 893–929.
- KARABIYIK, H., J.-P. URBAIN, AND J. WESTERLUND (2019): “CCE estimation of factor-augmented regression models with more factors than observables,” *Journal of Applied Econometrics*, 34, 268–284.
- KLEIBERGEN, F. R. AND R. PAAP (2006): “Generalized Reduced Rank Tests Using the Singular Value Decomposition,” *Journal of Econometrics*, 133, 97–126.
- KRUINIGER, H. (2008): “Not So Fixed Effects: Correlated Structural Breaks in Panel Data,” Mimeo.
- KUERSTEINER, G. AND I. R. PRUCHA (2013): “Limit Theory for Panel Data Models with Cross Sectional Dependence and Sequential Exogeneity,” *Journal of Econometrics*, 174, 107–126.
- MAIDMENT, D. R. AND S. P. MIAOU (1986): “Daily Water Use in Nine Cities,” *Water Resources Research*, 22, 845 – 851.
- MOON, H. R. AND M. WEIDNER (2015): “Linear Regression for Panel With Unknown Number of Factors as Interactive Fixed Effects,” *Econometrica*, 83, 1543–1579.
- NAUGES, C. AND A. THOMAS (2003): “Long-run Study of Residential Water Consumption,” *Environmental and Resource Economics*, 26, 25–43.
- PESARAN, M. H. (2006): “Estimation and Inference in Large Heterogeneous Panels with a Multifactor Error Structure,” *Econometrica*, 74, 967–1012.
- PESARAN, M. H., L. V. SMITH, AND T. YAMAGATA (2013): “Panel Unit Root Tests

- in the Presence of a Multifactor Error Structure,” *Journal of Econometrics*, 175, 94 – 115.
- ROBERTSON, D. AND V. SARAFIDIS (2015): “IV Estimation of Panels with Factor Residuals,” *Journal of Econometrics*, 185, 526–541.
- ROBERTSON, D., V. SARAFIDIS, AND J. WESTERLUND (2018): “Unit Root Inference in Generally Trending and Cross-Correlated Dynamic Panels,” *Journal of Business & Economic Statistics*, 36, 493–504.
- ROODMAN, D. (2009): “How To Do xtabond2: An Introduction to Difference and System GMM in Stata,” *The Stata Journal*, 9, 86 – 136.
- SARAFIDIS, V. AND T. J. WANSBEEK (2012): “Cross-sectional Dependence in Panel Data Analysis,” *Econometric Reviews*, 31, 483–531.
- STAIGER, D. AND J. H. STOCK (1997): “Instrumental Variables Regression with Weak Instruments,” *Econometrica*, 65, 557–586.
- WESTERLUND, J., Y. PETROVA, AND M. NORKUTE (2019): “CCE in Fixed-T Panels,” *Journal of Applied Econometrics*, 34, 746–761.
- WINDMEIJER, F. (2005): “A Finite Sample Correction for the Variance of Linear Efficient Two-Step GMM Estimators,” *Journal of Econometrics*, 126, 25–51.
- ZHU, H., V. SARAFIDIS, AND M. SILVAPULLE (2019): “A New Structural Break Test for Panels with Common Factors,” *Econometrics Journal*, forthcoming.

Appendix A. Monte Carlo Results

Table A1: Estimation results for $L = 1$ and $L_e = 1$ setup.

Designs				F1				F2				Fr				Fbic				
N	T	α	δ	Bias	RMSE	Std	Size													
α																				
200	4	4.0	.00	.02	.02	.06	.00	.05	.05	.02	.00	.02	.02	.06	.00	.04	.04	.07		
200	4	4.3	.00	.03	.03	.05	-.01	.07	.07	.02	.00	.03	.03	.07	.00	.04	.04	.08		
200	4	8.0	.00	.03	.03	.06	.00	.05	.05	.02	.00	.02	.02	.07	.00	.03	.03	.07		
200	4	8.3	.00	.03	.03	.07	.00	.06	.06	.03	.00	.03	.03	.07	.00	.03	.03	.08		
200	8	4.0	.00	.02	.02	.11	.00	.02	.02	.06	.00	.01	.01	.09	.00	.02	.02	.11		
200	8	4.3	-.01	.03	.03	.18	-.01	.03	.03	.09	-.01	.03	.03	.15	-.01	.04	.03	.18		
200	8	8.0	.00	.01	.01	.10	.00	.02	.02	.07	.00	.01	.01	.12	.00	.02	.02	.13		
200	8	8.3	.00	.02	.02	.13	.00	.02	.02	.07	.00	.02	.02	.13	.00	.02	.02	.13		
800	4	4.0	.00	.01	.01	.06	.00	.03	.03	.02	.00	.01	.01	.04	.00	.01	.01	.07		
800	4	4.3	.00	.02	.02	.05	.00	.04	.04	.02	.00	.02	.02	.06	.00	.02	.02	.06		
800	4	8.0	.00	.01	.01	.06	.00	.03	.03	.02	.00	.01	.01	.06	.00	.01	.01	.07		
800	4	8.3	.00	.01	.01	.05	.00	.04	.04	.02	.00	.01	.01	.05	.00	.02	.02	.06		
800	8	4.0	.00	.01	.01	.06	.00	.01	.01	.05	.00	.01	.01	.06	.00	.01	.01	.06		
800	8	4.3	.00	.01	.01	.06	.00	.01	.01	.05	.00	.01	.01	.06	.00	.01	.01	.06		
800	8	8.0	.00	.01	.01	.07	.00	.01	.01	.06	.00	.01	.01	.07	.00	.01	.01	.07		
800	8	8.3	.00	.01	.01	.06	.00	.01	.01	.05	.00	.01	.01	.07	.00	.01	.01	.07		
β																				
200	4	4.0	.00	.03	.03	.07	.00	.07	.07	.02	.00	.02	.02	.07	.00	.06	.06	.06		
200	4	4.3	.00	.03	.03	.06	.01	.08	.08	.02	.00	.03	.03	.07	.00	.04	.04	.06		
200	4	8.0	.00	.01	.01	.06	.00	.03	.03	.02	.00	.01	.01	.06	.00	.01	.01	.06		
200	4	8.3	.00	.01	.01	.05	.00	.04	.04	.03	.00	.01	.01	.06	.00	.02	.02	.06		
200	8	4.0	.00	.02	.02	.11	.00	.02	.02	.06	.00	.02	.02	.11	.01	.02	.02	.13		
200	8	4.3	.01	.04	.03	.19	.01	.03	.03	.10	.01	.03	.03	.16	.01	.04	.04	.19		
200	8	8.0	.00	.01	.01	.10	.00	.01	.01	.06	.00	.01	.01	.10	.00	.01	.01	.10		
200	8	8.3	.00	.02	.02	.13	.00	.02	.02	.07	.00	.01	.01	.12	.00	.02	.02	.14		
800	4	4.0	.00	.01	.01	.06	.00	.05	.05	.02	.00	.01	.01	.06	.00	.02	.02	.06		
800	4	4.3	.00	.02	.02	.05	.00	.05	.05	.03	.00	.02	.02	.05	.00	.02	.02	.05		
800	4	8.0	.00	.01	.01	.06	.00	.02	.02	.02	.00	.01	.01	.06	.00	.01	.01	.06		
800	4	8.3	.00	.01	.01	.05	.00	.03	.03	.02	.00	.01	.01	.05	.00	.01	.01	.06		
800	8	4.0	.00	.01	.01	.06	.00	.01	.01	.05	.00	.01	.01	.06	.00	.01	.01	.07		
800	8	4.3	.00	.01	.01	.06	.00	.01	.01	.05	.00	.01	.01	.06	.00	.01	.01	.06		
800	8	8.0	.00	.00	.00	.06	.00	.01	.01	.04	.00	.00	.00	.06	.00	.00	.00	.06		
800	8	8.3	.00	.01	.01	.06	.00	.01	.01	.05	.00	.01	.01	.06	.00	.01	.01	.06		

Table A2: Estimation results for $L = 2$ and $L_e = 2$ setup.

Designs				F1				F2				Fr				Fbic			
N	T	α	δ	Bias	RMSE	Std	Size												
α																			
200	4	4.0		-.02	.14	.14	.60	.00	.04	.04	.05	.00	.04	.04	.05	-.01	.09	.09	.09
200	4	4.3		-.05	.26	.26	.66	.00	.06	.06	.04	.00	.06	.06	.05	-.01	.22	.22	.09
200	4	8.0		-.02	.14	.13	.60	.00	.04	.04	.05	.00	.04	.04	.05	-.01	.09	.09	.10
200	4	8.3		-.03	.17	.17	.64	.00	.05	.05	.05	.00	.05	.05	.06	-.01	.15	.15	.11
200	8	4.0		-.03	.11	.11	.76	.00	.02	.02	.08	.00	.02	.02	.08	.00	.02	.02	.09
200	8	4.3		-.10	.33	.31	.83	.00	.03	.03	.11	.00	.03	.03	.11	.00	.04	.04	.11
200	8	8.0		-.02	.09	.09	.73	.00	.01	.01	.09	.00	.01	.01	.09	.00	.02	.02	.09
200	8	8.3		-.08	.20	.18	.79	.00	.02	.02	.10	.00	.02	.02	.10	.00	.02	.02	.11
800	4	4.0		-.02	.14	.14	.79	.00	.02	.02	.05	.00	.02	.02	.04	.00	.06	.06	.07
800	4	4.3		-.05	.25	.24	.80	.00	.03	.03	.06	.00	.03	.03	.05	-.01	.13	.13	.08
800	4	8.0		-.02	.13	.13	.78	.00	.02	.02	.05	.00	.02	.02	.05	.00	.04	.04	.08
800	4	8.3		-.04	.17	.17	.80	.00	.02	.02	.05	.00	.02	.02	.06	-.01	.12	.12	.08
800	8	4.0		-.02	.11	.10	.85	.00	.01	.01	.06	.00	.01	.01	.06	.00	.01	.01	.06
800	8	4.3		-.09	.31	.30	.91	.00	.01	.01	.06	.00	.01	.01	.06	.00	.01	.01	.07
800	8	8.0		-.02	.08	.08	.84	.00	.01	.01	.06	.00	.01	.01	.06	.00	.01	.01	.06
800	8	8.3		-.07	.17	.16	.87	.00	.01	.01	.06	.00	.01	.01	.07	.00	.01	.01	.08
β																			
200	4	4.0		.02	.14	.14	.42	.00	.06	.06	.06	.00	.06	.06	.05	.01	.21	.21	.08
200	4	4.3		.04	.26	.26	.54	.00	.07	.07	.05	.00	.08	.08	.06	.00	.28	.28	.08
200	4	8.0		.01	.05	.05	.32	.00	.03	.03	.05	.00	.03	.03	.06	.00	.07	.07	.06
200	4	8.3		.01	.07	.07	.35	.00	.04	.04	.05	.00	.04	.04	.04	.01	.13	.13	.07
200	8	4.0		.03	.13	.13	.71	.00	.02	.02	.09	.00	.02	.02	.09	.00	.02	.02	.09
200	8	4.3		.11	.39	.37	.81	.00	.03	.03	.11	.00	.03	.03	.11	.00	.04	.04	.12
200	8	8.0		.01	.05	.05	.61	.00	.01	.01	.09	.00	.01	.01	.09	.00	.01	.01	.10
200	8	8.3		.07	.19	.18	.74	.00	.02	.02	.10	.00	.02	.02	.10	.00	.02	.02	.11
800	4	4.0		.02	.14	.14	.64	.00	.03	.03	.05	.00	.03	.03	.05	.00	.18	.18	.07
800	4	4.3		.04	.24	.24	.72	.00	.03	.03	.05	.00	.04	.04	.05	.01	.12	.12	.07
800	4	8.0		.00	.04	.04	.53	.00	.01	.01	.06	.00	.01	.01	.05	.00	.13	.13	.07
800	4	8.3		.01	.08	.08	.56	.00	.02	.02	.05	.00	.02	.02	.06	.00	.12	.12	.07
800	8	4.0		.03	.12	.11	.81	.00	.01	.01	.05	.00	.01	.01	.06	.00	.01	.01	.06
800	8	4.3		.09	.37	.36	.90	.00	.01	.01	.06	.00	.01	.01	.06	.00	.01	.01	.07
800	8	8.0		.01	.04	.04	.74	.00	.01	.01	.05	.00	.01	.01	.06	.00	.01	.01	.06
800	8	8.3		.06	.16	.15	.83	.00	.01	.01	.07	.00	.01	.01	.06	.00	.01	.01	.07

Table A3: Model selection and specification testing results.

		$L = L_e = 1$						$L = L_e = 2$											
Designs		J-Statistic			BIC			ER L_e			J-Statistic			BIC			ER L_e		
N	T α δ	F1	F2	Fr	$\#\hat{L} = 1$	$\#\hat{L} = 2$	$\#\hat{L} = 3$	$\#\hat{L} = 1$	$\#\hat{L} = 2$	$\#\hat{L} = 3$	$\#\hat{L} = 1$	$\#\hat{L} = 2$	$\#\hat{L} = 3$	$\#\hat{L} = 1$	$\#\hat{L} = 2$	$\#\hat{L} = 3$	$\#\hat{L} = 1$	$\#\hat{L} = 2$	$\#\hat{L} = 3$
200	4.4.0.	.03	.01	.05	.98	.02	.02	.98	.00	.02	.16	.84	.16	.84	.16	.84	.16	.76	.08
200	4.4.3.	.03	.01	.04	.98	.02	.02	.97	.00	.02	.16	.84	.16	.84	.17	.84	.17	.76	.07
200	4.8.0.	.05	.02	.05	.98	.02	.02	.98	.00	.02	.20	.80	.20	.80	.17	.76	.17	.76	.08
200	4.8.3.	.04	.02	.04	.98	.02	.03	.97	.00	.03	.20	.80	.20	.80	.15	.77	.15	.77	.08
200	8.4.0.	.03	.01	.03	.96	.04	.00	1	.00	.00	.01	.99	.01	.99	.06	.93	.06	.93	.01
200	8.4.3.	.03	.01	.03	.96	.04	.00	1	.00	.00	.02	.98	.02	.98	.07	.93	.07	.93	.01
200	8.8.0.	.02	.01	.03	.97	.03	.00	1	.00	.00	.01	.99	.01	.99	.07	.92	.07	.92	.01
200	8.8.3.	.03	.01	.03	.97	.03	.00	1	.00	.00	.02	.98	.02	.98	.06	.93	.06	.93	.01
800	4.4.0.	.06	.02	.06	.99	.01	.01	.99	.00	.01	.05	.95	.05	.95	.07	.91	.07	.91	.03
800	4.4.3.	.04	.02	.04	1	.00	.01	.99	.00	.01	.05	.96	.04	.96	.07	.90	.07	.90	.03
800	4.8.0.	.04	.02	.04	1	.00	.01	.99	.00	.01	.06	.94	.06	.94	.08	.89	.08	.89	.03
800	4.8.3.	.05	.02	.06	1	.00	.01	.99	.00	.01	.06	.95	.05	.95	.07	.90	.07	.90	.03
800	8.4.0.	.05	.02	.04	1	.01	.00	1	.00	.00	.00	1	.05	.00	.01	.99	.01	.99	.00
800	8.4.3.	.04	.02	.04	.99	.01	.00	1	.00	.00	.00	1	.05	.00	.01	.99	.01	.99	.00
800	8.8.0.	.05	.02	.05	1	.00	.00	1	.00	.00	.00	1	.05	.00	.01	.99	.01	.99	.00
800	8.8.3.	.05	.02	.05	.99	.01	.00	1	.00	.00	.00	1	.05	.00	.01	.99	.01	.99	.00

A Linear Estimator for Factor-Augmented Fixed-T Panels with Endogenous Regressors

Supplementary Appendix

Artūras Juodis^a, Vasilis Sarafidis^{b,*}

^a*Faculty of Economics and Business, University of Groningen.*

^b*Department of Econometrics and Business Statistics, Monash University.*

Contents

S1 Extensions	2
S1.1 Unbalanced Panels	2
S1.2 Observed Factors	4
S1.3 Consistency Under an Alternative Set of Assumptions	5
S1.4 Relative Efficiency Comparison with the CCE Estimator	6
S2 Descriptive Statistics of Water Usage Data	7
S3 Further Finite Sample Evidence	9
S3.1 The Growth Ratio Statistic	9
S3.2 Residual Based Estimation of L	10
S3.3 $L = 1$ and $L_e = 2$	11
S3.4 Zero Mean Factor Loadings	12
S3.5 “Fr” with Estimated L_e	14
S3.6 Estimation with Additional Regressor	14
S3.7 Nonlinear GMM Estimators	14
S4 Proofs	32
S4.1 Auxiliary Results	32
S4.2 Main Theorems	33
S4.3 Additional Results	39

*Corresponding author. 900 Dandenong Road, Caulfield East, Victoria 3145, Australia. E.mail: vasilis.sarafidis@monash.edu

Contents

Section S1 of the present Supplementary Appendix analyses several extensions of the model analysed in the main paper. In particular, Section S1.1 discusses estimation of unbalanced panels and Section S1.2 accommodates the scenario where some common factors are observed. In addition, Section S1.3 demonstrates consistency of the proposed approach under an alternative set of assumptions, in which the factor loadings are treated as a sequence of constants. Lastly, Section S1.4 implements an analytic relative efficiency comparison between the CCE estimator proposed by Pesaran (2006) and the linear GMM estimator proposed in the present paper. We demonstrate that in a stylised model with a single exogenous regressor, these two estimators have the same asymptotic variance. Section S2 provides descriptive statistics for the data used in the empirical illustration. Section S3 reports additional Monte Carlo results. Finally, Section S4 provides proofs of theoretical results put forward in the main paper.

