

Linear regression for panel with unknown number of factors as interactive fixed effects

Hyungsik Roger Moon
Martin Weidner

The Institute for Fiscal Studies
Department of Economics, UCL

cemmap working paper CWP35/14

Linear Regression for Panel with Unknown Number of Factors as Interactive Fixed Effects*

Hyungsik Roger Moon[†] Martin Weidner[§]

August 18, 2014

Abstract

In this paper we study the least squares (LS) estimator in a linear panel regression model with *unknown* number of factors appearing as interactive fixed effects. Assuming that the number of factors used in estimation is larger than the true number of factors in the data we establish the limiting distribution of the LS estimator for the regression coefficients, as the number of time periods and the number of cross-sectional units jointly go to infinity. The main result of the paper is that under certain assumptions the limiting distribution of the LS estimator is independent of the number of factors used in the estimation, as long as this number is not underestimated. The important practical implication of this result is that for inference on the regression coefficients one does not necessarily need to estimate the number of interactive fixed effects consistently.

Keywords: Panel data, interactive fixed effects, factor models, perturbation theory of linear operators, random matrix theory.

JEL-Classification: C23, C33

*We thank the participants of the 2009 Cowles Summer Conference “Handling Dependence: Temporal, Cross-sectional, and Spatial” at Yale University, of the 2012 North American Summer Meeting of the Econometric Society at Northwestern University, of the 18th International Conference on Panel Data at the Banque de France, of the 2013 North American Winter Meeting of the Econometric Society in San Diego, of the 2014 Asia Meeting of Econometric Society in Taipei, of the 2014 Econometric Study Group Conference in Bristol, and of the econometrics seminars in USC and Toulouse for many interesting comments, and we thank Dukpa Kim, Tatsushi Oka, and Alexei Onatski for helpful discussions. We are also grateful for the comments and suggestions of James Stock, Elie Tamer, and anonymous referees. Moon acknowledges financial supports of the NSF via SES 0920903 and the faculty grant award of USC. Weidner acknowledges support from the Economic and Social Research Council through the ESRC Centre for Microdata Methods and Practice grant RES-589-28-0002.

[†]Department of Economics, University of Southern California, Los Angeles, CA 90089-0253. Email: moonr@usc.edu. Department of Economics, Yonsei University, Seoul, Korea.

[§]Corresponding author. Department of Economics, University College London, Gower Street, London WC1E 6BT, U.K., and CeMMaP. Email: m.weidner@ucl.ac.uk.

1 Introduction

Panel data models typically incorporate individual and time effects to control for heterogeneity in cross-section and over time. While often these individual and time effects enter the model additively, they can also be interacted multiplicatively, thus giving rise to so called interactive effects, which we also refer to as a factor structure. The multiplicative form captures the heterogeneity in the data more flexibly, since it allows for common time-varying shocks (factors) to affect the cross-sectional units with individual specific sensitivities (factor loadings).¹ It is this flexibility that motivated the discussion of interactive effects in the econometrics literature, e.g. Holtz-Eakin, Newey and Rosen (1988), Ahn, Lee and Schmidt (2001; 2013), Pesaran (2006), Bai (2009a; 2013), Zaffaroni (2009), Moon and Weidner (2013), and Lu and Su (2013).

Let N be the number of cross-sectional units, T be the number of time periods, K be the number of regressors, and R^0 be the true number of interactive fixed effects. We consider a linear regression model with observed outcomes Y , regressors X_k , and unobserved error structure ε , namely

$$Y = \sum_{k=1}^K \beta_k^0 X_k + \varepsilon, \quad \varepsilon = \lambda^0 f^{0'} + e, \quad (1.1)$$

where Y , X_k , ε and e are $N \times T$ matrices, λ^0 is an $N \times R^0$ matrix, f^0 is a $T \times R^0$ matrix, and the regression parameters β_k^0 are scalars — the superscript zero indicates the true value of the parameters. We write β for the K -vector of regression parameters, and we denote the components of the different matrices by Y_{it} , $X_{k,it}$, e_{it} , λ_{ir}^0 and f_{tr}^0 , where $i = 1, \dots, N$, $t = 1, \dots, T$, and $r = 1, \dots, R^0$. It is convenient to introduce the notation $\beta \cdot X := \sum_{k=1}^K \beta_k X_k$. All matrices, vectors and scalars in this paper are real valued.

We consider the interactive fixed effect specification, i.e. we treat λ^0 and f^0 as nuisance parameters, which are estimated jointly with the parameters of interest β .² The advantages of the fixed effects approach are for instance that it is semi-parametric, since no assumption on the distribution of the interactive effects needs to be made, and that the regressors can be arbitrarily correlated with the interactive effect parameters.

We study the least squares (LS) estimator of model (1.1), which minimizes the sum of squared residuals to estimate the unknown parameters β , λ and f .³ To our knowledge, this estimator was first discussed in Kiefer (1980). Under an asymptotic where N and T grow to infinity, the asymptotic properties of the LS estimator were derived in Bai (2009a) for strictly exogeneous regressors, and extended in Moon and Weidner (2013) to the case of pre-determined regressors.

An important restriction of these papers is that the number of factors R^0 is assumed

¹The conventional additive model can be interpreted as a two factor interactive fixed effects model.

²When we refer to interactive fixed effects we mean that both factors and factor loadings are treated as non-random parameters. Ahn, Lee and Schmidt (2001) take a hybrid approach in that they treat the factors as non-random, but the factor loadings as random. The common correlated effects estimator of Pesaran (2006) was introduced in a context, where both the factor loadings and the factors follow certain probability laws, but it exhibits many properties of a fixed effects estimator.

³The LS estimator is sometimes called “concentrated” least squares estimator in the literature, and in an earlier version of the paper we referred to it as the “Gaussian Quasi Maximum Likelihood Estimator”, since LS estimation is equivalent to maximizing a conditional Gaussian likelihood function. Note also that for fixed β the LS estimator for λ and f is simply the principal components estimator.

to be known. However, in many empirical applications there is no consensus about the exact number of factors in the data or in the relevant economic model. If R^0 is not known beforehand, then it may be estimated consistently,⁴ but difficulties in obtaining reliable estimates for the number of factors are well-documented in the literature (see, e.g., the simulation results in Onatski (2010), and also our empirical illustration in Section 5). Furthermore, in order to use the existing inference results on R^0 one still needs a good preliminary estimator for β , so that working out the asymptotic properties of the LS-estimator for $R \geq R^0$ is still useful when taking that route.

We investigate the asymptotic properties of the LS estimator when the true number of factors R^0 is unknown and R ($\geq R^0$) number of factors are used in the estimation.⁵ We denote this estimator by $\hat{\beta}_R$.

The main result of the paper, presented in Section 3, is that under certain assumptions the LS estimator $\hat{\beta}_R$ has the same limiting distribution as $\hat{\beta}_{R^0}$ for any $R \geq R^0$ under an asymptotic where both N and T become large, while R^0 and R are constant. This implies that the LS estimator $\hat{\beta}_R$ is asymptotically robust towards inclusion of extra interactive effects in the model, and within the LS estimation framework there is no asymptotic efficiency loss from choosing R larger than R^0 . The important empirical implication of our result is that the number of factors R^0 need not be known or estimated accurately to apply the LS estimator.

To derive this robustness result, we impose more restrictive conditions than those typically assumed with known R^0 . These include that the errors e_{it} are independent and identically (iid) normally distributed and that the regressors are composed of a “low-rank” strictly stationary component, a “high-rank” strictly stationary component, and a “high-rank” pre-determined component.⁶ Notice that while some of these restrictions are necessary for our robustness result, some of them (e.g. iid normality of e_{it}) are imposed for technical reasons, because in the proof we use certain results from the theory of random matrices that are currently only available in that case (see the discussion in Section 4.3). In the Monte Carlo simulations in Section 6, we consider DGPs that violate some technical conditions to demonstrate robustness of the result.

Under less restrictive assumptions we provide intermediate results that sequentially lead to the main result in Section 4 and Appendix A.3 and A.4. In Section 4.1 we show $\sqrt{\min(N, T)}$ -consistency of the LS estimator $\hat{\beta}_R$ as $N, T \rightarrow \infty$ under very mild regularity condition on X_{it} and e_{it} , and without imposing any assumptions on λ^0 and f^0 apart from $R \geq R^0$. We thus obtain consistency of the LS estimator not only for unknown number of factors, but also for weak factors,⁷ which is an important robustness result.

In Section 4.2 we derive an asymptotic expansion of the LS profile objective function that concentrates out f and λ , for the case $R = R^0$. Given that the profile objective function is a sum of eigenvalues of a covariance matrix, its quadratic approximation is challenging because the derivatives of the eigenvalues with respect to β are not generally known. We thus cannot use a conventional Taylor expansion, but instead apply the

⁴See the discussion in Bai (2009b), supplemental material, regarding estimation of R^0 .