S1. Extensions

S1.1. Unbalanced Panels

There is currently no literature on estimating unbalanced panels with a factor error structure when T is fixed. In certain cases, extending existing nonlinear GMM estimators to accommodate unbalanced panels requires additional, possibly nontrivial restrictions. For example, the quasi-long-differencing transformation employed by Ahn et al. (2013) requires that for all individuals at least L *common* time observations are available. In a model with weakly exogenous or endogenous regressors this requirement is more stringent in that the *last* L observations should be observed for all individuals.

In what follows, we assume that some data points (and, as a result, some of the moment conditions) are missing at random; that is, we assume that the commonly used *missing-at-random* assumption can be justified. To save space, we restrict our attention to the factor proxies constructed using multiple \mathbf{w}_i and a single \mathbf{v}_i .

Suppose that for each cross-sectional unit i the econometrician has data on the availability of the moment conditions (and weights), summarized by the $[(\zeta + L) \times 1]$ dimensional vector $\mathbf{e}_i = (\mathbf{e}'_{i,\zeta}, \mathbf{e}'_{i,w})'$. \mathbf{e}_i is a vector of indicator variables that take the value of 1 if the corresponding moment condition, or weight, is available and zero otherwise. We assume that all \mathbf{e}_i are i.i.d. draws from some distribution function that does not depend on data unconditionally.

The distinction between missing moment conditions ($\mathbf{e}_{i,\zeta}$) and missing approximated factors ($\mathbf{e}_{i,w}$) related to \mathbf{w}_i is made for the following reason: while the structure of $\mathbf{e}_{i,\zeta}$ depends on the observability of all data points for individual i , the structure of $\mathbf{e}_{i,w}$ is primarily determined by the observability of \mathbf{w}_i .

That is, even if for some individual i the vector \mathbf{w}_i is not observed, such that $\mathbf{e}_{i,w} = \mathbf{0}_\zeta$, the ‘nuisance’ parameter matrix $\mathbf{G}_{z,\lambda}$ can still be estimated using the subgroup of cross-sectional units for which \mathbf{w}_i is observed. As a result, the number of observations used in estimation to identify $\boldsymbol{\beta}$ can be different from that used to identify the nuisance parameters.

As an example, consider $w_i = y_{i,0}$ in the one factor model. If this initial observation is missing for individual i then $\mathbf{e}_{i,\zeta}$ has T zero entries, whereas $\mathbf{e}_{i,w}$ reduces to a scalar that equals zero; however, the remaining $(y_{i,1}, \dots, y_{i,T})$ observations for unit i can be still used in estimation to identify $\boldsymbol{\beta}$.

Thus, letting

$$\mathbf{N}_\zeta = \sum_{i=1}^N \mathbf{e}_{i,\zeta}; \quad \mathbf{N}_w = \sum_{i=1}^N \mathbf{e}_{i,w},$$

and defining \mathbf{N}_ζ^+ (and similarly \mathbf{N}_w^+) to be the $[\zeta \times 1]$ vector with typical k -th element given by the reciprocal of the k -th element of \mathbf{N}_ζ , the resulting vector of estimating equations in the case of unbalanced panels can be written as

$$\bar{\boldsymbol{\mu}}_N(\boldsymbol{\theta}) = \mathbf{N}_\zeta^+ \odot \sum_{i=1}^N \boldsymbol{\mu}_{i,\zeta}(\boldsymbol{\theta}) - \mathbf{N}_w^+ \odot \sum_{i=1}^N \boldsymbol{\mu}_{i,w}(\boldsymbol{\theta}), \quad (\text{S1})$$

where \odot denotes the pointwise (Hadamard) product,

$$\boldsymbol{\mu}_{i,\zeta}(\boldsymbol{\theta}) = \mathbf{e}_{i,\zeta} \odot (\mathbf{Z}'_i(\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta})), \quad (\text{S2})$$

and

$$\boldsymbol{\mu}_{i,w}(\boldsymbol{\theta}) = \mathbf{e}_{i,w} \odot (\mathbf{S}[(\mathbf{v}_i\mathbf{w}'_i) \otimes \mathbf{I}_d]\mathbf{g}). \quad (\text{S3})$$

Using this definition for $\bar{\boldsymbol{\mu}}_N(\boldsymbol{\theta})$, the GMM objective function is obtained in the same way as with balanced data.

S1.2. Observed Factors

Often the empirical practitioner may wish to include variables that are common across individual entities, such as interest rates, unemployment rate etc., which can be viewed as observed common factors. It turns out that the extension of our estimator to such scenario is straightforward. To illustrate, consider the following model:

$$y_{i,t} = \mathbf{x}'_{i,t}\boldsymbol{\beta} + \boldsymbol{\phi}'_i\mathbf{h}_t + \boldsymbol{\lambda}'_i\mathbf{f}_t + \varepsilon_{i,t}; \quad i = 1, \dots, N; \quad t = 1, \dots, T, \quad (\text{S4})$$

where \mathbf{h}_t denotes a $[n \times 1]$ vector of observed factors and $\boldsymbol{\phi}_i$ is the corresponding vector of (unobserved) factor loadings, whereas all other variables have been already defined. In vector form we obtain

$$\mathbf{y}_i = \mathbf{X}_i\boldsymbol{\beta} + \mathbf{H}\boldsymbol{\phi}_i + \mathbf{F}\boldsymbol{\lambda}_i + \varepsilon_i, \quad (\text{S5})$$

where \mathbf{H} has dimension $[T \times n]$. Under standard assumptions on $\boldsymbol{\phi}_i$ it can be shown that

$$\text{E}_{\mathcal{F}}[\mathbf{Z}'_i(\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta}_0 - \mathbf{H}\boldsymbol{\phi}_i - \mathbf{F}\boldsymbol{\lambda}_i)] = \mathbf{S}(\text{vec}(\text{E}_{\mathcal{F}}[\mathbf{z}_i(\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta}_0)'] - \mathbf{G}_{z,\phi}\mathbf{H}' - \mathbf{G}_{z,\lambda}\mathbf{F}')) = \mathbf{0}_{\zeta}, \quad (\text{S6})$$

where $\mathbf{G}_{z,\phi} = \text{E}_{\mathcal{F}}[\mathbf{z}_i\boldsymbol{\phi}'_i]$ is a $[d \times n]$ unknown population matrix that absorbs the unobserved covariances between instruments and factor loadings of the observed factors. As a result, the vector of estimating equations can be redefined as follows:

$$\bar{\boldsymbol{\mu}}_N(\boldsymbol{\theta}) \equiv \frac{1}{N} \sum_{i=1}^N \mathbf{Z}'_i(\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta}) - \mathbf{S}(\widehat{\mathbf{F}}_e \otimes \mathbf{I}_d)\mathbf{g} + \mathbf{S}(\mathbf{H} \otimes \mathbf{I}_d)\mathbf{g}_{\phi}, \quad (\text{S7})$$

where $\mathbf{g}_{\phi} = \text{vec}(\mathbf{G}_{z,\phi})$, and \mathbf{g} is defined as in the main text. One can proceed by minimising the GMM objective function in exactly the same way as previously.

S1.3. Consistency Under an Alternative Set of Assumptions

For convenience we repeat Assumption 2.1 from the main text:

Assumption S1.1. *The DGP for all i and t satisfies the following restrictions:*

- (a) $(\mathbf{X}_i, \boldsymbol{\varepsilon}_i, \boldsymbol{\lambda}_i)$ are identically distributed and independent across i , conditional on \mathcal{F} .
- (b) Each time-varying element $p_{i,t}^{(\cdot)}$ in $\mathbf{p}_{i,t} = \left(p_{i,t}^{(1)}, \dots, p_{i,t}^{(K+1)}\right)' \equiv (\mathbf{x}'_{i,t}, \varepsilon_{i,t})'$ satisfies $\mathbb{E} \left[\left| p_{i,t}^{(\cdot)} \right|^{4+\delta} \right] < \infty$ for all t .
- (c) Each time-invariant element $\lambda_i^{(\cdot)}$ in $\boldsymbol{\lambda}_i = (\lambda_i^{(1)}, \dots, \lambda_i^{(L)})'$ satisfies $\mathbb{E} \left[\left| \lambda_i^{(\cdot)} \right|^{4+\delta} \right] < \infty$.
- (d) $\mathbb{E}_{\mathcal{F}}[\varepsilon_{i,t} | \boldsymbol{\lambda}_i, \mathbf{x}_{i,1:\tau_1}^{(1)}, \dots, \mathbf{x}_{i,1:\tau_K}^{(K)}] = 0 \forall t$, for some positive integers τ_1, \dots, τ_K .

The above assumption implies that all moments of random variables do not depend on index i , i.e. individuals are identically distributed and independent across i , conditional on \mathcal{F} . In what follows, we show that this assumption can be relaxed, without affecting consistency of our estimator. To be specific, recall that the true parameter vector $\boldsymbol{\theta}_0 = (\boldsymbol{\beta}'_0, \mathbf{g}'_0)'$, where $\mathbf{g}_0 \equiv \text{vec} \left(\mathbf{G}_{z,\lambda_e} (\mathbf{A}_N^{-1})' \right)$ is a function of $\mathbb{E}_{\mathcal{F}}[\mathbf{z}_i \boldsymbol{\lambda}'_{i,e}]$.

Our next result, shows that $\widehat{\boldsymbol{\theta}}$ is consistent even if $\mathbb{E}_{\mathcal{F}}[\mathbf{z}_i \boldsymbol{\lambda}'_i] = \mathbf{G}_{z,\lambda}(i)$ is allowed to be individual-specific. Such possibility can be relevant if one treats time-invariant random variables (e.g. $\boldsymbol{\lambda}_i$), as a generic sequence of constants, as advocated in Hsiao et al. (2002). It is critical that in this case the following limit

$$\mathbf{G}_{z,\lambda} \equiv \text{plim}_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \mathbf{G}_{z,\lambda}(i), \quad (\text{S8})$$

is well defined. Below we summarize the main result of this section for $\boldsymbol{\theta}_0 = (\boldsymbol{\beta}'_0, \mathbf{g}'_0)'$ with $\mathbf{g}_0 = \text{vec}(N^{-1} \sum_{i=1}^N \mathbf{G}_{z,\lambda_e}(i) (\mathbf{A}_N^{-1})')$. Notice that in this case $\mathbf{G}_{z,\lambda_e}(i)$ is defined equivalently to \mathbf{G}_{z,λ_e} , in the main text.

Theorem S1. *Suppose that Assumptions 2.1, 3.1 and 3.2 in the main text hold true, and Eq. (S8) is well defined. Then, for $N \rightarrow \infty$, we have*

$$\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0 \xrightarrow{p} \mathbf{0}. \quad (\text{S9})$$

The proof of this theorem is provided in the corresponding section of this Appendix.

S1.4. Relative Efficiency Comparison with the CCE Estimator

As it was mentioned in the main text, comparing the relative efficiency between the CCE estimator and the proposed linear GMM estimator (when both are fixed T consistent) is not straightforward. Below we provide a stylized example with $K = 1$ and $L = L_e = 1$, and demonstrate that the linear GMM estimator with unity weighting matrix has asymptotic variance equal to that of the CCE estimator.

Let the DGP for $y_{i,t}$ be as in the main text with $K = 1$ and $L = 1$, i.e.

$$\mathbf{y}_i = \mathbf{x}_i\beta + \mathbf{f}\lambda_i + \boldsymbol{\varepsilon}_i, \quad (\text{S10})$$

and $\boldsymbol{\varepsilon}_i$ is homoskedastic and uncorrelated over time such that $\boldsymbol{\Sigma}_\varepsilon \equiv \text{E}[\boldsymbol{\varepsilon}_i\boldsymbol{\varepsilon}_i'] = \sigma_\varepsilon^2\mathbf{I}_T$. Furthermore, we assume that the DGP for $x_{i,t}$ is of the following form:

$$\mathbf{x}_i = \mathbf{f}\gamma_i + \mathbf{u}_i, \quad (\text{S11})$$

where \mathbf{x}_i is strictly exogenous with respect to $\boldsymbol{\varepsilon}_i$, in particular \mathbf{u}_i and $\boldsymbol{\varepsilon}_i$ are assumed to be independent of each other, as well as of all other stochastic quantities. A single factor proxy $\widehat{\mathbf{f}}$ is constructed as the cross-sectional average of \mathbf{x}_i , i.e. $\mathbf{v}_i = \mathbf{x}_i$ and $w_i = 1$ such that

$$\widehat{\mathbf{f}} = \bar{\mathbf{x}}. \quad (\text{S12})$$

Finally, denote by $\widehat{\beta}_{CCE}$ the pooled CCE estimator, and by $\widehat{\beta}_{GMM1}$ our proposed GMM estimator with $\boldsymbol{\Omega} = \mathbf{I}$, $\mathbf{z}_i = \mathbf{x}_i$ and $\mathbf{S} = \mathbf{I}_{T^2}$, i.e. all leads and lags of $x_{i,t}$ are used as instruments. Then the following result is obtained:

Proposition S1. *Consider the DGP in Eq. (S10)-(S11). Suppose that $\text{E}_{\mathcal{F}}[\gamma_i] = \gamma \neq 0$ and $\boldsymbol{\Sigma}_u \equiv \text{E}[\mathbf{u}_i\mathbf{u}_i'] = \sigma_u^2\mathbf{I}_T$. Then,*

$$\text{Avar}(\widehat{\beta}_{CCE}) = \text{Avar}(\widehat{\beta}_{GMM1}) = \frac{1}{T-1} \frac{\sigma_\varepsilon^2}{\sigma_u^2}. \quad (\text{S13})$$

The proof of this proposition is provided in the corresponding section of this Appendix. Upon inspecting the proof of this proposition it can be seen that sphericity of both $\boldsymbol{\varepsilon}_i$ and \mathbf{u}_i is essential for the final result. Without this restriction, any comparison is not straightforward, as the variance of the CCE estimator depends on the level of $\boldsymbol{\Sigma}_u$, whereas for the GMM estimator the higher order power of $\boldsymbol{\Sigma}_u$ plays an important role. Finally, notice that here we consider a highly stylized example with a diagonal weighting matrix $\boldsymbol{\Omega}$. The analytical comparison between the two approaches becomes nearly impossible for more general weighting matrices that do not have a Kronecker structure of the form $\boldsymbol{\Omega} = \mathbf{I}_T \otimes \boldsymbol{\Omega}_T$.

S2. Descriptive Statistics of Water Usage Data

The following table presents some descriptive statistics for the data. The average and median daily water usage in the sample equals roughly .567 and .515 kL, respectively. This indicates that water consumption is skewed to the right, which is plausible because there is no upper bound in water consumption (loosely speaking). The between standard deviation of daily water usage is larger than the within standard deviation, which implies that there is more variation in water consumption across households than over time, as expected. Not surprisingly, the same holds true for temperature. Interestingly, the opposite is true for rainfall, i.e. there appears to exist more variation in rainfall amounts over time than across households within the sample. Finally, for the price variable the between standard deviation is much smaller than the within deviation, i.e. the largest proportion of variation in price is due to the consecutive, year by year, upward changes set by the NSW Independent Pricing and Regulatory Tribunal (IPART). This implies that endogeneity of the price variable may not be such a big issue.