⁵For $R < R^0$ the LS estimator can be inconsistent, since then there are interactive fixed effects in the model which can be correlated with the regressors but are not controlled for in the estimation. We therefore restrict attention to the case $R \geq R^0$.

⁶The pre-determined component of the regressors allows for linear feedback of e_{it} into future realizations of $X_{k,it}$.

⁷See Onatski (2010; 2012) and Chudik, Pesaran and Tosetti (2011) for a discussion of “strong” vs. “weak” factors in factor models.

perturbation theory of linear operators to derive the approximation.

In Section 4.3 we provide an example that satisfies the typical assumptions imposed with known R^0 , so that $\widehat{\beta}_{R_0}$ is \sqrt{NT} consistent, but we show that $\widehat{\beta}_R$ with $R > R_0$ is only $\sqrt{\min(N, T)}$ consistent in that example. This shows that stronger conditions are required to derive our main result.

In Appendix A.3 we show faster than $\sqrt{\min(N, T)}$ -convergence of $\widehat{\beta}_R$ under assumptions that are less restrictive than those employed for the main result, in particular allowing for either cross-sectional or time-serial correlation of the errors e_{it} . In Appendix A.4 we provide an alternative version of our main result of asymptotic equivalence of $\widehat{\beta}_{R^0}$ and $\widehat{\beta}_R$, $R \geq R^0$, which is derived under high-level assumptions.

In Section 5 we follow Kim and Oka (2014) in employing the interactive fixed effects specification to study the effect of US divorce law reforms on divorce rates. This empirical example illustrates that the estimates for the coefficient β indeed become insensitive to the choice of R , once R is chosen sufficiently large, as expected from our theoretical results.

Section 6 contains Monte Carlo simulation results for a static panel model. For the simulations we consider a DGP that violates the iid normality restriction of the error term. The simulation results confirm our main result of the paper even with a relatively small sample size (e.g. $N = 100$, $T = 10$) and non-iid-normal errors. In the supplementary appendix, we report the Monte Carlo simulation results of an AR(1) panel model. It also confirms the robustness result in large samples, but in finite samples it shows more inefficiency than the static case. In general, one should expect some finite sample inefficiency from overestimating the number of factors when the sample size is small or the number of overfitted factors is large.

A few words on notation. The transpose of a matrix A is denoted by A' . For a column vectors v its Euclidean norm is defined by $\|v\| = \sqrt{v'v}$. For an $m \times n$ matrix A the Frobenius or Hilbert Schmidt norm is $\|A\|_{HS} = \sqrt{\text{Tr}(AA')}$, and the operator or spectral norm is $\|A\| = \max_{0 \neq v \in \mathbb{R}^n} \frac{\|Av\|}{\|v\|}$. Furthermore, we use $P_A = A(A'A)^\dagger A'$ and $M_A = \mathbb{1} - A(A'A)^\dagger A'$, where $\mathbb{1}$ is the $m \times m$ identity matrix, and $(A'A)^\dagger$ denotes some generalized inverse, in case A is not of full column rank. For square matrices B, C , we use $B > C$ (or $B \geq C$) to indicate that $B - C$ is positive (semi) definite. We use “wpal” for “with probability approaching one”.

2 Identification of β^0 , $\lambda^0 f^{0'}$, and R^0

In this section we provide a set of conditions under which the regression coefficient β^0 , the interactive fixed effects $\lambda^0 f^{0'}$, and the number of factors R^0 are determined uniquely by the data. Here, and throughout the whole paper, we treat λ and f as non-random parameters, i.e. all stochastics in the following are implicitly conditional on λ and f . Let $x_k = \text{vec}(X_k)$, the NT -vectorization of X_k , and let $x = (x_1, \dots, x_K)$, which is an $NT \times K$ matrix.

Assumption ID (Assumptions for Identification).

- (i) *The second moments of X_{it} and e_{it} exist for all i, t .*
- (ii) *$\mathbb{E}(e_{it}) = 0$, $\mathbb{E}(X_{it}e_{it}) = 0$, for all i, t .*
- (iii) *$\mathbb{E}[x'(M_F \otimes M_{\lambda^0})x] > 0$, for all $F \in \mathbb{R}^{T \times R}$.*

(iv) $R^0 := \text{rank}(\lambda^0 f^{0'}) \leq R$.

Theorem 2.1 (Identification). *Suppose that the Assumptions ID are satisfied. Then, β^0 , $\lambda^0 f^{0'}$, and R^0 are identified.*⁸

Assumption ID(i) imposes existence of second moments. Assumption ID(ii) is an exogeneity condition, which demands that x_{it} and e_{it} are not correlated contemporaneously, but allows for pre-determined regressors like lagged dependent variables. Assumption ID(iv) imposes that the true number of factors $R^0 := \text{rank}(\lambda^0 f^{0'})$ is bounded by a positive integer R , which cannot be too large (e.g. the trivial bound $R = N$ is not possible), since otherwise Assumption ID(iii) cannot be satisfied.

Assumption ID(iii) is a non-collinearity condition, which demands that the regressors have significant variation across i and over t after projecting out all variation that can be explained by the factor loadings λ^0 and by arbitrary factors $F \in \mathbb{R}^{T \times R}$. This generalizes the within variation assumption in the conventional panel regression with time invariant individual fixed effects, which in our notation reads $\mathbb{E}[x'(M_{1T} \otimes \mathbb{1}_N)x] > 0$.⁹ This conventional fixed effect assumption rules out time-invariant regressors. Similarly, Assumption ID(iii) rules out more general “low-rank regressors”,¹⁰ see our discussion of Assumption NC below.

3 Main Result

The estimator we investigate in this paper is the least squares (LS) estimator, which for a given choice of R reads¹¹

$$\left(\widehat{\beta}_R, \widehat{\Lambda}_R, \widehat{F}_R \right) \in \underset{\{\beta \in \mathbb{R}^K, \Lambda \in \mathbb{R}^{N \times R}, F \in \mathbb{R}^{T \times R}\}}{\text{argmin}} \left\| Y - \beta \cdot X - \Lambda F' \right\|_{HS}^2, \quad (3.1)$$

where $\|\cdot\|_{HS}$ refers to the Hilbert Schmidt norm, also called Frobenius norm. The objective function $\|Y - \beta \cdot X - \Lambda F'\|_{HS}^2$ is simply the sum of squared residuals. The estimator for β^0 can equivalently be defined by minimizing the profile objective function that concentrates out the R factors and the R factor loadings, namely

$$\widehat{\beta}_R = \underset{\beta \in \mathbb{R}^K}{\text{argmin}} \mathcal{L}_{NT}^R(\beta), \quad (3.2)$$

⁸Here, identification means that β^0 and $\lambda^0 f^{0'}$ can be uniquely recovered from the distribution of (Y, X) conditional on those parameters. Identification of the number of factors follows since $R^0 = \text{rank}(\lambda^0 f^{0'})$. The factor loadings and factors λ^0 and f^0 are not separately identified without further normalization restrictions, but the product $\lambda^0 f^{0'}$ is identified.

⁹The conventional panel regression with additive individual fixed effects and time effects requires a non-collinearity condition of the form $\mathbb{E}[x'(M_{1T} \otimes M_{1N})x] > 0$.

¹⁰We do not consider such “low-rank regressors” in this paper. Note also that Assumption A in Bai (2009a) is the sample version of our Assumption ID(iii).

¹¹The optimal $\widehat{\Lambda}_R$ and \widehat{F}_R in (3.1) are not unique, since the objective function is invariant under right-multiplication of Λ with a non-degenerate $R \times R$ matrix S , and simultaneous right-multiplication of F with $(S^{-1})'$. However, the column spaces of $\widehat{\Lambda}_R$ and \widehat{F}_R are uniquely determined.

with¹²

$$\begin{aligned}
\mathcal{L}_{NT}^R(\beta) &= \min_{\{\Lambda \in \mathbb{R}^{N \times R}, F \in \mathbb{R}^{T \times R}\}} \frac{1}{NT} \|Y - \beta \cdot X - \Lambda F'\|_{HS}^2 \\
&= \min_{F \in \mathbb{R}^{T \times R}} \frac{1}{NT} \text{Tr} [(Y - \beta \cdot X) M_F (Y - \beta \cdot X)'] \\
&= \frac{1}{NT} \sum_{r=R+1}^T \mu_r [(Y - \beta \cdot X)' (Y - \beta \cdot X)] , \tag{3.3}
\end{aligned}$$

where, $\mu_r(\cdot)$ is the r 'th largest eigenvalue of the matrix argument. Here, we first concentrated out Λ by use of its own first order condition. The resulting optimization problem for F is a principal components problem, so that the the optimal F is given by the R largest principal components of the $T \times T$ matrix $(Y - \beta \cdot X)' (Y - \beta \cdot X)$. At the optimum the projector M_F therefore exactly projects out the R largest eigenvalues of this matrix, which gives rise to the final formulation of the profile objective function as the sum over its $T - R$ smallest eigenvalues.¹³ We write $\mathcal{L}_{NT}^0(\beta)$ for $\mathcal{L}_{NT}^R(\beta)$, the profile objective function obtained for the true number of factors. Notice that we do not impose a compact parameter set for β .

Assumption SF (Strong Factor Assumption).

- (i) $0 < \text{plim}_{N,T \rightarrow \infty} \frac{1}{N} \lambda^{0'} \lambda^0 < \infty$.
- (ii) $0 < \text{plim}_{N,T \rightarrow \infty} \frac{1}{T} f^{0'} f^0 < \infty$.

Assumption NC (Non-Collinearity of X_k). Consider linear combinations $\alpha \cdot X = \sum_{k=1}^K \alpha_k X_k$ of the regressors X_k with K -vector α such that $\|\alpha\| = 1$. We assume that there exists a constant $b > 0$ such that

$$\min_{\{\alpha \in \mathbb{R}^K, \|\alpha\|=1\}} \sum_{r=R+R^0+1}^T \mu_r \left[\frac{(\alpha \cdot X)' (\alpha \cdot X)}{NT} \right] \geq b, \quad \text{wpa1.}$$

Assumption LL (Low Level Conditions for Main Result).

- (i) **Decomposition of Regressors:** $X_k = \bar{X}_k + \tilde{X}_k^{\text{str}} + \tilde{X}_k^{\text{weak}}$, for $k = 1, \dots, K$, where \bar{X}_k , \tilde{X}_k^{str} and $\tilde{X}_k^{\text{weak}}$ are $N \times T$ matrices, and
 - (i.a) **Low-Rank (strictly exogenous) Part of Regressors:** $\text{rank}(\bar{X}_k)$ is bounded as $N, T \rightarrow \infty$, and $\frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T \bar{X}_{k,it}^2 = \mathcal{O}_P(1)$.

¹²The profile objective function $\mathcal{L}_{NT}^R(\beta)$ need not be convex in β and can have multiple local minima. Depending on the dimension of β one should either perform an initial grid search or try multiple starting values for the optimization when calculating the global minimum $\hat{\beta}_R$ numerically. See also Section S.8 of the supplementary material.

¹³This last formulation of $\mathcal{L}_{NT}^R(\beta)$ is very convenient since it does not involve any explicit optimization over nuisance parameters. Numerical calculation of eigenvalues is very fast, so that the numerical evaluation of $\mathcal{L}_{NT}^R(\beta)$ is unproblematic for moderately large values of T . Since the model is symmetric under $N \leftrightarrow T$, $\Lambda \leftrightarrow F$, $Y \leftrightarrow Y'$, $X_k \leftrightarrow X_k'$ there also exists a dual formulation of $\mathcal{L}_{NT}^R(\beta)$ that involves solving an eigenvalue problem for an $N \times N$ matrix.

- (i.b) **High-Rank (strictly exogenous) Part of Regressors:** $\|\tilde{X}_k^{\text{str}}\| = \mathcal{O}_P(N^{3/4})$,
as can be justified e.g. by Lemma A.1 in the appendix.
- (i.c) **Weakly Exogenous Part of Regressors:** $\tilde{X}_{k,it}^{\text{weak}} = \sum_{\tau=1}^{t-1} \gamma_\tau e_{i,t-\tau}$, where the
real valued coefficients γ_τ satisfy $\sum_{\tau=1}^{\infty} |\gamma_\tau| < \infty$.
- (i.d) **Bounded Moments:** We assume that $\mathbb{E}|X_{k,it}|^2$, $\mathbb{E}|(M_{\lambda^0} X_k M_{f^0})_{it}|^{26}$, $\mathbb{E}|(M_{\lambda^0} X_k)_{it}|^8$
and $\mathbb{E}|(X_k M_{f^0})_{it}|^8$ are bounded uniformly over k, i, j, N and T .
- (ii) **Errors are iid Normal:** The error matrix e is independent of $\lambda^0, f^0, \bar{X}_k$, and
 \tilde{X}_k^{str} , $k = 1, \dots, K$, and its elements e_{it} are independent and identically distributed
as $\mathcal{N}(0, \sigma^2)$ across i and over t .
- (iii) **Number of Factors not Underestimated:** $R \geq R^0 := \text{rank}(\lambda^0 f^0)$.

Remarks

- (i) Assumption SF assumes that the factor f^0 and the factor loading λ^0 are strong. The strong factor assumption is regularly imposed in the literature on large N and T factor models, including Bai and Ng (2002), Stock and Watson (2002) and Bai (2009a).
- (ii) Assumption NC assumes that there exists significant sampling variation in the regressors after concentrating out $R + R^0$ factors (or factor loadings). It is a sample version of the identification Assumption ID(iii), and it is essentially equivalent to Assumption A of Bai (2009a), but avoids mentioning the unobserved loadings λ^0 .¹⁴
- (iii) Assumption NC is violated if there exists a linear combination $\alpha \cdot X$ of the regressors with $\alpha \neq 0$ and $\text{rank}(\alpha \cdot X) \leq R + R^0$, i.e. the assumption rules out “low-rank regressors” like time invariant regressors or cross-sectionally invariant regressors. These low-rank regressors require a special treatment in the interactive fixed effect model, see Bai (2009a) and Moon and Weidner (2013), and we do not consider them in the present paper. If one is not interested explicitly in their regression coefficients, then one can always eliminate the low-rank regressors by an appropriate projection of the data, e.g. subtraction of the time (or cross-sectional) means from the data eliminates all time-invariant (or cross-sectionally invariant) regressors.
- (iv) The norm restriction in Assumption LL(i.b) is a high level assumption. It is satisfied as long as $\tilde{X}_{k,it}^{\text{str}}$ is mean zero and weakly correlated across i and over t , for details see Appendix A.1 and Lemma A.1 there.
- (v) Assumption LL(i) imposes that each regressor consists of three parts: (a) a strictly exogenous low rank component, (b) a strictly exogenous component satisfying a norm restriction, and (c) a weakly exogenous component that follows a linear process with innovation given by the lagged error term e_{it} . For example, if $X_{k,it} \sim iid \mathcal{N}(\mu_k, \sigma_k^2)$, independent of e , then we have $\bar{X}_{k,it} = \mu_k$, $\tilde{X}_{k,it}^{\text{str}} \sim iid \mathcal{N}(0, \sigma_k^2)$ and $\tilde{X}_k^{\text{weak}} = 0$.

¹⁴By dropping the expected value from Assumption ID(iii) and replacing the zero lower bound by a positive constant one obtains $\inf_F [x'(M_F \otimes M_{\lambda^0})x/NT] \geq b > 0$, wpa1, which is equivalent to Assumption A of Bai (2009a), and can also be rewritten as $\min_{\|\alpha\|=1} \inf_F \text{Tr}[M_{\lambda^0}(\alpha \cdot X)'M_F(\alpha \cdot X)/NT] \geq b$. A slightly stronger version of the Assumption, which avoids mentioning the unobserved factor loading λ^0 , reads $\min_{\|\alpha\|=1} \inf_F \inf_\lambda \text{Tr}[M_\lambda(\alpha \cdot X)'M_F(\alpha \cdot X)/NT] \geq b$, where $F \in \mathbb{R}^{T \times R}$ and $\lambda \in \mathbb{R}^{N \times R^0}$, and this slightly stronger version is equivalent to Assumption NC.

Assumption LL(i) is also satisfied for a stationary panel VAR with interactive fixed effects as in Holtz-Eakin, Newey and Rosen (1988). A special case of this is a dynamic panel regression with fixed effects, where $Y_{it} = \beta Y_{i,t-1} + \lambda_i^{0'} f_t^0 + e_{it}$, with $|\beta| < 1$ and “infinite history”. In this case, we have $X_{it} = Y_{i,t-1} = \bar{X}_{it} + \tilde{X}_{it}^{\text{str}} + \tilde{X}_{it}^{\text{weak}}$, where $\bar{X}_{it} = \lambda_i^{0'} \sum_{\tau=1}^{\infty} \beta^{\tau-1} f_{t-\tau}^0$, $\tilde{X}_{it}^{\text{str}} = \sum_{\tau=t}^{\infty} \beta^{\tau-1} e_{i,t-\tau}$, and $\tilde{X}_{it}^{\text{weak}} = \sum_{\tau=0}^{t-1} \beta^{\tau-1} e_{i,t-\tau}$.