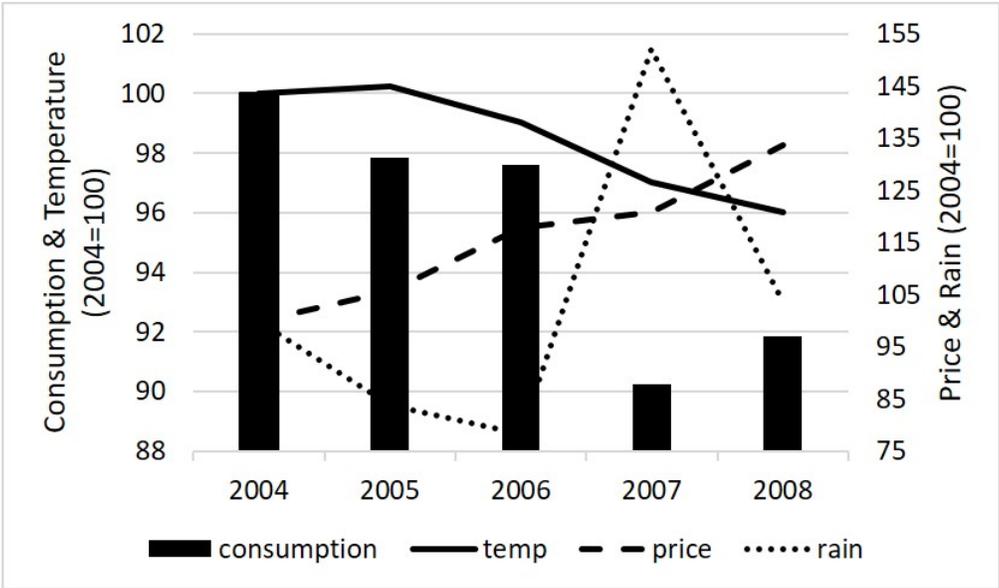
Table S1: Descriptive Statistics

		mean	median	st.dev.	10th perc.	90th perc.
logged cons.	overall	.567	.515	.328	.203	.977
	between			.303	-.117	.118
	within			.126	.219	.959
rain	overall	2.36	2.16	.739	1.54	3.50
	between			.350	-2.08	2.91
	within			.651	.681	1.05
temperature	overall	23.4	23.7	1.13	21.8	24.5
	between			1.04	21.4	24.2
	within			.437	-.601	.493
price	overall	1.35	1.37	.140	1.17	1.56
	between			.013	1.34	1.36
	within			.139	1.17	1.56

Figure S1 depicts the values of the cross-sectional averages of the variables included in the model, setting their corresponding 2004 values equal to 100. Therefore, the values of the variables from 2005 onwards are essentially percentage changes relative to the base year. To enhance visualization of the data, the values of water usage and

temperature are plotted with reference to the left vertical axis, while those of price and rain are plotted with reference to the right vertical axis. For instance, average water usage in 2007 was roughly 10% lower than 2004, and average price in 2008 was roughly 34% higher than 2004. During the period of our analysis, average daily temperature has followed a downward trend overall, whereas prices have steadily gone upwards every single year. At the same time, water usage experienced a significant drop in 2007 and remained much lower in 2008 relative the previous years. The rainfall variable exhibits large fluctuations over time, especially during 2006-2008, where a very dry year in 2006 was followed by a very wet one in 2007, followed again by a relatively dry year in 2008.

Figure S1: Water Consumption, Price and Temperature



S3. Further Finite Sample Evidence

S3.1. The Growth Ratio Statistic

In addition to the ER-statistics defined in the main paper, L_e can also be determined based on the Growth Ratio (GR) statistic of Ahn and Horenstein (2013). In particular, let $r_{\max} = \min(T, R) - 2$. The estimator of L_e based on the GR-statistic is defined as

$$\widehat{L}_e = \arg \max_{r \in \{1, \dots, r_{\max}\}} GR(r); \quad GR(r) = \frac{\ln(V(r-1)/V(r))}{\ln(V(R)/V(r+1))}, \quad (\text{S14})$$

where $V(r) = \sum_{k=r+1}^{\min(T,R)} \lambda_k(T^{-1} \widehat{\mathbf{F}}_R \widehat{\mathbf{F}}_R')$. Note that, by construction, the r_{\max} for the ER-statistic is always larger than that of the GR-statistic, due to the structure of the $V(r)$ function. Finally, one could also combine the ER and the GR criteria in the following hybrid form:

$$\widehat{L}_e = \arg \max_{r \in \{1, \dots, r_{\max}\}} GR^{(*)}(r); \quad GR^{(*)}(r) = V(r)/V(r-1), \quad (\text{S15})$$

where $r_{\max} = \min(T, R) - 1$, as in the case of the original ER statistic.

Remark 1. Note that the $V(r)$ function can also be used to construct a BIC criterion to select L_e . In particular, one may specify

$$S_N(r) = N \times N(V(r))^2 - \ln(N) \times h(r), \quad (\text{S16})$$

or

$$S_N(r) = N \times N(V(r)/V(r-1))^2 - \ln(N) \times h(r), \quad (\text{S17})$$

where $h(r) = C \times ((T-r)(R-r)) = \mathcal{O}(1)$, a strictly decreasing function of r with $0 < C < \infty$. Under the set of our assumptions, one can show that this BIC criterion can consistently estimate L_e . A detailed investigation of the finite-sample performance of this selection criterion, which needs also calibrate appropriate choices of C , is left for future research.

Table S2 presents model selection results on L_e for the GR-statistic, based on the setup considered in the main text. For simplicity, we restrict our attention to $T = 4$, as the results for $T = 8$ are mostly identical to those of the ER-statistic, which are reported in the main text. As we can see from Table S2, the gains of using the GR-statistic are especially pronounced in the two factor setup for smaller values of N . For example, for $N = 200$, $\alpha = 0.4$ and $\delta = 0$, the GR-statistic selects two factors in

87% of the replications, while the corresponding number for ER-statistic is 76% (84% if one combines $\widehat{L}_e = 2$ and $\widehat{L}_e = 3$). The main benefit of using the GR-statistic is mostly visible for smaller values of N , while in other setups the performance of the GR-statistic is similar to that of the ER-statistic. Overall, we may conclude that (with few exceptions) the GR-statistic would be preferred, when it is available.

S3.2. Residual Based Estimation of L

Given our estimator $\widehat{\boldsymbol{\beta}}$, define $\widehat{\mathbf{q}}_N = N^{-1} \sum_{i=1}^N \mathbf{Z}'_i(\mathbf{y}_i - \mathbf{X}_i\widehat{\boldsymbol{\beta}})$. Then $\text{plim}_{N \rightarrow \infty} \widehat{\mathbf{q}}_N = \mathbf{S} \text{vec}(\mathbf{G}_{z,\lambda} \mathbf{F}')$. Observe that the rank of the matrix inside the $\text{vec}(\cdot)$ operator is exactly L . However, depending on \mathbf{S} , only some elements of $\mathbf{G}_{z,\lambda} \mathbf{F}'$ can be consistently estimated. For instance, if $\mathbf{S} = \mathbf{I}$ we can recover fully the $\mathbf{G}_{z,\lambda} \mathbf{F}'$ component, thus the un-vec'd version of \mathbf{q}_N can be used as an input in the ER- (or GR-) statistic. If not all instruments are valid in all time periods (e.g. due to weak-exogeneity), then only a sub-matrix of $\mathbf{G}_{z,\lambda} \mathbf{F}'$ can be estimated consistently. For example, provided that a subset of \tilde{d} elements of \mathbf{z}_i are valid for \tilde{T} time periods, an un-vec'd $[\tilde{d} \times \tilde{T}]$ sub-matrix of $\widehat{\mathbf{q}}_N$ can be used to consistently estimate L , provided that $L < \min(\tilde{d}, \tilde{T})$. \widehat{L} estimated this way can be later used as an input for the non-linear estimators of Robertson and Sarafidis (2015) and Ahn et al. (2013).

In what follows, we illustrate how such estimator of L , based on the second-step ‘‘Fr’’ estimator as the plug-in estimator of $\boldsymbol{\beta}$ can be constructed for the DGP considered. In particular, take $\widehat{\boldsymbol{\varepsilon}}_i = \mathbf{y}_i - \mathbf{X}_i\widehat{\boldsymbol{\beta}}$ and $\mathbf{z}_i = (\mathbf{y}'_{i,-1}, \mathbf{x}'_i)'$. One can estimate L by constructing the following term:

$$\widehat{\mathbf{Q}}_N = \mathbf{S}_L \frac{1}{N} \sum_{i=1}^N \widehat{\boldsymbol{\varepsilon}}_i ((\mathbf{S}_R \otimes \mathbf{I}_2) \mathbf{z}_i)', \quad (\text{S18})$$

where \mathbf{S}_L is a selection matrix that picks the last $\lceil (T+1)/2 \rceil$ rows of the matrix that follows, while \mathbf{S}_R is a selection matrix that picks the first $\lceil T/2 \rceil$ elements of $\mathbf{y}_{i,-1}$ and \mathbf{x}_i , respectively. For example, for $T = 4$ $\widehat{\mathbf{Q}}_N$ is of the form:

$$\widehat{\mathbf{Q}}_N = \frac{1}{N} \sum_{i=1}^N \begin{pmatrix} \widehat{\boldsymbol{\varepsilon}}_{i,2} \\ \widehat{\boldsymbol{\varepsilon}}_{i,3} \\ \widehat{\boldsymbol{\varepsilon}}_{i,4} \end{pmatrix} \begin{pmatrix} y_{i,0} & y_{i,1} & x_{i,1} & x_{i,2} \end{pmatrix}, \quad (\text{S19})$$

while for $T = 8$

$$\widehat{\mathbf{Q}}_N = \frac{1}{N} \sum_{i=1}^N \begin{pmatrix} \widehat{\varepsilon}_{i,4} \\ \widehat{\varepsilon}_{i,5} \\ \widehat{\varepsilon}_{i,6} \\ \widehat{\varepsilon}_{i,7} \\ \widehat{\varepsilon}_{i,8} \end{pmatrix} \begin{pmatrix} y_{i,0} & y_{i,1} & y_{i,2} & y_{i,3} & x_{i,1} & x_{i,2} & x_{i,3} & x_{i,4} \end{pmatrix}. \quad (\text{S20})$$

$\widehat{\mathbf{Q}}_N$ has a total $O(T^2)$ number of elements, i.e. the same order as the total number of available moment conditions. Note that there are other ways of creating $\widehat{\mathbf{Q}}_N$. For example, one can fix the dimensions of either \mathbf{S}_L or \mathbf{S}_R , such that $\widehat{\mathbf{Q}}_N$ selects a fixed number of rows (columns) independently of T . In this case the total number of the elements can be of order $O(T)$. However, we believe that our specific choice is more natural, given that the total number of moment conditions available is of order $O(T^2)$.

The results are summarized in Table S3. Notice that for $T = 4$, the $\widehat{\mathbf{Q}}_N$ matrix can have at most rank 3, thus the ER criterion can only be used to choose between $L_e = 1$ and $L_e = 2$ (unless an appropriate mock eigenvalue strategy is used). From Table S3, we can see that for $L = L_e = 1$ the ER statistic performs well. This is true even for smaller values of N and T . In the two factor setup $L = L_e = 2$, the performance of the ER criterion deteriorates substantially. In particular, it appears that the values of N and T need to be substantially larger in order to ensure that the true number of factors is selected in the majority of the cases. We conjecture that this mainly due to the nature of the Monte Carlo setup considered at present, where the second factor in $y_{i,t}$ has mean zero and is uncorrelated with the factors in $x_{i,t}$. The implication is that the $\mathbf{G}_{z,\lambda}$ matrix has many zero or small values, which renders the second largest eigenvalue relatively small as compared to the largest. Limited (unreported) Monte Carlo results suggest that selection precision in the $L = L_e = 2$ setup improves substantially when the Growth Ratio criterion is used.

S3.3. $L = 1$ and $L_e = 2$

This section considers the same setup as in the main paper, however we allow for the possibility that factor proxies asymptotically identify more factors than present in the equation for $y_{i,t}$. In particular, we set

$$\lambda_{2,i}^y = 0; \quad \lambda_{2,i}^{v^2} \sim \mathcal{N}(1, 1), \quad (\text{S21})$$

for all i , such that $L = 1$ and $L_e = 2$. Note that in this case $P = 1$. Thus, the regularized GMM estimator will identify a two-dimensional factor space asymptotically, whereas

consistent estimation of $\boldsymbol{\beta}$ requires only identification of a one-dimensional subspace.

The estimation results are summarized in Tables S4-S5. As we can see, the performance of all four estimators is more than satisfactory in terms of both bias and RMSE. Fr has slightly larger dispersion compared to the ‘oracle’ estimator F1, which is expected given that Fr estimates one extra factor. For $T = 4$ Fbic slightly outperforms Fr in terms of RMSE, however these differences are eliminated for $T = 8$. The results on empirical size of all four estimators are similar to those obtained in the main text and therefore one can draw identical conclusions. The same applies for the properties of the J-statistic and model selection.

S3.4. Zero Mean Factor Loadings

This section studies the setup when the factor loadings corresponding to the first factor $f_{1,t}$ have zero mean for all variables of interest. In such case, factor proxies based solely on cross-sectional averages do not consistently estimate the column space of \mathbf{F} . Therefore, F1 and F2 become inconsistent and have a random “weak-instruments” limit.

The results are summarized in Tables S6-S7 and S8-S9 for $L = L_e = 1$ and $L = L_e = 2$, respectively. As expected, F1 and F2 perform poorly, although the J-statistic appears to have adequate power to detect violations of the null. On the other hand, Fr and Fbic continue to perform satisfactorily in term of bias, RMSE and size.

In comparison to the case where $\mu_\lambda = 1$, the performance of Fr and Fbic moderately deteriorates. The main reason is that when the factor loadings corresponding to the first factor have mean zero, the identification strength of our factor proxies is crucially determined by the value of $f_{1,0}$, which in our setup is drawn from the normal distribution with mean zero and unit variance. To see this, notice that for $L_e = 2$ the only informative factor proxies can be constructed from $\mathbf{v}_i^{(2)}$ and $\mathbf{w}_i = (1, y_{i,0})'$. In particular, the expectation of the product of these two variables is given by

$$\mathbf{E}_{\mathcal{F}}[\mathbf{v}_i^{(2)} \mathbf{w}_i'] = \mathbf{F} \mathbf{E}_{\mathcal{F}}[\lambda_i^2(1, y_{i,0})] = \mathbf{F} \begin{pmatrix} 0 & 0.6f_{1,0} \\ 1 & 0 \end{pmatrix}, \quad (\text{S22})$$

which implies that $\text{rk}(\mathbf{G}_{\gamma,w}) = 2$ a.s. However, given that $P(|0.6f_{1,0}| < 0.2) \approx 0.26$, the correlation determined by $\mathbf{G}_{\gamma,w}$ will be weak in a substantial fraction of replications. The same conclusion applies to the model with one factor, in which case we have

$$\mathbf{E}_{\mathcal{F}}[\mathbf{v}_i^{(2)} \mathbf{w}_i'] = \mathbf{F} \mathbf{E}_{\mathcal{F}}[\lambda_i^2(1, y_{i,0})] = \mathbf{F} \begin{pmatrix} 0 & 0.6f_{1,0} \end{pmatrix}. \quad (\text{S23})$$

Clearly, the value of $f_{1,0}$ determines the identification strength of our factor proxies. Using auxiliary simulations (not presented here), we have verified that when $f_{1,0}$ is fixed and sufficiently away from zero (e.g. $f_{1,0} = 1$), identification strength increases substantially.

In conclusion, the DGP of this section is generating semi-strong, rather than strong, identification for Fr. We believe that this case resembles close to the worst-case scenario. This is because, from the empirical point of view, it is highly unlikely that both $E[y_{i,t}] = 0$ and $E[x_{i,t}] = 0$ for all $t = 0, \dots, T$, which is precisely what the present setup implies.

S3.5. “Fr” with Estimated L_e

In the main text we presented results for “Fr”, the GMM estimator that uses the regularized version of $\widehat{\mathbf{F}}_R$ and treats L_e as known. Here we relax this assumption and instead consider the performance of Fr based on \widehat{L}_e , which is obtained using the ER-statistic. The corresponding estimator is denoted as “Fr+”.

The estimation results are summarized in Tables S10-S11. As we can see, the performance of Fr+ is very close to that of Fr and thereby it remains entirely satisfactory in all cases.

S3.6. Estimation with Additional Regressor

To illustrate how our proposed procedures scale up with an additional (weakly-) exogenous regressor, we consider the following ARDL(1,1) model:

$$y_{i,t} = \alpha y_{i,t-1} + \beta x_{i,t} + \beta_1 x_{i,t-1} + \sum_{r=1}^2 \lambda_{r,i}^y f_{r,t} + \varepsilon_{i,t}^y; \quad x_{i,t} = \delta y_{i,t-1} + \alpha_x x_{i,t-1} + \lambda_{1,i}^x f_{1,t} + \varepsilon_{i,t}^x.$$

For the sake of exposition we specify $\beta_1 = 0$, such that the population parameter values and parametrisation employed in the main text can be interpreted in exactly the same way here. The main difference, however, is that in this section all GMM estimators do not use knowledge of $\beta_1 = 0$, and instead estimate this coefficient. The results, summarized in Tables S12-S13 and S14-S15, are similar to those presented in the main text. The slight increase in RMSE is expected given that all estimators estimate an extra parameter that is redundant.

S3.7. Nonlinear GMM Estimators

This section investigates the performance of the nonlinear GMM estimators of Ahn et al. (2013) and Robertson and Sarafidis (2015), hereafter “ALS” and “FIVU” respectively. We focus on the one factor design, because the implementation of these non-linear estimators require larger values of T to be feasible. In particular, for $T = 4$ only the model with one factor can be evaluated. The results are summarized in Tables S16-S17. These results can serve as a rough point of reference, however they are not directly comparable with the linear GMM estimators proposed in the present paper, since they all use different moment conditions. The conclusions about the performance of the nonlinear GMM estimators are qualitatively similar to those discussed in Juodis and Sarafidis (2018). In particular, given the DGP we consider, the statistical properties of the FIVU estimator appear to be more satisfactory than those of the ALS

estimator. The interested reader is referred to Juodis and Sarafidis (2018) for a more detailed discussion.