- (vi) Assumption LL(i) is more restrictive than Assumption 5 in Moon and Weidner (2013), where R^0 is assumed to be known. However, it is more general than the restriction on the regressors in Pesaran (2006), where – in our notation – the decomposition $X_k = \bar{X}_k + \tilde{X}_k^{\text{str}}$ is imposed, but the lower rank component \bar{X}_k needs to satisfy further assumptions, and the weakly exogenous component $\tilde{X}_k^{\text{weak}}$ is not considered. Bai (2009a) requires no such decomposition, but imposes strict exogeneity of the regressors.
- (vii) Among the conditions in Assumption LL, the iid normality condition in Assumption LL(ii) may be the most restrictive. In Appendix A.4 we provide an alternative version of Theorem 3.1 that imposes more general high-level conditions. Verifying those high-level conditions requires results on the eigenvalues and eigenvectors of random covariance matrices, which can be verified for iid normal errors by using known results from the random matrix theory literature, see Section 4.3 for more details. We believe, however, that those high-level conditions and thus our main result hold more generally, and we explore non-normal and serially correlated errors in our Monte Carlo simulations below.

Theorem 3.1 (Main Result). *Let Assumption SF, NC and LL hold and consider a limit $N, T \rightarrow \infty$ with $N/T \rightarrow \kappa^2$, $0 < \kappa < \infty$. Then we have*

$$\sqrt{NT}(\hat{\beta}_R - \beta^0) = \sqrt{NT}(\hat{\beta}_{R^0} - \beta^0) + o_P(1).$$

Theorem 3.1 follows from Theorem A.3 and Lemma A.4 in the appendix, whose prove is given in the supplementary material. The theorem guarantees that the asymptotic distribution of $\hat{\beta}_R$, $R \geq R^0$, is identical to that of $\hat{\beta}_{R^0}$ in (3.4) below.

The limiting distribution of $\sqrt{NT}(\hat{\beta}_{R^0} - \beta^0)$ with known R^0 is available in the existing literature. According to Bai (2009a) and Moon and Weidner (2013),

$$\sqrt{NT}(\hat{\beta}_{R^0} - \beta^0) \Rightarrow \mathcal{N}(-\kappa \text{plim } W^{-1}B, \sigma^2 \text{plim } W^{-1}), \quad (3.4)$$

where W is the $K \times K$ matrix with elements $W_{k_1 k_2} = \frac{1}{NT} \text{Tr}(M_{\lambda^0} X_{k_1} M_{f^0} X'_{k_2})$, B is the K -vector with elements $B_k = \frac{1}{N} \text{Tr}[P_{f^0} \mathbb{E}(e' X_k)]$.¹⁵

The result (3.4) holds under the assumptions of Theorem 3.1 and also assuming that $\text{plim } W^{-1}B$ and $\text{plim } W^{-1}$ exist, where plim refers to the probability limit as $N, T \rightarrow \infty$. Note that Assumption NC guarantees that W is invertible asymptotically. The asymptotic

¹⁵The asymptotic distribution in (3.4) can also be derived from Corollary 4.3 below under more general conditions than in Assumption LL (see Moon and Weidner (2013) for details). Here we have used the homoscedasticity of e_{it} to simplify the structure of the asymptotic variance and bias. Bai (2009a) finds further asymptotic bias in $\hat{\beta}_{R^0}$ due to heteroscedasticity and correlation in e_{it} , which in our asymptotic result is ruled out by Assumption LL(ii), but is studied in our Monte Carlo simulations below. Moon and Weidner (2013) work out the additional asymptotic bias in $\hat{\beta}_{R^0}$ due to pre-determined regressors, which is allowed for in Theorem 3.1.

bias in (3.4) is an incidental parameter bias due to pre-determined regressors and is equal to zero for strictly exogenous regressors (for which $\mathbb{E}(e'X_k) = 0$); it generalizes the well-known Nickell (1981) bias of the within-group estimator for dynamic panel models.

Estimators for σ^2 , W and B are given by¹⁶

$$\hat{\sigma}_R^2 = \frac{1}{(N-R)(T-R) - K} \sum_{i=1}^N \sum_{t=1}^T (\hat{e}_{R,it})^2, \quad \widehat{W}_{R,k_1k_2} = \frac{1}{NT} \text{Tr} \left(M_{\widehat{\Lambda}_R} X_{k_1} M_{\widehat{F}_R} X'_{k_2} \right),$$

$$\widehat{B}_{R,k} = \sum_{t=1}^T \sum_{\tau=t+1}^{t+M} P_{\widehat{F}_R,t\tau} \left[\frac{1}{N} \sum_{i=1}^N \hat{e}_{R,it} X_{k,i\tau} \right],$$

where $\hat{e}_{R,it}$ denotes the $(i,t)^{th}$ element of $\widehat{e}_R = Y - \widehat{\beta}_R \cdot X - \widehat{\Lambda}_R \widehat{F}'_R$, and $P_{\widehat{F}_R,t\tau}$ denotes the $(t,\tau)^{th}$ element of $P_{\widehat{F}_R} = \mathbb{1}_T - M_{\widehat{F}_R} = \widehat{F}_R (\widehat{F}'_R \widehat{F}_R)^\dagger \widehat{F}'_R$, and $M \in \{1, 2, 3, \dots\}$ is a bandwidth parameter that also depends on the sample size N, T . Let \widehat{W}_R and \widehat{B}_R be the matrix and vector with elements \widehat{W}_{R,k_1k_2} and $\widehat{B}_{R,k}$, respectively.

The next theorem establishes the consistency of these estimators. Let $\lambda^{\text{red}} \in \mathbb{R}^{N \times (R-R^0)}$ and $f^{\text{red}} \in \mathbb{R}^{T \times (R-R^0)}$ be the leading $R-R^0$ principal components obtained from the $N \times T$ matrix $M_{\lambda^0} e M_{f^0}$, i.e. λ^{red} and f^{red} minimize the objective function $\|M_{\lambda^0} e M_{f^0} - \lambda^{\text{red}} f^{\text{red}'}\|_{HS}^2$, analogous to $\widehat{\Lambda}_R$ and \widehat{F}_R defined in (3.1).¹⁷

Theorem 3.2 (Consistency of Bias and Variance Estimators).

- (i) Let the conditions of Theorem 3.1 hold. Then we have $\|P_{\widehat{F}_R} - P_{[f^0, f^{\text{red}}]}\| = o_p(1)$, $\|P_{\widehat{\Lambda}_R} - P_{[\lambda^0, \lambda^{\text{red}}]}\| = o_p(1)$, $\hat{\sigma}_R^2 = \sigma^2 + o_p(1)$, and $\widehat{W}_R = W + o_p(1)$.
- (ii) In addition, let $X_{k,\cdot t} = (X_{k,1t}, \dots, X_{k,Nt})'$, and assume that (1) γ_τ in Assumption LL(i.c) satisfies $|\gamma_\tau| < c\tau^{-d}$ for some $c > 0$ and $d > 1$, (2) $\|\lambda_i^0\|$ and $\|f_t^0\|$ are uniformly bounded over i, t and N, T , (3) $\max_t \|X_{k,\cdot t}\| = \mathcal{O}_P(\sqrt{N} \log N)$,¹⁸ and (4) the bandwidth $M \rightarrow \infty$ such that $M(\log T)^2/T^{1/6} \rightarrow 0$. Then, we have $\widehat{B}_R = B + o_p(1)$.

Combining Theorems 3.1 and 3.2 and the asymptotic distribution in (3.4) allows inference on β , for $R \geq R^0$. In particular, the bias corrected estimator $\widehat{\beta}_R^{\text{BC}} = \widehat{\beta}_R + \frac{1}{T} \widehat{W}_R^{-1} \widehat{B}_R$ satisfies¹⁹

$$\sqrt{NT}(\widehat{\beta}_R^{\text{BC}} - \beta^0) \Rightarrow \mathcal{N}(0, \sigma^2 W^{-1}).$$

¹⁶The first factor in $\hat{\sigma}^2$ reflects the degree of freedom correction from estimating Λ , F and β , but could simply be chosen as $1/NT$ for the purpose of consistency. Note also that $P_{\widehat{F}_R,t\tau} = \mathcal{O}_P(1/T)$, which explains why no $1/T$ factor is required in the definition of $\widehat{B}_{R,k}$.

¹⁷The superscript ‘‘red’’ stands for redundant, because it turns out that λ^{red} and f^{red} are asymptotically close to the $R - R^0$ redundant principal components that are estimated in (3.1).

¹⁸The high-level assumption $\max_t \|X_{k,\cdot t}\| = \mathcal{O}_P(\sqrt{N} \log N)$ can be shown to be satisfied for the regressor component $\widetilde{X}_{k,it}^{\text{weak}}$ above, and can be justified for the other regressor components e.g. by assuming that \overline{X}_k and $\widetilde{X}_k^{\text{str}}$ are uniformly bounded.

¹⁹Instead of estimating the bias analytically one can use the result that the bias is of order T^{-1} and perform split panel bias correction as in Dhaene and Jochmans (2010), which instead of the conditions of Theorem 3.2(ii) only requires some stationary condition over time.

Heuristic Discussion of the Main Result

Intuitively, the inclusion of unnecessary factors in the LS estimation is similar to the inclusion of irrelevant regressors in an OLS regression. In the OLS case it is well known that if those irrelevant extra regressors are uncorrelated with the regressors of interest, then they have no effect on the asymptotic distribution of the regression coefficients of interest. It is therefore natural to expect that if the extra estimated factors in \widehat{F}_R are asymptotically uncorrelated with the regressors, then the result of Theorem 3.1 should hold. To explore this, remember that \widehat{F}_R is given by the first R principal components of the matrix $(Y - \widehat{\beta}_R \cdot X)'(Y - \widehat{\beta}_R \cdot X)$, and write

$$Y - \widehat{\beta}_R \cdot X = \lambda^0 f^{0'} + e - (\widehat{\beta}_R - \beta^0) \cdot X.$$

The strong factor assumption and the consistency of $\widehat{\beta}_R$ guarantee that the first R^0 principal components of $(Y - \widehat{\beta}_R \cdot X)'(Y - \widehat{\beta}_R \cdot X)$ are close to f^0 asymptotically, i.e. the true factors are correctly picked up by the principal component estimator. The additional $R - R^0$ principal components that are estimated for $R > R^0$ cannot pick up anymore true factors and are thus mostly determined by the remaining term $e - (\widehat{\beta}_R - \beta^0) \cdot X$. The key question for the properties of the extra estimated factors, and thus of $\widehat{\beta}_R$, is therefore whether the principal components obtained from $e - (\widehat{\beta}_R - \beta^0) \cdot X$ are dominated by e or by $(\widehat{\beta}_R - \beta^0) \cdot X$. Only if they are dominated by e can we expect the extra factors in \widehat{F}_R to be uncorrelated with X and thus the result in Theorem 3.1 to hold. The result on $P_{\widehat{F}_R}$ in Theorem 3.2 shows that the additional estimated factors are indeed close to f^{red} , i.e. are mostly determined by e , but this result is far from obvious a priori, as the following discussion shows.

Under our assumptions we have $\|e\| = \mathcal{O}_P(\sqrt{N})$ and $\|X_k\| = \mathcal{O}_P(\sqrt{NT})$ as N and T grow at the same rate. Thus, if the convergence rate of $\widehat{\beta}_R$ is faster than \sqrt{N} , i.e. $\|\widehat{\beta}_R - \beta^0\| = o_P(\sqrt{N})$, then we have $\|e\| \gg \left\| (\widehat{\beta}_R - \beta^0) \cdot X \right\|$ asymptotically, and we expect the extra \widehat{F}_R to be dominated by e . A crucial step in the derivation of Theorem 3.1 is therefore to show faster than \sqrt{N} convergence of $\widehat{\beta}_R$. Conversely, we expect counter examples to the main result to be such that the convergence rate of the estimator $\widehat{\beta}_R$ is not faster than \sqrt{N} , and we provide such a counter example – which, however, violates Assumptions LL – in Section 4.3 below. Whether the intuition about “inclusion of irrelevant regressors” carries over to the “inclusion of irrelevant factors” thus crucially depends on the convergence rate of $\widehat{\beta}_R$.

4 Asymptotic Theory and Discussion

Here we introduce key intermediate results on the way to deriving the main Theorem 3.1 stated above. These intermediate results may be useful independently of the main result, e.g. Moon and Weidner (2013) and Moon, Shum, and Weidner (2014) crucially use the results established in Section 4.2 for the case of known $R = R^0$. The assumptions introduced below are all implied by the low-level Assumptions LL above, according to Lemma A.4 in the appendix.

4.1 Consistency of $\widehat{\beta}_R$

Here we present a consistency result for $\widehat{\beta}_R$ under an arbitrary asymptotic $N, T \rightarrow \infty$, i.e. without the assumption that N and T grow at the same rate, which is imposed everywhere else in the paper. In addition to Assumption NC we require the following high level assumptions to obtain the result.

Assumption SN (Spectral Norm of X_k and e).

(i) $\|X_k\| = \mathcal{O}_P(\sqrt{NT})$, $k = 1, \dots, K$.

(ii) $\|e\| = \mathcal{O}_P(\sqrt{\max(N, T)})$.

Assumption EX (Weak Exogeneity of X_k). $\frac{1}{\sqrt{NT}} \text{Tr}(X_k e') = \mathcal{O}_P(1)$, $k = 1, \dots, K$.

Theorem 4.1. *Let Assumptions SN, EX and NC be satisfied and let $R \geq R^0$. For $N, T \rightarrow \infty$ we then have $\sqrt{\min(N, T)} (\widehat{\beta}_R - \beta^0) = \mathcal{O}_P(1)$.*

Remarks

- (i) One can justify Assumption SN(i) by use of the norm inequality $\|X_k\| \leq \|X_k\|_{HS}$ and the fact that $\|X_k\|_{HS}^2 = \sum_{i,t} X_{k,it}^2 = \mathcal{O}_P(NT)$, where the last step follows e.g. if $X_{k,it}$ has a uniformly bounded second moment.
- (ii) Assumption SN(ii) is a condition on the largest eigenvalue of the random covariance matrix $e'e$, which is often studied in the literature on random matrix theory, e.g. Geman (1980), Bai, Silverstein, Yin (1988), Yin, Bai, and Krishnaiah (1988), Silverstein (1989). The results in Latala (2005) show that $\|e\| = \mathcal{O}_P(\sqrt{\max(N, T)})$ if e has independent entries with mean zero and uniformly bounded fourth moment. Weak dependence of the entries e_{it} across i and over t is also permissible, see Appendix A.1
- (iii) Assumption EX requires exogeneity of the regressors X_k , allowing for pre-determined regressors, and some weak dependence of $X_{k,it}e_{it}$ across i and over t .²⁰
- (iv) The theorem imposes no restriction at all on f^0 and λ^0 , apart from the condition $R \geq \text{rank}(\lambda^0 f^{0'})$.²¹ In particular, the strong factor Assumption SF is not imposed here, i.e. consistency of $\widehat{\beta}_R$ holds independently of whether the factors are strong, weak, or not present at all. This is an important robustness result, which is new in the literature.
- (v) Under an asymptotic where N and T grow at the same rate, which is imposed everywhere else in the paper, Theorem 4.1 shows \sqrt{N} (or equivalently \sqrt{T}) consistency of the estimator $\widehat{\beta}_R$.
- (vi) \sqrt{N} consistency of $\widehat{\beta}_R$ implies that the residuals $Y - \widehat{\beta}_R \cdot X$ will be asymptotically close to $\lambda^0 f^{0'} + e$.²² This allows consistent estimation of R^0 under a strong factor Assumption SF by employing the known techniques on factor models without regressors (by applying, e.g. , Bai and Ng (2002) to $Y - \widehat{\beta}_R \cdot X$), as also discussed in Bai (2009b).²³

²⁰Note that $\frac{1}{\sqrt{NT}} \text{Tr}(X_k e') = \frac{1}{\sqrt{NT}} \sum_i \sum_t X_{k,it} e_{it}$.

²¹This is the main reason why we use a slightly different non-collinearity Assumption NC, which avoids mentioning λ^0 , compared to Bai (2009a).

²²In the sense that $\|(Y - \widehat{\beta}_R \cdot X) - (\lambda^0 f^{0'} + e)\| = \|(\widehat{\beta}_R - \beta) \cdot X\| = \mathcal{O}_P(\sqrt{N})$.

²³Bai (2009b) does not prove the required consistency and convergence rate of $\widehat{\beta}_R$, for $R \geq R^0$.

- (vii) Having a consistent estimator for R^0 , say \widehat{R} , one can calculate $\widehat{\beta}_{\widehat{R}}$, which will be asymptotically equal to $\widehat{\beta}_{R^0}$. In practice, however, the finite sample properties of the estimator $\widehat{\beta}_{\widehat{R}}$ crucially depend on the finite sample properties of \widehat{R} . Many recent papers have documented difficulties in obtaining reliable estimates for R^0 at finite sample (see, e.g., the simulation results of Onatski (2010) and Ahn and Horenstein (2013)), and those difficulties are also illustrated by our empirical example in Section 5.