Table S2: Selection Rates of L_e using the Growth Ratio statistic

Designs				$L_e = 1$		$L_e = 2$	
N	T	α	δ	$\#\hat{L}_e = 1$	$\#\hat{L}_e = 2$	$\#\hat{L}_e = 1$	$\#\hat{L}_e = 2$
200	4	0.4	0.0	0.96	0.04	0.13	0.87
200	4	0.4	0.3	0.96	0.04	0.13	0.87
200	4	0.8	0.0	0.95	0.05	0.13	0.87
200	4	0.8	0.3	0.96	0.04	0.12	0.88
800	4	0.4	0.0	0.97	0.03	0.04	0.96
800	4	0.4	0.3	0.98	0.02	0.05	0.95
800	4	0.8	0.0	0.99	0.01	0.04	0.96
800	4	0.8	0.3	0.98	0.02	0.04	0.96

Table S3: Selection Rates of L using the Eigenvalue Ratio statistic.

Designs				$L = 1$			$L = 2$		
N	T	α	δ	$\#\hat{L} = 1$	$\#\hat{L} = 2$	$\#\hat{L} = 3$	$\#\hat{L} = 1$	$\#\hat{L} = 2$	$\#\hat{L} = 3$
200	4	0.4	0.0	.92	.08	NA	.44	.56	NA
200	4	0.4	0.3	.95	.05	NA	.44	.56	NA
200	4	0.8	0.0	.83	.17	NA	.47	.53	NA
200	4	0.8	0.3	.88	.12	NA	.47	.53	NA
200	8	0.4	0.0	.99	.00	.00	.43	.55	.02
200	8	0.4	0.3	1	.00	.00	.39	.61	.01
200	8	0.8	0.0	.98	.00	.02	.50	.46	.05
200	8	0.8	0.3	.99	.00	.01	.48	.51	.02
800	4	0.4	0.0	.98	.02	NA	.34	.66	NA
800	4	0.4	0.3	.98	.02	NA	.32	.68	NA
800	4	0.8	0.0	.93	.08	NA	.41	.59	NA
800	4	0.8	0.3	.94	.06	NA	.41	.59	NA
800	8	0.4	0.0	1	.00	.00	.27	.73	.00
800	8	0.4	0.3	1	.00	.00	.23	.77	.00
800	8	0.8	0.0	1	.00	.00	.33	.66	.01
800	8	0.8	0.3	1	.00	.00	.30	.70	.00

Table S4: Estimation results for $L = 1$ and $L_e = 2$ setup.

Designs				F1				F2				Fr				Fbic			
N	T	α	δ	Bias	RMSE	Std	Size												
α																			
200	4	4.0		.00	.02	.02	.06	.00	.04	.04	.06	.00	.04	.04	.06	.00	.03	.03	.07
200	4	4.3		.00	.03	.03	.07	.00	.06	.06	.05	.00	.06	.06	.05	.00	.05	.05	.08
200	4	8.0		.00	.03	.03	.07	.00	.05	.05	.07	.00	.04	.04	.06	.00	.03	.03	.08
200	4	8.3		.00	.03	.03	.06	.00	.05	.05	.06	-.01	.05	.05	.05	.00	.04	.04	.08
200	8	4.0		.00	.02	.02	.11	.00	.02	.02	.09	.00	.02	.02	.10	.00	.02	.02	.12
200	8	4.3		-.01	.03	.03	.17	-.01	.03	.03	.13	-.01	.03	.03	.12	-.01	.04	.04	.18
200	8	8.0		.00	.01	.01	.12	.00	.02	.02	.09	.00	.02	.02	.09	.00	.02	.02	.12
200	8	8.3		.00	.02	.02	.13	.00	.02	.02	.11	.00	.02	.02	.11	-.01	.02	.02	.13
800	4	4.0		.00	.01	.01	.06	.00	.02	.02	.05	.00	.02	.02	.06	.00	.01	.01	.06
800	4	4.3		.00	.02	.02	.05	.00	.03	.03	.05	.00	.03	.03	.05	.00	.02	.02	.06
800	4	8.0		.00	.01	.01	.06	.00	.02	.02	.05	.00	.02	.02	.05	.00	.01	.01	.06
800	4	8.3		.00	.01	.01	.06	.00	.03	.03	.06	.00	.03	.03	.05	.00	.02	.02	.06
800	8	4.0		.00	.01	.01	.07	.00	.01	.01	.05	.00	.01	.01	.05	.00	.01	.01	.07
800	8	4.3		.00	.01	.01	.08	.00	.01	.01	.07	.00	.01	.01	.07	.00	.01	.01	.08
800	8	8.0		.00	.01	.01	.07	.00	.01	.01	.06	.00	.01	.01	.06	.00	.01	.01	.07
800	8	8.3		.00	.01	.01	.06	.00	.01	.01	.06	.00	.01	.01	.06	.00	.01	.01	.07
β																			
200	4	4.0		.00	.02	.02	.07	.01	.07	.07	.05	.01	.06	.06	.05	.00	.06	.06	.07
200	4	4.3		.00	.03	.03	.06	.00	.07	.07	.05	.00	.07	.07	.05	.00	.05	.05	.07
200	4	8.0		.00	.01	.01	.07	.00	.03	.03	.05	.00	.03	.03	.05	.00	.02	.02	.07
200	4	8.3		.00	.02	.02	.07	.00	.04	.04	.05	.00	.04	.04	.05	.00	.02	.02	.07
200	8	4.0		.00	.02	.02	.12	.00	.02	.02	.10	.00	.02	.02	.10	.01	.02	.02	.14
200	8	4.3		.01	.04	.04	.19	.01	.04	.03	.12	.01	.04	.03	.12	.01	.05	.04	.20
200	8	8.0		.00	.01	.01	.10	.00	.01	.01	.09	.00	.01	.01	.08	.00	.01	.01	.11
200	8	8.3		.00	.02	.02	.14	.00	.02	.02	.11	.00	.02	.02	.12	.01	.02	.02	.14
800	4	4.0		.00	.01	.01	.05	.00	.04	.04	.05	.00	.03	.03	.05	.00	.02	.02	.05
800	4	4.3		.00	.02	.02	.06	.00	.04	.04	.05	.00	.04	.04	.05	.00	.02	.02	.06
800	4	8.0		.00	.01	.01	.06	.00	.02	.02	.05	.00	.02	.02	.05	.00	.01	.01	.07
800	4	8.3		.00	.01	.01	.05	.00	.02	.02	.05	.00	.02	.02	.05	.00	.01	.01	.06
800	8	4.0		.00	.01	.01	.06	.00	.01	.01	.06	.00	.01	.01	.05	.00	.01	.01	.06
800	8	4.3		.00	.01	.01	.08	.00	.01	.01	.07	.00	.01	.01	.07	.00	.01	.01	.07
800	8	8.0		.00	.00	.00	.07	.00	.01	.01	.06	.00	.01	.01	.06	.00	.00	.00	.07
800	8	8.3		.00	.01	.01	.07	.00	.01	.01	.06	.00	.01	.01	.06	.00	.01	.01	.07

Table S5: Model selection and specification testing results for $L = 1$ and $L_e = 2$ setup.

Designs		J Statistic			BIC		ER L_e				
N	T	α	δ	F1	F2	Fr	$\#\hat{L} = 1$	$\#\hat{L} = 2$	$\#\hat{L}_e = 1$	$\#\hat{L}_e = 2$	$\#\hat{L}_e = 3$
200	4	.4	.0	.04	.04	.04	.97	.03	.18	.75	.07
200	4	.4	.3	.04	.04	.04	.97	.03	.19	.74	.07
200	4	.8	.0	.04	.04	.04	.98	.02	.18	.76	.07
200	4	.8	.3	.04	.03	.03	.97	.03	.16	.76	.07
200	8	.4	.0	.03	.03	.03	.97	.03	.07	.93	.00
200	8	.4	.3	.03	.03	.03	.97	.03	.07	.93	.01
200	8	.8	.0	.03	.03	.03	.97	.03	.07	.93	.01
200	8	.8	.3	.02	.03	.03	.97	.03	.06	.94	.01
800	4	.4	.0	.05	.05	.06	1	.01	.08	.89	.02
800	4	.4	.3	.04	.05	.05	.99	.01	.08	.89	.03
800	4	.8	.0	.04	.03	.04	1	.01	.09	.89	.03
800	4	.8	.3	.05	.05	.05	.99	.01	.08	.90	.02
800	8	.4	.0	.05	.05	.05	1	.00	.01	.99	.00
800	8	.4	.3	.05	.04	.05	1	.01	.01	.99	.00
800	8	.8	.0	.04	.04	.05	1	.00	.01	.99	.00
800	8	.8	.3	.05	.05	.05	1	.01	.01	.99	.00

Table S6: Estimation results for $L = 1$ and $L_e = 1$ setup with zero-mean factor loadings.

Designs		F1				F2				Fr				Fbic					
N	T	α	δ	Bias	RMSE	Std	Size	Bias	RMSE	Std	Size	Bias	RMSE	Std	Size	Bias	RMSE	Std	Size
α																			
200	4	4.0		-.03	.09	.08	.21	-.03	.19	.18	.06	.00	.03	.03	.08	.00	.05	.05	.08
200	4	4.3		-.09	.16	.14	.31	-.07	.25	.24	.08	-.01	.05	.05	.09	-.01	.07	.07	.09
200	4	8.0		-.02	.09	.09	.20	-.03	.19	.19	.05	.00	.03	.03	.08	.00	.04	.04	.07
200	4	8.3		-.03	.12	.11	.23	-.04	.23	.22	.05	.00	.04	.04	.08	-.01	.08	.08	.08
200	8	4.0		-.06	.09	.06	.50	-.03	.07	.06	.25	-.01	.02	.02	.16	-.01	.02	.02	.15
200	8	4.3		-.23	.27	.14	.82	-.15	.20	.13	.53	-.03	.07	.06	.27	-.04	.07	.06	.25
200	8	8.0		-.02	.06	.05	.32	-.01	.04	.04	.13	.00	.02	.02	.12	.00	.02	.02	.12
200	8	8.3		-.10	.15	.11	.52	-.06	.11	.09	.22	-.01	.04	.04	.15	-.01	.03	.03	.15
800	4	4.0		-.03	.08	.08	.30	-.04	.21	.21	.07	.00	.02	.02	.07	.00	.03	.03	.07
800	4	4.3		-.08	.16	.14	.38	-.08	.30	.29	.09	.00	.03	.03	.06	.00	.04	.04	.06
800	4	8.0		-.02	.08	.08	.28	-.03	.21	.20	.06	.00	.02	.02	.06	.00	.03	.03	.06
800	4	8.3		-.03	.12	.11	.29	-.05	.29	.28	.06	.00	.02	.02	.07	.00	.03	.03	.07
800	8	4.0		-.05	.07	.06	.51	-.03	.06	.05	.24	.00	.01	.01	.08	.00	.01	.01	.07
800	8	4.3		-.21	.25	.14	.80	-.14	.19	.13	.53	-.01	.04	.04	.13	-.01	.03	.03	.10
800	8	8.0		-.02	.05	.04	.36	-.01	.03	.03	.14	.00	.01	.01	.09	.00	.01	.01	.08
800	8	8.3		-.08	.13	.10	.52	-.04	.09	.08	.23	.00	.02	.02	.09	.00	.01	.01	.06
β																			
200	4	4.0		.04	.10	.09	.20	.04	.26	.26	.03	.01	.03	.03	.07	.01	.06	.06	.07
200	4	4.3		.09	.16	.13	.27	.07	.28	.27	.06	.01	.05	.05	.08	.01	.08	.08	.07
200	4	8.0		.01	.04	.04	.11	.01	.11	.11	.02	.00	.01	.01	.06	.00	.02	.02	.06
200	4	8.3		.02	.05	.05	.12	.02	.15	.15	.02	.00	.02	.02	.07	.00	.05	.05	.06
200	8	4.0		.09	.12	.08	.63	.05	.09	.07	.29	.01	.03	.03	.16	.01	.03	.03	.18
200	8	4.3		.25	.29	.15	.84	.17	.22	.14	.54	.04	.08	.07	.28	.04	.07	.06	.27
200	8	8.0		.02	.04	.03	.34	.01	.03	.02	.12	.00	.01	.01	.11	.00	.01	.01	.11
200	8	8.3		.09	.12	.09	.59	.05	.08	.07	.23	.01	.03	.03	.16	.01	.02	.02	.16
800	4	4.0		.04	.11	.10	.26	.06	.29	.28	.04	.00	.02	.02	.06	.00	.04	.04	.05
800	4	4.3		.08	.16	.13	.34	.07	.34	.33	.07	.00	.03	.03	.07	.00	.04	.04	.06
800	4	8.0		.01	.05	.05	.16	.01	.16	.16	.03	.00	.01	.01	.05	.00	.02	.02	.05
800	4	8.3		.01	.06	.06	.19	.02	.19	.18	.03	.00	.01	.01	.06	.00	.03	.03	.06
800	8	4.0		.07	.10	.07	.59	.04	.07	.06	.25	.00	.01	.01	.08	.00	.01	.01	.06
800	8	4.3		.22	.26	.15	.82	.15	.21	.14	.54	.01	.04	.04	.13	.01	.03	.03	.09
800	8	8.0		.02	.03	.02	.36	.01	.02	.02	.10	.00	.01	.01	.06	.00	.01	.01	.05
800	8	8.3		.06	.10	.08	.55	.04	.07	.06	.21	.00	.01	.01	.07	.00	.01	.01	.06

Table S7: Model selection and specification testing results for $L = 1$ and $L_e = 1$ setup.

Designs		J Statistic			BIC		ER L_e				
N	T	α	δ	F1	F2	Fr	$\#\hat{L} = 1$	$\#\hat{L} = 2$	$\#\hat{L}_e = 1$	$\#\hat{L}_e = 2$	$\#\hat{L}_e = 3$
200	4	.4	.0	.49	.10	.08	.97	.03	.71	.06	.23
200	4	.4	.3	.53	.10	.09	.97	.03	.72	.07	.21
200	4	.8	.0	.44	.09	.09	.97	.03	.71	.06	.22
200	4	.8	.3	.47	.07	.07	.97	.03	.72	.05	.23
200	8	.4	.0	.57	.25	.06	.91	.09	.87	.02	.12
200	8	.4	.3	.69	.35	.08	.89	.11	.86	.02	.12
200	8	.8	.0	.54	.20	.07	.91	.09	.86	.02	.13
200	8	.8	.3	.60	.26	.06	.92	.08	.86	.02	.12
800	4	.4	.0	.68	.16	.09	.98	.02	.84	.03	.13
800	4	.4	.3	.72	.18	.10	.98	.02	.85	.03	.13
800	4	.8	.0	.64	.14	.07	.99	.01	.84	.02	.14
800	4	.8	.3	.66	.14	.07	.99	.01	.83	.03	.14
800	8	.4	.0	.85	.55	.09	.94	.06	.94	.01	.06
800	8	.4	.3	.93	.73	.13	.91	.09	.93	.01	.07
800	8	.8	.0	.82	.44	.09	.94	.06	.93	.01	.07
800	8	.8	.3	.87	.56	.10	.92	.08	.94	.01	.06

Table S8: Estimation results for $L = 2$ and $L_e = 2$ setup with zero-mean factor loadings.

Designs				F1				F2				Fr				Fbic			
N	T	α	δ	Bias	RMSE	Std	Size												
α																			
200	4	4.0		-.04	.11	.11	.33	-.03	.15	.15	.16	.00	.06	.06	.06	-.01	.13	.13	.10
200	4	4.3		-.10	.20	.17	.41	-.07	.22	.21	.19	-.01	.09	.09	.07	-.02	.20	.20	.13
200	4	8.0		-.03	.12	.12	.35	-.03	.19	.19	.14	-.01	.06	.06	.06	-.01	.17	.17	.13
200	4	8.3		-.05	.15	.14	.39	-.04	.18	.18	.15	-.01	.08	.08	.06	-.02	.16	.16	.14
200	8	4.0		-.07	.10	.08	.59	-.04	.07	.06	.39	-.01	.03	.03	.13	-.01	.03	.03	.13
200	8	4.3		-.26	.31	.17	.80	-.18	.23	.14	.67	-.04	.08	.07	.22	-.04	.08	.07	.22
200	8	8.0		-.03	.08	.07	.48	-.02	.05	.05	.24	.00	.02	.02	.11	.00	.02	.02	.11
200	8	8.3		-.13	.18	.13	.66	-.07	.12	.10	.41	-.01	.03	.03	.15	-.01	.04	.03	.13
800	4	4.0		-.04	.11	.11	.47	-.04	.16	.15	.23	.00	.04	.04	.06	-.01	.09	.09	.08
800	4	4.3		-.11	.21	.18	.54	-.08	.22	.21	.29	-.01	.06	.06	.07	-.01	.17	.17	.08
800	4	8.0		-.03	.12	.12	.50	-.03	.17	.17	.21	.00	.04	.04	.06	.00	.10	.10	.09
800	4	8.3		-.04	.15	.15	.54	-.04	.22	.22	.21	.00	.05	.05	.06	-.01	.11	.11	.10
800	8	4.0		-.06	.10	.07	.66	-.04	.07	.05	.43	.00	.01	.01	.08	.00	.01	.01	.07
800	8	4.3		-.24	.29	.16	.84	-.16	.21	.13	.69	-.01	.04	.04	.13	-.01	.04	.03	.11
800	8	8.0		-.03	.07	.06	.56	-.01	.04	.04	.33	.00	.01	.01	.08	.00	.01	.01	.08
800	8	8.3		-.10	.16	.12	.68	-.06	.10	.09	.43	.00	.02	.02	.09	.00	.02	.02	.08
β																			
200	4	4.0		.05	.11	.10	.27	.07	.25	.24	.12	.01	.08	.08	.06	.02	.21	.21	.08
200	4	4.3		.10	.18	.15	.34	.08	.29	.28	.16	.01	.11	.11	.07	.02	.24	.24	.10
200	4	8.0		.01	.04	.04	.16	.02	.13	.13	.07	.00	.05	.05	.06	.00	.12	.12	.07
200	4	8.3		.02	.06	.06	.18	.02	.14	.14	.08	.01	.06	.06	.06	.01	.11	.11	.09
200	8	4.0		.11	.14	.10	.69	.08	.11	.08	.50	.01	.03	.03	.13	.01	.03	.03	.13
200	8	4.3		.28	.34	.19	.82	.20	.26	.16	.69	.04	.09	.08	.22	.04	.09	.08	.23
200	8	8.0		.03	.05	.04	.45	.02	.03	.03	.25	.00	.01	.01	.11	.00	.01	.01	.11
200	8	8.3		.11	.15	.10	.68	.07	.11	.08	.45	.01	.03	.03	.14	.01	.03	.03	.13
800	4	4.0		.05	.11	.10	.42	.07	.28	.27	.19	.01	.08	.08	.06	.02	.23	.23	.07
800	4	4.3		.10	.20	.17	.49	.09	.27	.25	.24	.01	.07	.07	.07	.00	.30	.30	.08
800	4	8.0		.01	.04	.04	.32	.02	.14	.14	.14	.00	.04	.04	.06	.01	.11	.11	.06
800	4	8.3		.02	.06	.06	.34	.03	.15	.15	.15	.00	.05	.05	.06	.01	.08	.08	.07
800	8	4.0		.09	.12	.08	.76	.06	.09	.07	.50	.00	.02	.02	.08	.00	.01	.01	.07
800	8	4.3		.26	.31	.17	.85	.18	.23	.15	.71	.01	.05	.05	.12	.01	.04	.04	.11
800	8	8.0		.02	.04	.03	.56	.01	.03	.02	.30	.00	.01	.01	.08	.00	.01	.01	.07
800	8	8.3		.09	.13	.09	.70	.05	.09	.07	.46	.00	.02	.02	.09	.00	.02	.01	.08

Table S9: Model selection and specification testing results for $L = 2$ and $L_e = 2$ setup.

Designs		J Statistic			BIC		ER L_e				
N	T	α	δ	F1	F2	Fr	$\#\hat{L} = 1$	$\#\hat{L} = 2$	$\#\hat{L}_e = 1$	$\#\hat{L}_e = 2$	$\#\hat{L}_e = 3$
200	4	.4	.0	.87	.28	.06	.27	.73	.30	.52	.18
200	4	.4	.3	.89	.30	.07	.27	.73	.30	.50	.21
200	4	.8	.0	.84	.25	.05	.33	.67	.29	.52	.19
200	4	.8	.3	.85	.24	.05	.33	.67	.29	.52	.19
200	8	.4	.0	.99	.53	.08	.04	.96	.30	.64	.06
200	8	.4	.3	.99	.66	.12	.05	.95	.30	.65	.05
200	8	.8	.0	.98	.50	.07	.05	.96	.30	.65	.05
200	8	.8	.3	.99	.55	.08	.04	.96	.29	.66	.05
800	4	.4	.0	.98	.47	.09	.09	.91	.26	.65	.09
800	4	.4	.3	.98	.52	.08	.08	.92	.25	.66	.09
800	4	.8	.0	.97	.41	.07	.12	.88	.25	.66	.09
800	4	.8	.3	.97	.43	.07	.10	.90	.24	.66	.10
800	8	.4	.0	1	.81	.11	.00	1	.22	.77	.01
800	8	.4	.3	1	.89	.14	.00	1	.23	.75	.02
800	8	.8	.0	1	.78	.10	.00	1	.22	.77	.01
800	8	.8	.3	1	.81	.12	.00	1	.22	.77	.01

Table S10: Estimation results for $L = 1$ and $L_e = 1$ setup for the “Fr” estimator based on \hat{L}_e .

Designs				α				β				
N	T	α	δ	Bias	RMSE	Std	Size	Bias	RMSE	Std	Size	J
200	4	.4	.0	.00	.02	.02	.06	.00	.03	.03	.07	.05
200	4	.4	.3	.00	.03	.03	.07	.00	.03	.03	.07	.04
200	4	.8	.0	.00	.02	.02	.07	.00	.01	.01	.06	.05
200	4	.8	.3	.00	.03	.03	.07	.00	.02	.02	.06	.04
200	8	.4	.0	.00	.01	.01	.09	.00	.02	.02	.11	.03
200	8	.4	.3	-.01	.03	.03	.15	.01	.03	.03	.16	.03
200	8	.8	.0	.00	.01	.01	.12	.00	.01	.01	.10	.03
200	8	.8	.3	.00	.02	.02	.13	.00	.01	.01	.12	.03
800	4	.4	.0	.00	.01	.01	.04	.00	.01	.01	.06	.05
800	4	.4	.3	.00	.02	.02	.06	.00	.02	.02	.05	.04
800	4	.8	.0	.00	.01	.01	.06	.00	.01	.01	.06	.04
800	4	.8	.3	.00	.01	.01	.05	.00	.01	.01	.05	.06
800	8	.4	.0	.00	.01	.01	.06	.00	.01	.01	.06	.04
800	8	.4	.3	.00	.01	.01	.06	.00	.01	.01	.06	.04
800	8	.8	.0	.00	.01	.01	.07	.00	.00	.00	.06	.05
800	8	.8	.3	.00	.01	.01	.07	.00	.01	.01	.06	.05

Table S11: Estimation results for $L = 2$ and $L_e = 2$ setup for the “Fr” estimator based on \hat{L}_e .

Designs				α				β				
N	T	α	δ	Bias	RMSE	Std	Size	Bias	RMSE	Std	Size	J
200	4	.4	.0	.00	.04	.04	.08	.00	.06	.06	.07	.15
200	4	.4	.3	.00	.06	.06	.08	.00	.07	.07	.07	.15
200	4	.8	.0	.00	.04	.04	.08	.00	.03	.03	.06	.14
200	4	.8	.3	.00	.05	.05	.09	.00	.04	.04	.05	.13
200	8	.4	.0	.00	.02	.02	.10	.00	.02	.02	.11	.08
200	8	.4	.3	.00	.03	.03	.13	.00	.04	.04	.13	.08
200	8	.8	.0	.00	.02	.02	.10	.00	.01	.01	.10	.09
200	8	.8	.3	.00	.02	.02	.12	.00	.02	.02	.11	.08
800	4	.4	.0	.00	.02	.02	.06	.00	.03	.03	.06	.11
800	4	.4	.3	.00	.03	.03	.07	.00	.04	.04	.06	.11
800	4	.8	.0	.00	.02	.02	.07	.00	.01	.01	.06	.10
800	4	.8	.3	.00	.02	.02	.08	.00	.02	.02	.06	.12
800	8	.4	.0	.00	.01	.01	.06	.00	.01	.01	.06	.06
800	8	.4	.3	.00	.01	.01	.07	.00	.01	.01	.06	.05
800	8	.8	.0	.00	.01	.01	.06	.00	.01	.01	.06	.05
800	8	.8	.3	.00	.01	.01	.07	.00	.01	.01	.07	.05

Table S12: Estimation results for $L = 1$ and $L_e = 1$ setup with $x_{i,t-1}$ as a regressor.

Designs		F1				F2				Fr				Fbic					
N	T	α	δ	Bias	RMSE	Std	Size	Bias	RMSE	Std	Size	Bias	RMSE	Std	Size	Bias	RMSE	Std	Size
α																			
200	4	.4	.0	.00	.04	.04	.06	.00	.09	.09	.02	.00	.04	.04	.06	.00	.05	.05	.07
200	4	.4	.3	.00	.04	.04	.06	.00	.09	.09	.02	.00	.04	.04	.06	.00	.05	.05	.07
200	4	.8	.0	.00	.03	.03	.06	.00	.06	.06	.02	.00	.03	.03	.06	.00	.04	.04	.07
200	4	.8	.3	.00	.03	.03	.07	.00	.08	.08	.02	.00	.03	.03	.07	.00	.04	.04	.07
200	8	.4	.0	.00	.03	.03	.12	.00	.03	.03	.07	.00	.02	.02	.12	.00	.03	.03	.12
200	8	.4	.3	-.01	.04	.04	.15	.00	.03	.03	.08	-.01	.03	.03	.14	-.01	.04	.04	.14
200	8	.8	.0	.00	.02	.02	.11	.00	.02	.02	.08	.00	.02	.02	.11	.00	.02	.02	.12
200	8	.8	.3	.00	.02	.02	.12	.00	.02	.02	.07	.00	.02	.02	.12	.00	.02	.02	.12
800	4	.4	.0	.00	.02	.02	.06	.00	.05	.05	.02	.00	.02	.02	.05	.00	.02	.02	.06
800	4	.4	.3	.00	.02	.02	.06	.00	.06	.06	.02	.00	.02	.02	.06	.00	.02	.02	.06
800	4	.8	.0	.00	.02	.02	.06	.00	.03	.03	.03	.00	.01	.01	.07	.00	.02	.02	.07
800	4	.8	.3	.00	.01	.01	.04	.00	.04	.04	.02	.00	.01	.01	.05	.00	.02	.02	.06
800	8	.4	.0	.00	.01	.01	.07	.00	.01	.01	.05	.00	.01	.01	.07	.00	.01	.01	.06
800	8	.4	.3	.00	.01	.01	.07	.00	.02	.02	.05	.00	.01	.01	.07	.00	.01	.01	.06
800	8	.8	.0	.00	.01	.01	.07	.00	.01	.01	.05	.00	.01	.01	.07	.00	.01	.01	.07
800	8	.8	.3	.00	.01	.01	.06	.00	.01	.01	.05	.00	.01	.01	.07	.00	.01	.01	.07
β																			
200	4	.4	.0	.00	.03	.03	.06	.01	.12	.12	.01	.00	.03	.03	.06	.00	.05	.05	.06
200	4	.4	.3	.00	.04	.04	.07	.00	.17	.17	.02	.00	.04	.04	.06	.00	.07	.07	.06
200	4	.8	.0	.00	.02	.02	.06	.00	.06	.06	.02	.00	.02	.02	.06	.00	.02	.02	.06
200	4	.8	.3	.00	.02	.02	.06	.00	.08	.08	.01	.00	.02	.02	.07	.00	.05	.05	.06
200	8	.4	.0	.01	.02	.02	.13	.00	.02	.02	.06	.00	.02	.02	.12	.01	.03	.02	.16
200	8	.4	.3	.02	.05	.04	.22	.01	.04	.04	.09	.01	.04	.04	.17	.03	.06	.05	.24
200	8	.8	.0	.00	.01	.01	.11	.00	.01	.01	.07	.00	.01	.01	.10	.00	.01	.01	.11
200	8	.8	.3	.01	.02	.02	.14	.00	.02	.02	.07	.00	.02	.02	.13	.01	.03	.03	.15
800	4	.4	.0	.00	.02	.02	.05	.00	.07	.07	.01	.00	.02	.02	.05	.00	.02	.02	.05
800	4	.4	.3	.00	.02	.02	.05	.00	.09	.09	.02	.00	.02	.02	.05	.00	.04	.03	.05
800	4	.8	.0	.00	.01	.01	.05	.00	.03	.03	.02	.00	.01	.01	.06	.00	.01	.01	.06
800	4	.8	.3	.00	.01	.01	.06	.00	.05	.05	.02	.00	.01	.01	.05	.00	.01	.01	.06
800	8	.4	.0	.00	.01	.01	.07	.00	.01	.01	.05	.00	.01	.01	.06	.00	.01	.01	.07
800	8	.4	.3	.00	.01	.01	.07	.00	.02	.02	.04	.00	.01	.01	.05	.00	.01	.01	.07
800	8	.8	.0	.00	.01	.01	.05	.00	.01	.01	.04	.00	.01	.01	.06	.00	.01	.01	.06
800	8	.8	.3	.00	.01	.01	.07	.00	.01	.01	.06	.00	.01	.01	.07	.00	.01	.01	.07

Table S13: Estimation results for $L = 1$ and $L_e = 1$ setup with $x_{i,t-1}$ as a regressor.

Designs		F1				F2				Fr				Fbic					
N	T	α	δ	Bias	RMSE	Std	Size	Bias	RMSE	Std	Size	Bias	RMSE	Std	Size	Bias	RMSE	Std	Size
β_1																			
200	4	4.0		.00	.04	.04	.06	.00	.13	.13	.02	.00	.04	.04	.06	.00	.05	.05	.07
200	4	4.3		.00	.05	.05	.07	.00	.15	.15	.02	.00	.05	.05	.07	.00	.07	.07	.07
200	4	8.0		.00	.02	.02	.07	.00	.05	.05	.04	.00	.02	.02	.07	.00	.02	.02	.07
200	4	8.3		.00	.02	.02	.06	.00	.07	.07	.02	.00	.02	.02	.06	.00	.04	.04	.07
200	8	4.0		-.01	.03	.03	.14	.00	.03	.03	.08	.00	.03	.03	.12	-.01	.03	.03	.14
200	8	4.3		-.01	.05	.05	.18	-.01	.05	.05	.08	-.01	.05	.05	.16	-.02	.06	.05	.18
200	8	8.0		.00	.01	.01	.12	.00	.02	.02	.07	.00	.01	.01	.12	.00	.02	.02	.12
200	8	8.3		.00	.02	.02	.13	.00	.02	.02	.07	.00	.02	.02	.11	-.01	.03	.03	.14
800	4	4.0		.00	.02	.02	.05	.00	.07	.07	.02	.00	.02	.02	.06	.00	.03	.03	.05
800	4	4.3		.00	.03	.03	.05	.00	.09	.09	.02	.00	.02	.02	.05	.00	.04	.04	.05
800	4	8.0		.00	.01	.01	.05	.00	.02	.02	.03	.00	.01	.01	.05	.00	.01	.01	.05
800	4	8.3		.00	.01	.01	.05	.00	.03	.03	.03	.00	.01	.01	.05	.00	.01	.01	.05
800	8	4.0		.00	.01	.01	.06	.00	.02	.02	.05	.00	.01	.01	.07	.00	.01	.01	.07
800	8	4.3		.00	.02	.02	.07	.00	.02	.02	.05	.00	.02	.02	.07	.00	.02	.02	.06
800	8	8.0		.00	.01	.01	.07	.00	.01	.01	.05	.00	.01	.01	.07	.00	.01	.01	.07
800	8	8.3		.00	.01	.01	.07	.00	.01	.01	.05	.00	.01	.01	.07	.00	.01	.01	.07

Table S14: Estimation results for $L = 2$ and $L_e = 2$ setup with $x_{i,t-1}$ as a regressor.