4.2 Quadratic Approximation of $\mathcal{L}_{NT}^0(\beta)(:= \mathcal{L}_{NT}^{R^0}(\beta))$

To derive the limiting distribution of $\widehat{\beta}_R$, we study the asymptotic properties of the profile objective function $\mathcal{L}_{NT}^R(\beta)$ around β^0 . The expression in (3.3) cannot easily be discussed by analytic means, since no explicit formula for the eigenvalues of a matrix is available. In particular, a standard Taylor expansion of $\mathcal{L}_{NT}^R(\beta)$ around β^0 cannot easily be derived. Here, we consider the case of known $R = R^0$ and we perform a joint expansion of the corresponding profile objective function $\mathcal{L}_{NT}^0(\beta)$ in the regression parameters β and in the idiosyncratic error terms e . To perform this joint expansion we apply the perturbation theory of linear operators (e.g., Kato (1980)). We thereby obtain an approximate quadratic expansion of $\mathcal{L}_{NT}^0(\beta)$ in β , which can be used to derive the first order asymptotic theory of the LS estimator $\widehat{\beta}_{R^0}$, see Appendix A.2 for details. In addition to the $K \times K$ matrix W already defined in Section 3 we now also define

$$\begin{aligned} C_k^{(1)} &= \frac{1}{\sqrt{NT}} \text{Tr}(M_{\lambda^0} X_k M_{f^0} e'), \\ C_k^{(2)} &= -\frac{1}{\sqrt{NT}} \left[\text{Tr}(e M_{f^0} e' M_{\lambda^0} X_k f^0 (f^{0'} f^0)^{-1} (\lambda^{0'} \lambda^0)^{-1} \lambda^{0'}) \right. \\ &\quad + \text{Tr}(e' M_{\lambda^0} e M_{f^0} X_k' \lambda^0 (\lambda^{0'} \lambda^0)^{-1} (f^{0'} f^0)^{-1} f^{0'}) \\ &\quad \left. + \text{Tr}(e' M_{\lambda^0} X_k M_{f^0} e' \lambda^0 (\lambda^{0'} \lambda^0)^{-1} (f^{0'} f^0)^{-1} f^{0'}) \right]. \end{aligned} \quad (4.1)$$

Let $C^{(1)}$ and $C^{(2)}$ be the K -vectors with elements $C_k^{(1)}$ and $C_k^{(2)}$, respectively.

Theorem 4.2. *Let Assumptions SF and SN be satisfied. Suppose that $N, T \rightarrow \infty$ with $N/T \rightarrow \kappa^2$, $0 < \kappa < \infty$. Then we have*

$$\mathcal{L}_{NT}^0(\beta) = \mathcal{L}_{NT}^0(\beta^0) - \frac{2}{\sqrt{NT}} (\beta - \beta^0)' (C^{(1)} + C^{(2)}) + (\beta - \beta^0)' W (\beta - \beta^0) + \mathcal{L}_{NT}^{0,\text{rem}}(\beta),$$

where the remainder term $\mathcal{L}_{NT}^{0,\text{rem}}(\beta)$ satisfies for any sequence $c_{NT} \rightarrow 0$

$$\sup_{\{\beta: \|\beta - \beta^0\| \leq c_{NT}\}} \frac{|\mathcal{L}_{NT}^{0,\text{rem}}(\beta)|}{\left(1 + \sqrt{NT} \|\beta - \beta^0\|\right)^2} = o_p\left(\frac{1}{NT}\right).$$

The bound on remainder²⁴ in Theorem 4.2 is such that it has no effect on the first order

²⁴The expansion in Theorem 4.2 contains a term that is linear in β and linear in e ($C^{(1)}$ term), a term that is

asymptotic theory of $\widehat{\beta}_{R^0}$, as stated in the following corollary (see also Andrews (1999)).

Corollary 4.3. *Let Assumptions SF, SN, EX and NC be satisfied. In the limit $N, T \rightarrow \infty$ with $N/T \rightarrow \kappa^2$, $0 < \kappa < \infty$, we then have $\sqrt{NT} \left(\widehat{\beta}_{R^0} - \beta^0 \right) = W^{-1} \left(C^{(1)} + C^{(2)} \right) + o_P \left(1 + \|C^{(1)}\| \right)$. If we furthermore assume that $C^{(1)} = \mathcal{O}_P(1)$, then we obtain*

$$\sqrt{NT} \left(\widehat{\beta}_{R^0} - \beta^0 \right) = W^{-1} \left(C^{(1)} + C^{(2)} \right) + o_P(1) = \mathcal{O}_P(1).$$

Note that our assumptions already guarantee $C^{(2)} = \mathcal{O}_P(1)$ and that W is invertible with $W^{-1} = \mathcal{O}_P(1)$, so this need not be explicitly assumed in Corollary 4.3.

Remarks

- (i) More details on the expansion of $\mathcal{L}_{NT}^0(\beta)$ are provided in Appendix A.2 and the formal proofs can be found in Section S.2 of the supplementary appendix.
- (ii) Corollary 4.3 allows to replicate the results in Bai (2009a) and Moon and Weidner (2013) on the asymptotic distribution of $\widehat{\beta}_{R^0}$, including the result in formula (3.4) above.²⁵ The assumptions of the corollary do not restrict the regressor to be strictly exogenous and do not impose Assumption LL.
- (iii) If one weakens Assumption SN(ii) to $\|e\| = o_P(N^{2/3})$, then Theorem 4.2 still continues to hold. If we assume that $C^{(2)} = \mathcal{O}_P(1)$, then Corollary 4.3 also holds under this weaker condition on $\|e\|$.

4.3 Remarks on Deriving the Convergence Rate and Asymptotic Distribution of $\widehat{\beta}_R$ for $R > R^0$.

An example that motivates stronger restrictions

The results in Bai (2009a) and Corollary 4.3 above show that under appropriate assumptions the estimator $\widehat{\beta}_R$ is \sqrt{NT} -consistent for $R = R^0$. For $R > R^0$ we know from Theorem 4.1 that $\widehat{\beta}_R$ is \sqrt{N} consistent as N and T grow at the same rate, but we have not shown faster than \sqrt{N} convergence of $\widehat{\beta}_R$ for $R > R^0$, yet, which according to the heuristic discussion at the end of Section 3 is a very important intermediate step to obtain our main result.²⁶ However, one might not obtain a faster than \sqrt{N} convergence rate of $\widehat{\beta}_R$ for $R > R^0$ without imposing further restrictions, as the following example shows.

linear in β and quadratic in e ($C^{(2)}$ term), and a term that is quadratic in β (W term). All higher order terms of the expansion are contained in the remainder term $\mathcal{L}_{NT}^{0,\text{rem}}(\beta)$.

²⁵Let ρ , $D(\cdot)$, D_0 , D_Z , B_0 and C_0 be the notation used in Assumption A and Theorem 3 of Bai (2009a), and let Bai's assumptions be satisfied. Then, our κ , W , $C^{(1)}$ and $C^{(2)}$ satisfy $\kappa = \rho^{-1/2}$, $W = D(f^0) \rightarrow_p D > 0$, $C^{(1)} \rightarrow_d \mathcal{N}(0, D_Z)$ and $W^{-1}C^{(2)} \rightarrow_p \rho^{1/2}B_0 + \rho^{-1/2}C_0$. Corollary 4.3 can therefore be used to replicate Theorem 3 in Bai (2009a). For more details and extensions of this we refer to Moon and Weidner (2013).

²⁶One reason why $\widehat{\beta}_R$ might only converge at \sqrt{N} rate, but not faster, are weak factors (both for $R > R^0$ and for $R = R^0$). A weak factor (see e.g. Onatski (2010; 2012) and Chudik, Pesaran and Tosetti (2011)) might not be picked up at all or might only be estimated very inaccurately by the principal components estimator \widehat{F}_R , in which case that factor is not properly accounted for in the LS estimation procedure. If this happens and the weak factor is correlated with the regressors, then there is some uncorrected weak endogeneity problem, and $\widehat{\beta}_R$ will only converge at \sqrt{N} rate. We do not consider the issue of weak factors any further in this paper.

Example. Let $R^0 = 0$ (no true factors) and $K = 1$ (one regressor). The true model reads $Y_{it} = \beta^0 X_{it} + e_{it}$, and we consider the following data generating process (DGP)

$$X_{it} = a\tilde{X}_{it} + \lambda_{x,i}f_{x,t}, \quad e = \left(\mathbb{1}_N + c \frac{\lambda_x \lambda'_x}{N} \right) u \left(\mathbb{1}_T + c \frac{f_x f'_x}{T} \right),$$

where e and u are $N \times T$ matrices with entries e_{it} and u_{it} , respectively, and λ_x is an N -vector with entries $\lambda_{x,i}$, and f_x is a T -vector with entries $f_{x,t}$. Let \tilde{X}_{it} and u_{it} be mutually independent iid standard normally distributed random variables. Let $\lambda_{x,i} \in \mathcal{B}$ and $f_{x,t} \in \mathcal{B}$ be non-random sequences with bounded range $\mathcal{B} \subset \mathbb{R}$ such that $\frac{1}{N} \sum_{i=1}^N \lambda_{x,i}^2 \rightarrow 1$ and $\frac{1}{T} \sum_{t=1}^T f_{x,t}^2 \rightarrow 1$ asymptotically.²⁷ Consider $N, T \rightarrow \infty$ such that $N/T \rightarrow \kappa^2$, $0 < \kappa < \infty$, and let $0 < a < (1/2)^{2/3} \min(\kappa^2, \kappa^{-2})$ and $c \geq \frac{(2+\sqrt{2})(1+\kappa)(1+\sqrt{3}a^{-1/4})}{\min(1,\kappa)[1/2-a^{3/2} \max(\kappa, \kappa^{-1})]}$.²⁸ For this DGP one can show that $\hat{\beta}_1$, the LS-estimator with $R = 1 > R^0$, only converges at a rate of \sqrt{N} to β^0 , but not faster.