Designs				F1				F2				Fr				Fbic			
N	T	α	δ	Bias	RMSE	Std	Size												
α																			
200	4	4.0		-.05	.26	.25	.69	.00	.06	.06	.06	.00	.06	.06	.06	-.01	.23	.23	.11
200	4	4.3		-.07	.30	.29	.67	.00	.07	.07	.04	.00	.08	.08	.04	.00	.35	.35	.11
200	4	8.0		-.04	.19	.19	.67	.00	.05	.05	.05	.00	.05	.05	.06	-.01	.15	.15	.13
200	4	8.3		-.05	.21	.20	.67	.00	.05	.05	.04	.00	.06	.06	.05	-.03	.21	.20	.14
200	8	4.0		-.06	.21	.20	.83	.00	.02	.02	.09	.00	.02	.02	.09	.00	.03	.03	.10
200	8	4.3		-.10	.30	.29	.80	.00	.03	.03	.09	.00	.03	.03	.10	-.01	.06	.06	.11
200	8	8.0		-.03	.11	.11	.77	.00	.02	.02	.09	.00	.02	.02	.09	.00	.02	.02	.10
200	8	8.3		-.09	.20	.18	.80	.00	.02	.02	.09	.00	.02	.02	.09	.00	.02	.02	.10
800	4	4.0		-.05	.27	.26	.84	.00	.03	.03	.06	.00	.03	.03	.05	.00	.15	.15	.10
800	4	4.3		-.06	.29	.28	.82	.00	.04	.04	.05	.00	.04	.04	.05	.00	.26	.26	.09
800	4	8.0		-.04	.18	.18	.81	.00	.02	.02	.04	.00	.02	.02	.05	-.01	.17	.17	.09
800	4	8.3		-.06	.22	.21	.82	.00	.03	.03	.05	.00	.03	.03	.06	-.01	.16	.16	.10
800	8	4.0		-.06	.21	.20	.91	.00	.01	.01	.06	.00	.01	.01	.06	.00	.01	.01	.07
800	8	4.3		-.09	.28	.27	.90	.00	.01	.01	.05	.00	.01	.01	.05	.00	.01	.01	.06
800	8	8.0		-.03	.11	.10	.87	.00	.01	.01	.06	.00	.01	.01	.06	.00	.01	.01	.06
800	8	8.3		-.08	.19	.17	.89	.00	.01	.01	.06	.00	.01	.01	.06	.00	.01	.01	.08
β																			
200	4	4.0		-.01	.15	.15	.36	.01	.09	.09	.05	.01	.10	.10	.05	.03	.54	.54	.08
200	4	4.3		.02	.23	.23	.37	.00	.13	.13	.05	.00	.13	.13	.05	.03	.66	.66	.08
200	4	8.0		.00	.06	.06	.20	.00	.05	.05	.04	.00	.05	.05	.04	.00	.26	.26	.06
200	4	8.3		.00	.09	.09	.25	.00	.06	.06	.05	.00	.07	.07	.04	.02	.23	.23	.08
200	8	4.0		.00	.14	.14	.65	.00	.02	.02	.08	.00	.02	.02	.09	.00	.03	.03	.10
200	8	4.3		.05	.40	.40	.78	.01	.04	.04	.11	.01	.04	.04	.11	.02	.13	.13	.16
200	8	8.0		-.01	.06	.06	.50	.00	.02	.02	.10	.00	.02	.02	.10	.00	.02	.02	.10
200	8	8.3		.00	.19	.19	.69	.00	.02	.02	.10	.00	.02	.02	.10	.00	.04	.04	.12
800	4	4.0		.00	.15	.15	.56	.00	.06	.06	.05	.00	.06	.06	.05	.01	.61	.61	.08
800	4	4.3		.01	.20	.20	.59	.00	.07	.07	.04	.00	.07	.07	.04	-.02	.73	.73	.07
800	4	8.0		.00	.05	.05	.40	.00	.03	.03	.05	.00	.03	.03	.04	.00	.15	.15	.06
800	4	8.3		.00	.08	.08	.47	.00	.03	.03	.05	.00	.03	.03	.05	.00	.19	.19	.07
800	8	4.0		-.02	.10	.10	.73	.00	.01	.01	.05	.00	.01	.01	.05	.00	.01	.01	.05
800	8	4.3		.03	.31	.31	.83	.00	.02	.02	.06	.00	.02	.02	.06	.00	.02	.02	.07
800	8	8.0		-.01	.04	.04	.61	.00	.01	.01	.06	.00	.01	.01	.06	.00	.01	.01	.06
800	8	8.3		.00	.14	.14	.77	.00	.01	.01	.07	.00	.01	.01	.07	.00	.01	.01	.07

Table S15: Estimation results for $L = 2$ and $L_e = 2$ setup with $x_{i,t-1}$ as a regressor.

Designs		F1				F2				Fr				Fbic					
N	T	α	δ	Bias	RMSE	Std	Size	Bias	RMSE	Std	Size	Bias	RMSE	Std	Size	Bias	RMSE	Std	Size
β_1																			
200	4	4.0		.06	.29	.28	.61	.00	.09	.09	.05	.00	.09	.09	.05	.00	.42	.42	.09
200	4	4.3		.04	.38	.38	.61	.00	.12	.12	.05	.00	.11	.11	.05	-.02	.62	.62	.10
200	4	8.0		.02	.09	.08	.44	.00	.04	.04	.04	.00	.04	.04	.05	.00	.15	.15	.08
200	4	8.3		.02	.11	.11	.47	.00	.05	.05	.04	.00	.05	.05	.04	.00	.15	.15	.08
200	8	4.0		.06	.29	.28	.82	.00	.03	.03	.10	.00	.03	.03	.10	.00	.04	.04	.10
200	8	4.3		.06	.43	.43	.80	-.01	.04	.04	.10	-.01	.04	.04	.11	-.01	.09	.09	.13
200	8	8.0		.02	.09	.09	.68	.00	.02	.02	.09	.00	.02	.02	.10	.00	.02	.02	.11
200	8	8.3		.08	.26	.25	.78	.00	.02	.02	.10	.00	.02	.02	.10	.00	.03	.03	.12
800	4	4.0		.05	.29	.29	.79	.00	.05	.05	.05	.00	.05	.05	.05	.00	.26	.26	.08
800	4	4.3		.04	.37	.37	.81	.00	.06	.06	.05	.00	.06	.06	.05	.01	.47	.47	.08
800	4	8.0		.02	.08	.07	.68	.00	.02	.02	.05	.00	.02	.02	.05	.00	.08	.08	.07
800	4	8.3		.03	.11	.10	.68	.00	.03	.03	.05	.00	.03	.03	.06	.00	.13	.13	.08
800	8	4.0		.08	.27	.26	.87	.00	.01	.01	.05	.00	.01	.01	.06	.00	.01	.01	.06
800	8	4.3		.06	.40	.40	.90	.00	.02	.02	.07	.00	.02	.02	.06	.00	.02	.02	.06
800	8	8.0		.03	.08	.08	.80	.00	.01	.01	.06	.00	.01	.01	.06	.00	.01	.01	.06
800	8	8.3		.07	.22	.21	.88	.00	.01	.01	.07	.00	.01	.01	.06	.00	.01	.01	.06

Table S16: Estimation results for $L = 1$ and $L_e = 1$ setup for the nonlinear GMM ALS estimator.

Designs				α				β				
N	T	α	δ	Bias	RMSE	Std	Size	Bias	RMSE	Std	Size	J
200	4	.4	.0	-.03	.10	.09	.20	.01	.11	.11	.12	.13
200	4	.4	.3	-.07	.19	.18	.24	-.05	.25	.24	.18	.18
200	4	.8	.0	-.11	.29	.27	.26	-.02	.11	.11	.17	.19
200	4	.8	.3	-.19	.42	.37	.32	-.05	.18	.18	.25	.23
200	8	.4	.0	-.03	.09	.08	.40	.02	.07	.07	.31	.17
200	8	.4	.3	-.06	.15	.14	.43	.02	.13	.13	.37	.20
200	8	.8	.0	-.06	.17	.16	.41	.00	.05	.05	.28	.19
200	8	.8	.3	-.13	.27	.24	.48	.00	.10	.10	.33	.27
800	4	.4	.0	-.02	.08	.08	.14	.02	.15	.15	.11	.12
800	4	.4	.3	-.05	.15	.14	.19	-.04	.24	.24	.17	.16
800	4	.8	.0	-.09	.29	.28	.22	-.01	.16	.16	.18	.18
800	4	.8	.3	-.17	.41	.37	.27	-.03	.23	.23	.25	.23
800	8	.4	.0	-.01	.06	.06	.20	.01	.05	.05	.14	.15
800	8	.4	.3	-.05	.13	.13	.25	.01	.11	.11	.22	.20
800	8	.8	.0	-.04	.15	.15	.25	.00	.04	.04	.17	.19
800	8	.8	.3	-.11	.26	.24	.32	-.01	.10	.10	.26	.28

Table S17: Estimation results for $L = 1$ and $L_e = 1$ setup for the nonlinear GMM FIVU estimator.

Designs				α				β				
N	T	α	δ	Bias	RMSE	Std	Size	Bias	RMSE	Std	Size	J
200	4	.4	.0	.00	.05	.05	.10	.01	.09	.09	.09	.04
200	4	.4	.3	-.01	.10	.10	.11	.02	.12	.12	.10	.06
200	4	.8	.0	.00	.05	.05	.09	.00	.04	.04	.07	.04
200	4	.8	.3	-.01	.07	.07	.10	.00	.05	.05	.07	.05
200	8	.4	.0	.00	.02	.02	.18	.00	.03	.03	.15	.04
200	8	.4	.3	-.01	.07	.07	.21	.01	.08	.08	.19	.06
200	8	.8	.0	.00	.02	.02	.19	.00	.01	.01	.15	.04
200	8	.8	.3	.00	.05	.05	.21	.00	.05	.05	.16	.04
800	4	.4	.0	.00	.04	.04	.07	.01	.10	.10	.07	.07
800	4	.4	.3	-.01	.08	.08	.08	.01	.09	.09	.07	.06
800	4	.8	.0	.00	.05	.05	.07	.00	.04	.04	.07	.06
800	4	.8	.3	.00	.05	.05	.08	.00	.05	.05	.07	.06
800	8	.4	.0	.00	.02	.02	.09	.00	.03	.03	.08	.06
800	8	.4	.3	-.01	.06	.06	.09	.01	.07	.07	.08	.07
800	8	.8	.0	.00	.01	.01	.09	.00	.00	.00	.08	.05
800	8	.8	.3	.00	.04	.04	.10	.00	.04	.04	.08	.07

S4. Proofs

S4.1. Auxiliary Results

The following lemma is a restricted version of Lemma 1 in Andrews (2005).

Lemma 1. *Let $h(\cdot)$ be a scalar function such that $\mathbb{E}[h(\mathbf{H}_i)^2] < \infty$, where matrix \mathbf{H}_i satisfies conditional independence as in Assumption 2.1. Then:*

$$\frac{1}{N} \sum_{i=1}^N h(\mathbf{H}_i) \xrightarrow{p} \mathbb{E}_{\mathcal{F}}[h(\mathbf{H}_i)]. \quad (\text{S24})$$

Proof. First of all, observe that the existence of $\mathbb{E}_{\mathcal{F}}[h(\mathbf{H}_i)]$ is guaranteed by $\mathbb{E}[h(\mathbf{H}_i)^2] < \infty$. Next let $\tilde{h}_i \equiv h(\mathbf{H}_i) - \mathbb{E}_{\mathcal{F}}[h(\mathbf{H}_i)]$. By construction:

$$\mathbb{E}[\tilde{h}_i] = 0; \quad \mathbb{E}[(\tilde{h}_i)^2] < \infty. \quad (\text{S25})$$

Furthermore,

$$\begin{aligned} \mathbb{E} \left[\left(\frac{1}{N} \sum_{i=1}^N \tilde{h}_i \right)^2 \right] &= \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N \mathbb{E}[\tilde{h}_i \tilde{h}_j] \\ &= \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N \mathbb{E}[\mathbb{E}_{\mathcal{F}}[\tilde{h}_i \tilde{h}_j]] \\ &= \frac{1}{N^2} \sum_{i=1}^N \mathbb{E}[\mathbb{E}_{\mathcal{F}}[\tilde{h}_i \tilde{h}_i]] \\ &= \frac{1}{N} \mathbb{E}[(\tilde{h}_i)^2]. \end{aligned}$$

The third equality is obtained using the conditional independence assumption. Since the final expression is finite by assumption, the proof is concluded using the Chebyshev's inequality. \square

Lemma 2 (Central Limit Theorem).

Let $h(\cdot)$ be a scalar function such that $\mathbb{E}[|h(\mathbf{H}_i)|^{2+\delta}] < \infty$, where $\delta > 0$ and \mathbf{H}_i satisfies conditional independence as in Assumption 2.1. Then:

$$\frac{1}{\sqrt{N}} \sum_{i=1}^N (h(\mathbf{H}_i) - \mathbb{E}_{\mathcal{F}}[h(\mathbf{H}_i)]) \xrightarrow{d} (\mathbb{E}_{\mathcal{F}}[(h(\mathbf{H}_i) - \mathbb{E}_{\mathcal{F}}[h(\mathbf{H}_i)])^2])^{1/2} \times \pi \quad (\text{stably}), \quad (\text{S26})$$

where $\pi \sim N(0, 1)$ is independent of $\mathbb{E}_{\mathcal{F}}[(h(\mathbf{H}_i) - \mathbb{E}_{\mathcal{F}}[h(\mathbf{H}_i)])^2]$.

Proof. To prove this result we verify all conditions in Theorem 3.2 and Corollary 3.1 in Hall and Heyde (1980). As previously we denote $\tilde{h}_i = h(\mathbf{H}_i) - \mathbb{E}_{\mathcal{F}}[h(\mathbf{H}_i)]$. For $i \geq 1$, let \mathcal{A}_i denote the σ -field generated by \mathcal{F} and $(\tilde{h}_1, \dots, \tilde{h}_i)$. In this case $\mathbb{E}[\tilde{h}_i | \mathcal{A}_{i-1}] = \mathbb{E}[\tilde{h}_i] = 0$, and $\mathcal{A}_i \subseteq \mathcal{A}_{i+1}$ by construction. Thus $\{\tilde{h}_i, \mathcal{A}_i : i \geq 1\}$ is a Martingale Difference Sequence (MDS).

It remains to check the sufficient conditions in Corollary 3.1 of Hall and Heyde (1980). The main requirement is that the conditional Lindeberg's condition:

$$N^{-1} \sum_{i=1}^N \mathbb{E}[(\tilde{h}(\mathbf{H}_i))^2 \mathbf{1}(|\tilde{h}(\mathbf{H}_i)| > \sqrt{N}\epsilon) | \mathcal{A}_i] = \mathbb{E}_{\mathcal{F}}[(\tilde{h}(\mathbf{H}_i))^2 \mathbf{1}(|\tilde{h}(\mathbf{H}_i)| > \sqrt{N}\epsilon)] \xrightarrow{p} 0, \quad (\text{S27})$$

as $N \rightarrow \infty$, is satisfied. Given that the conditional Lyapunov's condition implies the conditional Lindeberg's condition, it is sufficient that $\mathbb{E}_{\mathcal{F}}[|h(\mathbf{H}_i)|^{2+\delta}] < \infty$, which is satisfied by assumption. Finally, consider the conditional variance:

$$V_N^2 = N^{-1} \sum_{i=1}^N \mathbb{E}[(\tilde{h}(\mathbf{H}_i))^2 | \mathcal{A}_i] = N^{-1} \sum_{i=1}^N \mathbb{E}_{\mathcal{F}}[(\tilde{h}(\mathbf{H}_i))^2] = \mathbb{E}_{\mathcal{F}}[(\tilde{h}(\mathbf{H}_i))^2].$$

This completes the proof. \square

5.2. Main Theorems

Proof of Theorem 1.

The GMM estimator is given by

$$\hat{\boldsymbol{\theta}} = \left(\hat{\boldsymbol{\Gamma}}' \boldsymbol{\Omega}_N \hat{\boldsymbol{\Gamma}} \right)^{-1} \hat{\boldsymbol{\Gamma}}' \boldsymbol{\Omega}_N \bar{\boldsymbol{\mu}}_N(\mathbf{0}), \quad (\text{S28})$$

which can be expanded in the usual way:

$$\sqrt{N} \left(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0 \right) = \left(\hat{\boldsymbol{\Gamma}}' \boldsymbol{\Omega}_N \hat{\boldsymbol{\Gamma}} \right)^{-1} \hat{\boldsymbol{\Gamma}}' \boldsymbol{\Omega}_N \sqrt{N} \bar{\boldsymbol{\mu}}_N(\boldsymbol{\theta}_0). \quad (\text{S29})$$

By assumption $\boldsymbol{\Omega}_N \xrightarrow{p} \boldsymbol{\Omega}$, where $\boldsymbol{\Omega}$ is a full-rank \mathcal{F} -measurable matrix. The sample Jacobian matrix has a well-defined asymptotic limit

$$\hat{\boldsymbol{\Gamma}} \xrightarrow{p} \boldsymbol{\Gamma}, \quad (\text{S30})$$

with $\boldsymbol{\Gamma}$ defined in Assumption 3.2. Next we establish convergence in distribution of $\sqrt{N} \bar{\boldsymbol{\mu}}_N(\boldsymbol{\theta}_0)$. At first we expand

$$\sqrt{N} \bar{\boldsymbol{\mu}}_N(\boldsymbol{\theta}_0) = \frac{1}{\sqrt{N}} \sum_{i=1}^N \boldsymbol{\mu}_i(\boldsymbol{\theta}_0) + o_P(1) = \frac{1}{\sqrt{N}} \sum_{i=1}^N \mathbf{S} \text{vec}(\boldsymbol{\Xi}_i) + o_P(1), \quad (\text{S31})$$

where

$$\begin{aligned}\boldsymbol{\Xi}_i &= \left(\mathbf{z}_i \boldsymbol{\varepsilon}'_i + (\mathbf{z}_i \boldsymbol{\lambda}'_i - \mathbb{E}_{\mathcal{F}}[\mathbf{z}_i \boldsymbol{\lambda}'_i]) \mathbf{F}' - \mathbf{G}_{z, \lambda_e} (\mathbf{A}_N^{-1})' \boldsymbol{\Psi}'_i \right) \\ &= \left(\mathbf{z}_i \boldsymbol{\varepsilon}'_i + (\mathbf{z}_i \boldsymbol{\lambda}'_i - \mathbb{E}_{\mathcal{F}}[\mathbf{z}_i \boldsymbol{\lambda}'_i]) \mathbf{F}' - \mathbf{G}_{z, \lambda_e} (\mathbf{A}_N^{-1})' \mathbf{B}_N (\mathbf{I}_T \otimes \boldsymbol{\psi}_i) \right).\end{aligned}\quad (\text{S32})$$

From the above decomposition we can express $\boldsymbol{\mu}_i(\boldsymbol{\theta}_0)$ as follows:

$$\boldsymbol{\mu}_i(\boldsymbol{\theta}_0) = \mathbf{R}_N \boldsymbol{\xi}_i, \quad (\text{S33})$$

where

$$\boldsymbol{\xi}_i = \begin{pmatrix} \mathbf{S} \text{vec}(\mathbf{z}_i \boldsymbol{\varepsilon}'_i) \\ \mathbf{S} \text{vec}((\mathbf{z}_i \boldsymbol{\lambda}'_i - \mathbb{E}_{\mathcal{F}}[\mathbf{z}_i \boldsymbol{\lambda}'_i]) \mathbf{F}') \\ \boldsymbol{\psi}_i \end{pmatrix}, \quad (\text{S34})$$

is a $[2\zeta + q \times 1]$ vector, while \mathbf{R}_N is a $[\zeta \times 2\zeta + q]$ selection matrix defined implicitly as a function of $\mathbf{G}_{z, \lambda_e} (\mathbf{A}_N^{-1})' \mathbf{B}_N$. In particular, by Assumption 3.1 $\mathbf{R}_N \xrightarrow{p} \mathbf{R}$, where the limiting matrix is \mathcal{F} -measurable. Notice that all elements in $\boldsymbol{\xi}_i$ satisfy the regularity conditions of Lemma 1, thus

$$\frac{1}{N} \sum_{i=1}^N \boldsymbol{\xi}_i \boldsymbol{\xi}'_i \xrightarrow{p} \mathbb{E}_{\mathcal{F}}[\boldsymbol{\xi}_i \boldsymbol{\xi}'_i]. \quad (\text{S35})$$

In addition Assumption 3.2 ensures that:

$$\frac{1}{N} \sum_{i=1}^N \boldsymbol{\mu}_i(\boldsymbol{\theta}_0) \boldsymbol{\mu}_i(\boldsymbol{\theta}_0)' = \mathbf{R}_N \left(\frac{1}{N} \sum_{i=1}^N \boldsymbol{\xi}_i \boldsymbol{\xi}'_i \right) \mathbf{R}'_N \xrightarrow{p} \boldsymbol{\Delta}, \quad (\text{S36})$$

has full column rank ζ .

Furthermore, by Assumptions 2.1 and 3.1 $\mathbb{E}_{\mathcal{F}}[\boldsymbol{\xi}_i] = \mathbf{0}_{2\zeta+q}$, and every element in $\boldsymbol{\xi}_i$ has a finite $4 + \delta$ absolute moment. Thus for all $h = 1, \dots, 2\zeta + q$ element-wise $\xi_i^{(h)}$ satisfies the conditions of Lemma 2, so that

$$\frac{1}{\sqrt{N}} \sum_{i=1}^N \xi_i^{(h)} \xrightarrow{d} (\mathbb{E}_{\mathcal{F}}[(\xi_i^{(h)})^2])^{1/2} \times \pi_h \quad (\text{stably}). \quad (\text{S37})$$

Convergence holds for all $h = 1, \dots, 2\zeta + q$ elements of $\boldsymbol{\xi}_i$. Using the Cramér-Wold device we establish joint convergence for the vector:

$$\sqrt{N} \bar{\boldsymbol{\mu}}_N(\boldsymbol{\theta}_0) \xrightarrow{d} \boldsymbol{\Delta}^{1/2} \times \boldsymbol{\pi} \quad (\text{stably}), \quad (\text{S38})$$

where $\boldsymbol{\pi} \sim N(\mathbf{0}_{\zeta}, \boldsymbol{\Gamma}_{\zeta})$. Given that $\boldsymbol{\Omega}$, $\boldsymbol{\Gamma}$, $\boldsymbol{\Delta}$ are \mathcal{F} measurable and independent of $\boldsymbol{\pi}$, we can invoke the Continuous Mapping Theorem (CMT) to conclude that:

$$\sqrt{N} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \xrightarrow{d} [(\boldsymbol{\Gamma}' \boldsymbol{\Omega} \boldsymbol{\Gamma})^{-1} \boldsymbol{\Gamma}' \boldsymbol{\Omega}] \boldsymbol{\Delta}^{1/2} \boldsymbol{\pi} \quad (\text{stably}). \quad (\text{S39})$$

The final result follows because stable convergence implies \mathcal{F} -stable convergence. \square

Proof of Theorem 2.

We prove the main result using the equivalent steps in Theorem 1 of Bai (2003). In particular, let \mathbf{V}_N denote the $[L_e \times L_e]$ diagonal matrix with the eigenvalues of $[T \times T]$ matrix $\widehat{\mathbf{X}} = T^{-1}\widehat{\mathbf{F}}_R\widehat{\mathbf{F}}_R'$ on the diagonal in decreasing order. By the definition of eigenvalues and eigenvectors $T^{-1}\widehat{\mathbf{X}}\tilde{\mathbf{F}} = \tilde{\mathbf{F}}\mathbf{V}_N$. Let $\mathbf{A}_N \equiv (\mathbf{A}_R\mathbf{A}'_R) \left(\frac{\tilde{\mathbf{F}}'\tilde{\mathbf{F}}}{T} \right) \mathbf{V}_N^{-1}$, and $\|\cdot\|$ denote the Frobenius norm. In what follows we ignore the (R) subscript on all $\overline{\boldsymbol{\psi}}_s^{(R)}$, and simply denote these elements by $\overline{\boldsymbol{\psi}}_s$. We break down the proof into four steps.

Step 1. By Lemma 1 $\widehat{\mathbf{F}}_R \xrightarrow{p} \mathbf{F}_e\mathbf{A}_R$, thus $\widehat{\mathbf{X}} \xrightarrow{p} \mathbf{X} = \mathbf{F}_e\mathbf{A}_R\mathbf{A}'_R\mathbf{F}'_e$. Given that $\widehat{\mathbf{X}}$ is a matrix of fixed dimensions, by the continuity of eigenvalues we also have $\mathbf{V}_N \xrightarrow{p} \mathbf{V}$. Here \mathbf{V} denotes the corresponding diagonal matrix with the L_e largest eigenvalues of $\mathbf{F}_e\mathbf{A}_R\mathbf{A}'_R\mathbf{F}'_e$ in decreasing order.

Step 2. Consider the following identity:

$$\begin{aligned} \tilde{\mathbf{f}}_t - \mathbf{A}'_N\mathbf{f}_{t,e} &= \mathbf{V}_N^{-1} \frac{1}{T} \left(\sum_{s=1}^T \tilde{\mathbf{f}}_s(\overline{\boldsymbol{\psi}}_s'\mathbf{A}'_R\mathbf{A}_R\overline{\boldsymbol{\psi}}_t) + \sum_{s=1}^T \tilde{\mathbf{f}}_s\mathbf{f}'_{s,e}\mathbf{A}_R\overline{\boldsymbol{\psi}}_t + \sum_{s=1}^T \tilde{\mathbf{f}}_s\mathbf{f}'_{t,e}\mathbf{A}_R\overline{\boldsymbol{\psi}}_s \right) \\ &= \mathbf{V}_N^{-1} \frac{1}{T} (\mathbf{A}_{1t} + \mathbf{A}_{2t} + \mathbf{A}_{3t}), \end{aligned} \quad (\text{S40})$$

for all $t = 1, \dots, T$. From the above equation, we have

$$\|(\tilde{\mathbf{f}}_t - \mathbf{A}'_N\mathbf{f}_{t,e})\| \leq \frac{1}{T} \|\mathbf{V}_N^{-1}\| (\|\mathbf{A}_{1t}\| + \|\mathbf{A}_{2t}\| + \|\mathbf{A}_{3t}\|). \quad (\text{S41})$$

Consider each term separately:

$$\begin{aligned} \|\mathbf{A}_{1t}\| &= \left\| \sum_{s=1}^T \tilde{\mathbf{f}}_s(\overline{\boldsymbol{\psi}}_s'\mathbf{A}'_R\mathbf{A}_R\overline{\boldsymbol{\psi}}_t) \right\| \\ &\leq \|\mathbf{A}_R\overline{\boldsymbol{\psi}}_t\| \left\| \sum_{s=1}^T \tilde{\mathbf{f}}_s\overline{\boldsymbol{\psi}}_s'\mathbf{A}'_R \right\| \\ &\leq \|\mathbf{A}_R\overline{\boldsymbol{\psi}}_t\| \sqrt{\sum_{s=1}^T \|\tilde{\mathbf{f}}_s\|^2} \sqrt{\sum_{s=1}^T \|(\mathbf{A}_R\overline{\boldsymbol{\psi}}_s)\|^2} \\ &= O_P(N^{-1/2})O_P(1)O_P(N^{-1/2}). \end{aligned}$$

Here the $O_P(N^{-1/2})$ results follow from the fact that for all $s = 1, \dots, T$, $\overline{\boldsymbol{\psi}}_s$ are zero-mean vectors that satisfy the assumptions of Lemma 2. On the other hand, $\sum_{s=1}^T \|\tilde{\mathbf{f}}_s\|^2 =$

T by construction. Similarly, for \mathbf{A}_{2t} we have

$$\begin{aligned}
\|\mathbf{A}_{2t}\| &= \left\| \sum_{s=1}^T \tilde{\mathbf{f}}_s \mathbf{f}'_{s,e} \mathbf{A}_R \bar{\boldsymbol{\psi}}_t \right\| \\
&\leq \|\mathbf{A}_R \bar{\boldsymbol{\psi}}_t\| \left\| \sum_{s=1}^T \tilde{\mathbf{f}}_s \mathbf{f}'_{s,e} \right\| \\
&\leq \|\mathbf{A}_R \bar{\boldsymbol{\psi}}_t\| \sqrt{\sum_{s=1}^T \|\tilde{\mathbf{f}}_s\|^2} \sqrt{\sum_{s=1}^T \|\mathbf{f}_{s,e}\|^2} \\
&= O_P(N^{-1/2}) O_P(1) O_P(1).
\end{aligned}$$

For \mathbf{A}_{3t} we obtain

$$\begin{aligned}
\|\mathbf{A}_{3t}\| &= \left\| \sum_{s=1}^T \tilde{\mathbf{f}}_s \mathbf{f}'_{t,e} \mathbf{A}_R \bar{\boldsymbol{\psi}}_s \right\| \\
&\leq \|\mathbf{f}_t\| \left\| \sum_{s=1}^T \tilde{\mathbf{f}}_s \bar{\boldsymbol{\psi}}'_s \mathbf{A}'_R \right\| \\
&\leq \|\mathbf{f}_{t,e}\| \sqrt{\sum_{s=1}^T \|\tilde{\mathbf{f}}_s\|^2} \sqrt{\sum_{s=1}^T \|(\mathbf{A}_R \bar{\boldsymbol{\psi}}_s)\|^2} \\
&= O_P(1) O_P(1) O_P(N^{-1/2}).
\end{aligned}$$

Next observe that $\|\mathbf{V}_N\| = O_P(1)$, and $\text{rk}(\mathbf{V}_N) = \text{rk}(\mathbf{V}) = L_e$ implies $\|\mathbf{V}_N^{-1}\| = O_P(1)$; see e.g. Andrews (1987). Combining these results, it follows that $\|(\tilde{\mathbf{f}}_t - \mathbf{A}'_N \mathbf{f}_t)\| = O_P(N^{-1/2}) = o_P(1)$.

Step 3. Under Assumption 3.3, following analogous steps to those in Proposition 1 of Bai (2003) and using rates from Step 2, one can establish that \mathbf{A}_N has a well-defined asymptotic limit, i.e. $\mathbf{A}_N \xrightarrow{p} \mathbf{A}$. Here \mathbf{A} is a function of $(\mathbf{A}_R \mathbf{A}'_R)$ and $(\mathbf{F}'_e \mathbf{F}_e)$, see Proposition 1 of Bai (2003). Hence, \mathbf{A} is \mathcal{F} -measurable.

Step 4. Combining Steps 2 and 3 it follows that

$$\tilde{\mathbf{f}}_t \xrightarrow{p} \mathbf{A}' \mathbf{f}_{t,e}, \quad t = 1, \dots, T. \quad (\text{S42})$$

From Step 1 we know that \mathbf{A}_{1t} is of lower order, so that

$$\sqrt{N}(\tilde{\mathbf{f}}_t - \mathbf{A}'_N \mathbf{f}_{t,e}) = \mathbf{V}_N^{-1} \frac{\sqrt{N}}{T} \left(\sum_{s=1}^T \tilde{\mathbf{f}}_s \mathbf{f}'_{s,e} \mathbf{A}_R \bar{\boldsymbol{\psi}}_t + \sum_{s=1}^T \tilde{\mathbf{f}}_s \mathbf{f}'_{t,e} \mathbf{A}_R \bar{\boldsymbol{\psi}}_s \right) + O_P(N^{-1/2}). \quad (\text{S43})$$

Finally, observe that for any t , we can expand:

$$\sqrt{N}(\tilde{\mathbf{f}}_t - \mathbf{A}'_N \mathbf{f}_{t,e}) = \mathbf{B}_N^{(t)} \frac{1}{\sqrt{N}} \bar{\boldsymbol{\psi}}, \quad (\text{S44})$$

where

$$\mathbf{B}_N^{(t)} = \frac{1}{T} \mathbf{V}_N^{-1} \left(\tilde{\mathbf{f}}_1 \mathbf{f}'_{1,e} \mathbf{A}_R, \dots, (\tilde{\mathbf{f}}_t \mathbf{f}'_{t,e} \mathbf{A}_R + \sum_{s=1}^T \tilde{\mathbf{f}}_s \mathbf{f}'_{s,e} \mathbf{A}_R), \dots, \tilde{\mathbf{f}}_T \mathbf{f}'_{T,e} \mathbf{A}_R \right). \quad (\text{S45})$$

Using this notation we can see that

$$\mathbf{B}_N = \left(\mathbf{B}_N^{(1)}, \dots, \mathbf{B}_N^{(t)}, \dots, \mathbf{B}_N^{(T)} \right). \quad (\text{S46})$$

Here every block element is an $[L_e \times TL_e]$ matrix. Finally, from Eq. (S42) it follows that \mathbf{B}_N has a well-defined \mathcal{F} -measurable limit \mathbf{B} . This completes the proof. \square

Proof of Theorem 3.

As in Theorem 1, we express

$$\hat{\boldsymbol{\theta}} = \left(\hat{\boldsymbol{\Gamma}}' \boldsymbol{\Omega}_N \hat{\boldsymbol{\Gamma}} \right)^{-1} \hat{\boldsymbol{\Gamma}}' \boldsymbol{\Omega}_N \bar{\boldsymbol{\mu}}_N(\mathbf{0}). \quad (\text{S47})$$

Evaluating the expression above at $\boldsymbol{\Omega}_N = \mathbf{I}_\zeta$, and using properties of partitioned matrices, one obtains

$$\hat{\boldsymbol{\beta}} = \left(\hat{\boldsymbol{\Gamma}}'_\beta \hat{\boldsymbol{\Gamma}}_\beta - \hat{\boldsymbol{\Gamma}}'_\beta \hat{\boldsymbol{\Gamma}}_g \left(\hat{\boldsymbol{\Gamma}}'_g \hat{\boldsymbol{\Gamma}}_g \right)^{-1} \hat{\boldsymbol{\Gamma}}'_g \hat{\boldsymbol{\Gamma}}_\beta \right)^{-1} \left(\hat{\boldsymbol{\Gamma}}'_\beta - \hat{\boldsymbol{\Gamma}}'_\beta \hat{\boldsymbol{\Gamma}}_g \left(\hat{\boldsymbol{\Gamma}}'_g \hat{\boldsymbol{\Gamma}}_g \right)^{-1} \hat{\boldsymbol{\Gamma}}'_g \right) \bar{\boldsymbol{\mu}}_N(\mathbf{0}), \quad (\text{S48})$$

Analogously to Theorem 1, we maintain the assumption that

$$\hat{\boldsymbol{\Gamma}}_\beta \xrightarrow{p} \boldsymbol{\Gamma}_\beta. \quad (\text{S49})$$

Next, let $\dot{\boldsymbol{\theta}} = (\boldsymbol{\beta}'_0, \mathbf{0}')'$. Using this definition, Eq. (S48) simplifies to:

$$\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0 = \left(\hat{\boldsymbol{\Gamma}}'_\beta \hat{\boldsymbol{\Gamma}}_\beta - \hat{\boldsymbol{\Gamma}}'_\beta \hat{\boldsymbol{\Gamma}}_g \left(\hat{\boldsymbol{\Gamma}}'_g \hat{\boldsymbol{\Gamma}}_g \right)^{-1} \hat{\boldsymbol{\Gamma}}'_g \hat{\boldsymbol{\Gamma}}_\beta \right)^{-1} \left(\hat{\boldsymbol{\Gamma}}'_\beta - \hat{\boldsymbol{\Gamma}}'_\beta \hat{\boldsymbol{\Gamma}}_g \left(\hat{\boldsymbol{\Gamma}}'_g \hat{\boldsymbol{\Gamma}}_g \right)^{-1} \hat{\boldsymbol{\Gamma}}'_g \right) \bar{\boldsymbol{\mu}}_N(\dot{\boldsymbol{\theta}}). \quad (\text{S50})$$

The last component can be expanded as follows:

$$\bar{\boldsymbol{\mu}}_N(\dot{\boldsymbol{\theta}}) = \frac{1}{N} \sum_{i=1}^N \boldsymbol{\mu}_i(\dot{\boldsymbol{\theta}}) = \frac{1}{N} \sum_{i=1}^N \mathbf{S} \text{vec}(\mathbf{z}_i \boldsymbol{\varepsilon}'_i + (\mathbf{z}_i \boldsymbol{\lambda}'_i) \mathbf{F}'). \quad (\text{S51})$$