The proof of the last statement is provided in the supplementary material. The DGP in this example satisfies all the assumptions imposed in Corollary 4.3 to derive the limiting distribution of the LS-estimator for $R = R^0$, including \sqrt{NT} -consistency of $\hat{\beta}_R$ for $R = R^0$ ($=0$ in this example). It also satisfies all the regularity conditions imposed in Bai (2009a).²⁹ The aspect that is special about this DGP is that λ_x and f_x feature both in X_{it} and in the second moment structure of e_{it} . The heuristic discussion at the end of Section 3 provides some intuition why this can be problematic, because the leading principal components obtained from only the error matrix e will have a strong sample correlation with X_{it} for this DGP.

Faster than \sqrt{N} convergence of $\hat{\beta}_R$

In Appendix A.3, we summarize our results on faster than \sqrt{N} convergence of $\hat{\beta}_R$ for $R \geq R^0$. The above example shows that this requires more restrictive assumptions than those imposed for the analysis of the case $R = R^0$ above, but the assumption that we impose for this intermediate results are still significantly weaker than the Assumptions LL required for our main result above, in particular either cross-sectional correlation or time-series correlation of e_{it} are still allowed.

In that appendix we also provide one set of assumptions (Assumption DX-2) for faster than \sqrt{N} convergence such that no additional conditions on e are required, but where the regressors are restricted to essentially be lagged dependent variables in an AR(p) model with factors.

²⁷We could also allow λ_x and f_x to be random (but independent of e and \tilde{X}) and we could let the range of \mathcal{B} be unbounded. We only assume non-random λ_x and f_x to guarantee that the DGP satisfies Assumption D of Bai (2009a), namely that X and e are independent (otherwise we only have mean-independence, i.e. $\mathbb{E}(e|X) = 0$). Similarly, we only assume bounded \mathcal{B} to satisfy the restrictions on e_{it} imposed in Assumption C of Bai (2009a).

²⁸The bounds on the constants a and c imposed here are sufficient, but not necessary for the result of no faster than \sqrt{N} convergence of $\hat{\beta}_1$. Simulation evidence suggests that this result holds for a much larger range of a, c values.

²⁹See Section S.9 in the supplementary material for details.

On the role of the iid normality of e_{it}

We establish the asymptotic equivalence of $\widehat{\beta}_R$ and $\widehat{\beta}_{R^0}$ in Theorem 3.1 by showing that the LS objective function $\mathcal{L}_{NT}^R(\beta)$ can, up to a constant, be uniformly well approximated by $\mathcal{L}_{NT}^0(\beta)$ in shrinking neighborhoods around the true parameter. For this, we need not only the faster than \sqrt{N} convergence rate of $\widehat{\beta}_R$, but also require the Assumption EV in Appendix A.4. This is a high-level assumption on the eigenvalues and eigenvectors of the random covariance matrices EE' and $E'E$, where $E = M_{\lambda^0}eM_{f^0}$. The assumption essentially requires the eigenvalues of those matrices to be sufficiently separated from each other, as well as the eigenvectors of those matrices to be sufficiently uncorrelated with the regressors X_k , and with eP_{f^0} and $P_{\lambda^0}e$.

We use the iid normality of e_{it} to verify those high-level conditions in Section S.4.2 of the supplementary appendix. There are three reasons why we can currently only verify those conditions for iid normal errors:

- (i) The random matrix theory literature studies the eigenvalues and eigenvectors of random covariance matrices of the form ee' and $e'e$, while we have to deal with the additional projectors M_{λ^0} and M_{f^0} in the random covariance matrices. These additional projections stem from integrating out the true factors and factor loadings of the model. If the error distribution is *iid* normal, and independent from λ^0 and f^0 , then these projections are unproblematic, since the distribution of e is rotationally invariant from the left and right in that case, so that the projections are mathematically equivalent to a reduction of the sample size by R^0 in both panel dimensions.
- (ii) In the iid normal case one can furthermore use the invariance of the distribution of e under orthonormal rotations from the left and from the right to also fully characterize the distribution of the eigenvectors of EE' and $E'E$.³⁰ The conjecture in the random matrix theory literature is that the limiting distribution of the eigenvectors of a random covariance matrix is “distribution free”, i.e. is independent of the particular distribution of e_{it} , see, e.g., Silverstein (1990) and Bai (1999). However, we are not currently aware of a formulation and corresponding proof of this conjecture that is sufficient for our purposes, i.e. that would allow us to verify our high-level Assumption EV more generally.
- (iii) We also require certain properties of the eigenvalues of EE' and $E'E$. Eigenvalues are studied more intensely than eigenvectors in the random matrix theory literature, and it is well-known that the properly normalized empirical distribution of the eigenvalues (the so called empirical spectral distribution) of an *iid* sample covariance matrix converges to the Marčenko-Pastur-law (Marčenko and Pastur (1967)) for asymptotics where N and T grow at the same rate. This result does not require normality, and results on the limiting spectral distribution are also known for non-iid matrices. However, to check our high-level Assumption EV we also need results on the convergence rate of the empirical spectral distribution to its limit law, which is an ongoing research subject in the literature, e.g. Bai (1993), Bai, Miao and Yao (2004), Götze and Tikhomirov (2010), and we are currently only aware of results on this convergence rate for the case of either iid or iid normal errors. To verify the high-level assumption we furthermore use a result from Johnstone (2001) and Soshnikov (2002) that shows that the properly normalized few largest eigenvalues of

³⁰Rotational invariance implies that the distribution of the normalized eigenvectors is given by the Haar measure of a rotation group manifold.

EE' and EE' converge to the Tracy-Widom law, and to our knowledge this result is not established for error distributions that are not iid normal.

In spite of these severe mathematical challenges, we believe that in principle our high-level Assumption EV could be verified for more general error distributions, implying that our main result of asymptotic equivalence of $\hat{\beta}_R$ and $\hat{\beta}_{R0}$ holds more generally. This is also supported by our Monte Carlo simulations, where we explore non-independent and non-normal error distributions.

5 Empirical Illustration

As an illustrative empirical example, we estimate the dynamic effects of unilateral divorce law reforms on the state-wise divorce rates in the US. The impact of the divorce law reform has been studied by many researches (e.g., Allen (1992), Peters (1986; 1992), Gray (1998), Friedberg (1998), Wolfers (2006), and Kim and Oka (2014)). In this section we revisit this topic, extending Wolfers (2006) and Kim and Oka (2014) by controlling for interactive fixed effects and also a lagged dependent variable.

Let Y_{it} denote the number of divorces per 1000 people in state i at time t , and let D_i denote the year in which state i introduced the unilateral divorce law, i.e. before year D_i state i had a consent divorce law, while from D_i onwards state i had a unilateral “no-fault” divorce law, which lowers the barrier for divorce. The goal is to estimate the dynamic effects of this law change on the divorce rate. The empirical model we estimate is

$$Y_{it} = \beta_0 Y_{i,t-1} + \sum_{k=1}^8 \beta_k X_{k,it} + \alpha_i + \gamma_i t + \delta_i t^2 + \mu_t + \lambda_i' f_t + e_{it}, \quad (5.1)$$

where we follow Wolfers (2006) in defining the regressors as bi-annual dummies:

$$\begin{aligned} X_{k,it} &= \mathbb{1}\{D_i + 2(k-1) \leq t \leq D_i + 2k - 1\}, \quad \text{for } k = 1, \dots, 7, \\ X_{8,it} &= \mathbb{1}\{D_i + 2(k-1) \leq t\}. \end{aligned}$$

The dummy variable and quadratic trend specification $\alpha_i + \gamma_i t + \delta_i t^2 + \mu_t$ is also used in Friedberg (1998) and Wolfers (2006). The additional interactive fixed effects $\lambda_i' f_t$ were added in Kim and Oka (2014) to control for additional unobserved heterogeneity in the divorce rate, e.g. due to social, cultural or demographic factors. We extend the specification further by adding a lagged dependent variable $Y_{i,t-1}$ to control for state dependence of the divorce rate, but we also report results without $Y_{i,t-1}$ below. We use the dataset of Kim and Oka (2014),³¹ which is a balanced panel of $N = 48$ states over $T = 33$ years, leaving $T = 32$ time periods if the lagged dependent variable is included.