Thus, it is easy to see that

$$\mathbb{E}_{\mathcal{F}}[\boldsymbol{\mu}_i(\dot{\boldsymbol{\theta}})] = \mathbf{S} \text{vec}(\mathbb{E}_{\mathcal{F}}[\mathbf{z}_i \boldsymbol{\lambda}'_i] \mathbf{F}') = \mathbf{S} \text{vec}(\mathbf{G}_{z,\lambda} \mathbf{F}'). \quad (\text{S52})$$

Furthermore, all elements of $\boldsymbol{\mu}_i(\dot{\boldsymbol{\theta}})$ satisfy the sufficient conditions of Lemma 1. As a result, using the CMT we have

$$\bar{\boldsymbol{\mu}}_N(\dot{\boldsymbol{\theta}}) \xrightarrow{p} \mathbf{S} \text{vec}(\mathbf{G}_{z,\lambda} \mathbf{F}'). \quad (\text{S53})$$

It remains to investigate the properties of the term $\widehat{\boldsymbol{\Gamma}}_g$, where

$$\widehat{\boldsymbol{\Gamma}}_g = \mathbf{S} \left(\left(\frac{1}{N} \sum_{i=1}^N \mathbf{v}_i \mathbf{w}'_i \right) \otimes \mathbf{I}_d \right). \quad (\text{S54})$$

By construction,

$$\mathbf{v}_i \mathbf{w}'_i = \mathbf{F}_e \boldsymbol{\gamma}_i \mathbf{w}'_i + \mathbf{u}_i \mathbf{w}'_i = \mathbf{F}_e \mathbf{C} \mathbf{D}' + (\mathbf{F}_e (\boldsymbol{\gamma}_i \mathbf{w}'_i - \mathbf{C} \mathbf{D}') + \mathbf{u}_i \mathbf{w}'_i),$$

from $\mathbb{E}_{\mathcal{F}}[\boldsymbol{\gamma}_i \mathbf{w}'_i] = \mathbf{C} \mathbf{D}'$. Furthermore, as Eq. (S48) is invariant to non-singular multiplication of $\widehat{\boldsymbol{\Gamma}}_g$, we consider the following transformed matrix:

$$\widehat{\boldsymbol{\Xi}} = \widehat{\boldsymbol{\Gamma}}_g \left((\mathbf{D} \quad \sqrt{N} \mathbf{D}_{\perp}) \otimes \mathbf{I}_d \right) = (\widehat{\boldsymbol{\Xi}}_q \quad \widehat{\boldsymbol{\Xi}}_{L_e-q}). \quad (\text{S55})$$

From this decomposition it follows that

$$\widehat{\boldsymbol{\Xi}}_q = \mathbf{S} \left(\left(\frac{1}{N} \sum_{i=1}^N \mathbf{v}_i \mathbf{w}'_i \mathbf{D} \right) \otimes \mathbf{I}_d \right) \xrightarrow{p} \mathbf{S} (\mathbf{F}_e \mathbf{C} (\mathbf{D}' \mathbf{D}) \otimes \mathbf{I}_d), \quad (\text{S56})$$

which has full column-rank a.s. by construction. On the other hand, the remaining component is given by

$$\widehat{\boldsymbol{\Xi}}_{L_e-q} = \mathbf{S} \left(\left(\frac{1}{\sqrt{N}} \sum_{i=1}^N \mathbf{v}_i \mathbf{w}'_i \mathbf{D}_{\perp} \right) \otimes \mathbf{I}_d \right), \quad (\text{S57})$$

which is a mean zero random matrix, such that

$$\mathbb{E}_{\mathcal{F}}[\mathbf{v}_i \mathbf{w}'_i \mathbf{D}_{\perp}] = \mathbf{O}_{T \times L_e-q}. \quad (\text{S58})$$

Furthermore, for each i all elements of $\mathbf{v}_i \mathbf{w}'_i \mathbf{D}_{\perp}$ satisfy the sufficient conditions of Lemma 2. Thus, using the Cramér-Wold device, we establish the joint convergence, i.e.

$$\frac{1}{\sqrt{N}} \sum_{i=1}^N \text{vec}(\mathbf{v}_i \mathbf{w}'_i \mathbf{D}_{\perp}) \xrightarrow{d} (\mathbb{E}_{\mathcal{F}}[\text{vec}(\mathbf{v}_i \mathbf{w}'_i \mathbf{D}_{\perp}) \text{vec}(\mathbf{v}_i \mathbf{w}'_i \mathbf{D}_{\perp})'])^{1/2} \times \boldsymbol{\pi}_F \quad (\text{stably}). \quad (\text{S59})$$

The final result of this theorem follows from the CMT. \square

S4.3. Additional Results

Proof of Theorem S1.

As in Theorem 1, we express

$$\left(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0\right) = \left(\widehat{\boldsymbol{\Gamma}}' \boldsymbol{\Omega}_N \widehat{\boldsymbol{\Gamma}}\right)^{-1} \widehat{\boldsymbol{\Gamma}}' \boldsymbol{\Omega}_N \bar{\boldsymbol{\mu}}_N(\boldsymbol{\theta}_0). \quad (\text{S60})$$

By assumption $\boldsymbol{\Omega}_N \xrightarrow{p} \boldsymbol{\Omega}$, where $\boldsymbol{\Omega}$ is a full rank \mathcal{F} -measurable matrix. The sample Jacobian matrix has a well-defined asymptotic limit

$$\widehat{\boldsymbol{\Gamma}} \xrightarrow{p} \boldsymbol{\Gamma}, \quad (\text{S61})$$

where $\boldsymbol{\Gamma}$ is defined in Assumption 3.2. Next we show that

$$\bar{\boldsymbol{\mu}}_N(\boldsymbol{\theta}_0) \xrightarrow{p} \mathbf{0}_\zeta. \quad (\text{S62})$$

At first, we expand

$$\bar{\boldsymbol{\mu}}_N(\boldsymbol{\theta}_0) = \frac{1}{N} \sum_{i=1}^N \boldsymbol{\mu}_i(\boldsymbol{\theta}_0) + o_P(1) = \frac{1}{N} \sum_{i=1}^N \mathbf{S} \text{vec}(\boldsymbol{\Xi}_i) + o_P(1), \quad (\text{S63})$$

where in this setup

$$\boldsymbol{\Xi}_i = \left(\mathbf{z}_i \boldsymbol{\varepsilon}'_i + (\mathbf{z}_i \boldsymbol{\lambda}'_i - \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{\mathcal{F}}[\mathbf{z}_i \boldsymbol{\lambda}'_i]) \mathbf{F}' - \frac{1}{N} \sum_{i=1}^N \mathbf{G}_{z, \lambda_e}(i) (\mathbf{A}_N^{-1})' \mathbf{B}_N(\mathbf{I}_T \otimes \boldsymbol{\psi}_i) \right). \quad (\text{S64})$$

From the above decomposition one can express $\boldsymbol{\mu}_i(\boldsymbol{\theta}_0)$ as follows:

$$\boldsymbol{\mu}_i(\boldsymbol{\theta}_0) = \mathbf{R}_N \boldsymbol{\xi}_i, \quad (\text{S65})$$

where

$$\boldsymbol{\xi}_i = \begin{pmatrix} \mathbf{S} \text{vec}(\mathbf{z}_i \boldsymbol{\varepsilon}'_i) \\ \mathbf{S} \text{vec}((\mathbf{z}_i \boldsymbol{\lambda}'_i - \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{\mathcal{F}}[\mathbf{z}_i \boldsymbol{\lambda}'_i]) \mathbf{F}') \\ \boldsymbol{\psi}_i \end{pmatrix}. \quad (\text{S66})$$

Finally, provided that $\text{plim}_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \mathbf{G}_{z, \lambda}(i)$ exists, $\mathbf{R}_N \rightarrow \mathbf{R}$, while

$$\frac{1}{N} \sum_{i=1}^N \boldsymbol{\xi}_i \xrightarrow{p} \mathbf{0}_{2\zeta+q}, \quad (\text{S67})$$

after appropriately extending Lemma 1. By the CMT, it follows that

$$\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0 \xrightarrow{p} \mathbf{0}. \quad (\text{S68})$$

□

Proof of Proposition S1.

Provided that $E[\gamma_i] \neq 0$ (i.e. Assumption 3.1 (c) is satisfied), then based on the results in Westerlund et al. (2019) we conclude that

$$\text{Avar}(\widehat{\beta}_{CCE}) = \frac{\sigma_\varepsilon^2}{\text{tr}(E[\mathbf{u}_i \mathbf{u}_i'] \mathbf{M}_f)}. \quad (\text{S69})$$

Here $\mathbf{M}_f = \mathbf{I}_T - \mathbf{f}(\mathbf{f}'\mathbf{f})^{-1}\mathbf{f}'$, is the orthogonal projection matrix of the space spanned by \mathbf{f} . If we further assume that the idiosyncratic component \mathbf{u}_i is uncorrelated over time, then $\boldsymbol{\Sigma}_u \equiv E[\mathbf{u}_i \mathbf{u}_i'] = \sigma_u^2 \mathbf{I}_T$, so that:

$$\text{Avar}(\widehat{\beta}_{CCE}) = \frac{\sigma_\varepsilon^2}{\sigma_u^2 \text{tr}(\mathbf{M}_f)} = \frac{1}{T-1} \frac{\sigma_\varepsilon^2}{\sigma_u^2}. \quad (\text{S70})$$

Here we have made use of the fact that \mathbf{M}_f is a orthogonal projection matrix, so that $\text{tr}(\mathbf{M}_f) = \text{rk}(\mathbf{M}_f) = T - 1$.

Next we consider our GMM estimator. It is straightforward to show that

$$\boldsymbol{\Gamma}_\beta = \text{vec}(\boldsymbol{\Sigma}_u + \mathbf{f} E[\gamma_i^2] \mathbf{f}') \quad (\text{S71})$$

$$\boldsymbol{\Gamma}_g = (\mathbf{f} \gamma \otimes \mathbf{I}_T). \quad (\text{S72})$$

Next, let $\boldsymbol{\Pi} = \boldsymbol{\Gamma}'_\beta - \boldsymbol{\Gamma}'_\beta \boldsymbol{\Gamma}_g (\boldsymbol{\Gamma}'_g \boldsymbol{\Gamma}_g)^{-1} \boldsymbol{\Gamma}'_g$. Hence, the asymptotic variance of the GMM estimator for β is of the following form:

$$\text{Avar}(\widehat{\beta}_{GMM}) = (\boldsymbol{\Pi} \boldsymbol{\Gamma}_\beta)^{-1} \boldsymbol{\Pi} \boldsymbol{\Delta} \boldsymbol{\Pi}' (\boldsymbol{\Pi} \boldsymbol{\Gamma}_\beta)^{-1}. \quad (\text{S73})$$

In our case, we have

$$\boldsymbol{\Pi} = \boldsymbol{\Gamma}'_\beta (\mathbf{M}_f \otimes \mathbf{I}_T) = \text{vec}(\boldsymbol{\Sigma}_u)' (\mathbf{M}_f \otimes \mathbf{I}_T), \quad (\text{S74})$$

where the second equality follows using the properties of the Kronecker product and the $\text{vec}(\cdot)$ operator. Likewise,

$$\begin{aligned} \boldsymbol{\Pi} \boldsymbol{\Gamma}_\beta &= \text{vec}(\boldsymbol{\Sigma}_u)' (\mathbf{M}_f \otimes \mathbf{I}_T) \text{vec}(\boldsymbol{\Sigma}_u) \\ &= \text{tr}(\boldsymbol{\Sigma}'_u \boldsymbol{\Sigma}_u \mathbf{M}_f) \\ &= \text{tr}(\boldsymbol{\Sigma}_u \boldsymbol{\Sigma}_u \mathbf{M}_f), \end{aligned}$$

by the symmetry of the corresponding matrices. Finally, consider the term $\boldsymbol{\Pi} \boldsymbol{\Delta} \boldsymbol{\Pi}'$. By construction,

$$\boldsymbol{\Delta} = E_{\mathcal{F}}[\boldsymbol{\mu}_i(\boldsymbol{\theta}_0) \boldsymbol{\mu}_i(\boldsymbol{\theta}_0)'] = E_{\mathcal{F}}[\text{vec}(\mathbf{Q}_i) \text{vec}(\mathbf{Q}_i)'], \quad (\text{S75})$$

where

$$\mathbf{Q}_i = \mathbf{x}_i \boldsymbol{\varepsilon}_i' + \mathbf{x}_i \lambda_i \mathbf{f}' - \mathbb{E}_{\mathcal{F}}[\mathbf{x}_i \lambda_i] \gamma^{-1} (\gamma_i \mathbf{f}' + \mathbf{u}_i'). \quad (\text{S76})$$

As we are only interested in $\boldsymbol{\Pi} \boldsymbol{\Delta} \boldsymbol{\Pi}'$ and not in $\boldsymbol{\Delta}$ itself, only the contributions in \mathbf{Q}_i that are not proportional to \mathbf{f} matter for the final result. That is,

$$\boldsymbol{\Pi} \boldsymbol{\Delta} \boldsymbol{\Pi}' = \text{vec}(\boldsymbol{\Sigma}_{\mathbf{u}})' (\mathbf{M}_{\mathbf{f}} \otimes \mathbf{I}_T) \boldsymbol{\Delta}_{\perp} (\mathbf{M}_{\mathbf{f}} \otimes \mathbf{I}_T) \text{vec}(\boldsymbol{\Sigma}_{\mathbf{u}}).$$

Using the fact that both $\boldsymbol{\varepsilon}_i$ and \mathbf{u}_i are independent and that $\mathbb{E}_{\mathcal{F}}[\mathbf{x}_i \lambda_i] = \mathbf{f} \mathbb{E}[\gamma_i \lambda_i]$, we conclude that

$$\boldsymbol{\Delta}_{\perp} = \mathbf{I}_T \otimes \sigma_{\varepsilon}^2 (\mathbf{f} \mathbb{E}[\gamma_i^2] \mathbf{f}' + \boldsymbol{\Sigma}_{\mathbf{u}}) + (\mathbf{I}_T \otimes \mathbf{f} \mathbb{E}[\gamma_i \lambda_i]) \boldsymbol{\Sigma}_{\mathbf{u}} (\mathbf{I}_T \otimes \mathbf{f}' \mathbb{E}[\gamma_i \lambda_i]).$$

If we further assume that $\boldsymbol{\Sigma}_{\mathbf{u}} = \sigma_u^2 \mathbf{I}_T$, then $\boldsymbol{\Pi} \boldsymbol{\Delta} \boldsymbol{\Pi}'$ becomes free from the terms that are additive in factors, such that

$$\boldsymbol{\Pi} \boldsymbol{\Delta} \boldsymbol{\Pi}' = (\sigma_u^2)^3 \text{tr}(\mathbf{M}_{\mathbf{f}}) = (\sigma_u^2)^3 (T - 1).$$

As in this case we also have $\boldsymbol{\Pi} \boldsymbol{\Gamma}_{\beta} = (\sigma_u^2)^2 (T - 1)$, we conclude that:

$$\text{Avar}(\widehat{\beta}_{GMM}) = \frac{1}{T-1} \frac{\sigma_{\varepsilon}^2}{\sigma_u^2}. \quad (\text{S77})$$

□

References

- AHN, S. C. AND A. R. HORENSTEIN (2013): “Eigenvalue Ratio Test for the Number of Factors,” *Econometrica*, 81, 1203–1227.
- AHN, S. C., Y. H. LEE, AND P. SCHMIDT (2013): “Panel Data Models with Multiple Time-varying Individual Effects,” *Journal of Econometrics*, 174, 1–14.
- ANDREWS, D. W. K. (1987): “Asymptotic Results for Generalized Wald Tests,” *Econometric Theory*, 3, pp. 348–358.
- (2005): “Cross-Section Regression with Common Shocks,” *Econometrica*, 73, 1551–1585.
- BAI, J. (2003): “Inferential Theory for Factor Models of Large Dimensions,” *Econometrica*, 71, 135–171.
- HALL, P. AND C. C. HEYDE (1980): *Martingale Limit Theory and Its Application*, Probability and Mathematical Statistics, Academic Press.

- HSIAO, C., M. H. PESARAN, AND A. K. TAHMISCIOLU (2002): “Maximum Likelihood Estimation of Fixed Effects Dynamic Panel Data Models Covering Short Time Periods,” *Journal of Econometrics*, 109, 107–150.
- JUODIS, A. AND V. SARAFIDIS (2018): “Fixed T Dynamic Panel Data Estimators with Multifactor Errors,” *Econometric Reviews*, 37, 893–929.
- PESARAN, M. H. (2006): “Estimation and Inference in Large Heterogeneous Panels with a Multifactor Error Structure,” *Econometrica*, 74, 967–1012.
- ROBERTSON, D. AND V. SARAFIDIS (2015): “IV Estimation of Panels with Factor Residuals,” *Journal of Econometrics*, 185, 526–541.
- WESTERLUND, J., Y. PETROVA, AND M. NORKUTE (2019): “CCE in Fixed-T Panels,” *Journal of Applied Econometrics*, 34, 746–761.