For estimation we first eliminate α_i , γ_i , δ_i and μ_t from the model by projecting the outcome variable and all regressors accordingly, e.g. $\tilde{Y} = M_{1_N} Y M_{(1_T, \mathbf{t}, \mathbf{t}^2)}$, where 1_N and 1_T are N - and T -vectors, respectively, with all entries equal to one, and \mathbf{t} and \mathbf{t}^2 are T -vectors with entries t and t^2 , respectively. The model after projection reads $\tilde{Y}_{it} = \beta_0 \tilde{Y}_{i,t-1} + \sum_{k=1}^8 \beta_k \tilde{X}_{k,it} + \tilde{\lambda}_i' \tilde{f}_t + \tilde{e}_{it}$, which is exactly the model we have studied so far in

³¹The data is available from <http://qed.econ.queensu.ca/jae/2014-v29.2/kim-oka/>

this paper.³² We will use the LS estimator described above to estimate this model. The projection reduces the effective sample size to $N = 48 - 1 = 47$ and $T = 32 - 3 = 29$, which should be accounted for when calculating standard errors, e.g. in the formula for $\hat{\sigma}_R^2$ above (degree of freedom correction). Our theoretical results are still applicable.³³

We need to decide on a number of factors R when implementing the LS estimator. As already mentioned in the last remark in Section 4.1 above, we can apply known techniques from the literature on factor models without regressors to obtain a consistent estimator of R^0 . To do so we choose a maximum number of factors of $R_{\max} = 9$ to obtain the preliminary estimate $\hat{\beta}_{R_{\max}}$ and then calculate the residuals $\hat{u}_{it} = \tilde{Y}_{it} - \hat{\beta}_{R_{\max},0} \tilde{Y}_{i,t-1} - \sum_{k=1}^8 \beta_{R_{\max},k} \tilde{X}_{k,it}$. We then apply the IC, PC and BIC3 criteria of Bai and Ng (2002),³⁴ the criterion described in Onatski (2010), and the ER and GR criteria of Ahn and Horenstein (2013) to \hat{u} .³⁵ Most of these criteria also require specification of R_{\max} , and we continue to use $R_{\max} = 9$. The corresponding estimation results for R are presented in Table 1. In addition, we also report the log scree plot, i.e. the sorted eigenvalues of $\hat{u}'\hat{u}$ in Figure 1.

The log scree plot already shows that it is not obvious how to decompose the eigenvalue spectrum into a few larger eigenvalues stemming from factors and the remaining smaller eigenvalues stemming from the idiosyncratic error term.³⁶ This problem is also reflected in the very different estimates for R that one obtains from the various criteria. It might appear that IC1, IC3, PC1, PC2 and PC3 all agree on $\hat{R} = 9$, but this is simply $\hat{R} = R_{\max}$, and if we choose $R_{\max} = 10$, then all these criteria deliver $\hat{R} = 10$, so this should not be considered a reliable estimate.

On the other hand, our asymptotic theory suggests, that the exact choice of R in the estimation of $\hat{\beta}_R$ should not matter too much, as long as R is chosen large enough to cover all relevant factors. Table 2 contains the estimation results for the bias corrected $\hat{\beta}_R$ for $R \in \{0, 1, \dots, 9\}$. Table 3 contains estimates if the lagged dependent variable is not included into the model.³⁷ For all reported estimates we perform bias correction and standard error estimation as described in Bai (2009a) and Moon and Weidner (2013).³⁸

³²To construct $\tilde{Y}_{i,t-1}$ we first apply the lag-operator and then apply the projections M_{1N} and $M_{(1T,t,t^2)}$.

³³If e_{it} is iid normal, then \tilde{e}_{it} is not, but one can apply appropriate orthogonal rotations in N - and T -space such that \tilde{e}_{it} becomes iid normal again, although with sample size reduced to $N = 47$ and $T = 29$. The rotation has no effect on the LS estimator, i.e. it does not matter whether we work in the original or the rotated frame.

³⁴Following Onatski (2010) and Ahn and Horenstein (2013) we report only BIC3 among the AIC and BIC criteria of Bai and Ng (2002).

³⁵To include $R = 0$ as a possible outcome for the Ahn and Horenstein (2013) criterion, we use the mock eigenvalue used in their simulations.

³⁶The first largest eigenvalue is 2.2 times larger than the second eigenvalue, the second is 1.6 times larger than the third, the third is 1.9 times larger than fourth. So the largest view eigenvalues are larger than the remaining ones, and the strong factor assumption might not be completely inappropriate here. However, deciding on a cutoff between factor and non-factor eigenvalues is difficult.

³⁷The result for $R = 7$ in Table 3 should be equal to column (6) in Table III of Kim and Oka (2014). The discrepancy is explained by a coding error in their bias computation. Note also that the result for $R = 0$ in Table 3 does not match the one in Wolfers (2006), because he uses WLS with state population weights, while we use OLS for simplicity. Kim and Oka (2014) estimate both WLS and OLS and find that the difference between the resulting estimates becomes insignificant, once a sufficient number of interactive fixed effects is controlled for.

³⁸We correct for the biases due to heteroscedasticity in both panel dimensions worked out in Bai (2009a), as well as for the dynamic bias worked out in Moon and Weidner (2013). For the latter we use the formula for $\hat{B}_{R,k}$ above, with bandwidth $M = 2$. For the standard error estimation we allow for heteroscedasticity in both panel dimensions, also following Bai (2009a) and Moon and Weidner (2013). The bias and standard error formulas in

Criterion: \hat{R}	Criterion: \hat{R}	Criterion: \hat{R}
IC1: 9	PC1: 9	Onatski: 1
IC2: 7	PC2: 9	ER: 1
IC3: 9	PC3: 9	GR: 3
BIC3: 6		

Table 1: Estimated number of factors in the residuals \hat{u} , using different criteria for estimation and $R_{\max} = 9$. The IC, PC and BIC criteria are described in Bai and Ng (2002), the ER and GR criteria are from Ahn and Horenstein (2013), and we also use the criterion of Onatski (2010).

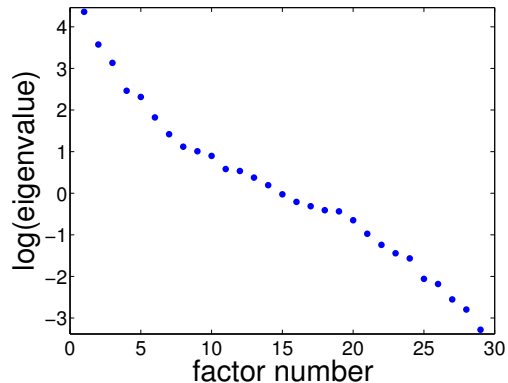


Figure 1: Log scree plot. The natural logarithm of the sorted eigenvalues (corresponding to the principal components, or factors) of $\hat{u}'\hat{u}$ are plotted.

	R = 0	R = 1	R = 2	R = 3	R = 4	R = 5	R = 6	R = 7	R = 8	R = 9
lagged Y	0.432** (4.84)	0.623** (15.38)	0.573** (13.81)	0.411** (8.69)	0.369** (8.19)	0.191** (4.21)	0.137** (2.93)	0.154** (3.24)	0.063 (1.31)	-0.026 (-0.53)
years 1-2	0.043 (0.48)	0.089 (1.79)	0.098 (1.93)	0.105 (1.80)	0.112 (1.90)	0.043 (0.70)	0.087 (1.45)	0.064 (1.08)	0.089 (1.50)	0.039 (0.68)
years 3-4	0.016 (0.18)	0.116* (2.15)	0.147** (2.83)	0.214** (3.31)	0.242** (3.47)	0.170* (2.21)	0.206* (2.53)	0.162* (1.98)	0.204* (2.41)	0.149 (1.61)
years 5-6	-0.040 (-0.41)	0.058 (0.82)	0.102 (1.53)	0.165* (2.01)	0.183* (1.99)	0.115 (1.19)	0.179 (1.84)	0.125 (1.30)	0.148 (1.49)	0.221* (2.00)
years 7-8	-0.010 (-0.08)	0.072 (0.80)	0.114 (1.19)	0.190 (1.64)	0.177 (1.46)	0.140 (1.16)	0.163 (1.34)	0.082 (0.67)	0.093 (0.73)	0.153 (1.15)
years 9-10	-0.126 (-0.84)	0.043 (0.40)	0.041 (0.37)	0.112 (0.86)	0.119 (0.87)	0.013 (0.09)	0.048 (0.34)	-0.032 (-0.23)	0.011 (0.08)	0.054 (0.36)
years 11-12	-0.122 (-0.71)	0.088 (0.70)	0.062 (0.48)	0.122 (0.81)	0.109 (0.69)	0.000 (0.00)	0.042 (0.25)	-0.018 (-0.11)	-0.015 (-0.09)	0.025 (0.14)
years 12-14	-0.122 (-0.59)	0.163 (1.09)	0.097 (0.64)	0.143 (0.83)	0.109 (0.61)	-0.029 (-0.15)	0.017 (0.08)	-0.032 (-0.17)	-0.045 (-0.24)	-0.040 (-0.21)
years 15+	-0.004 (-0.02)	0.301 (1.59)	0.216 (1.15)	0.272 (1.33)	0.232 (1.09)	0.102 (0.46)	0.130 (0.56)	0.081 (0.37)	0.042 (0.19)	0.028 (0.13)

Table 2: Dynamic effects of divorce law reform. We report bias corrected LS-estimates for the regression coefficients in model (5.1). Each column corresponds to a different number of factors $R \in \{0, 1, \dots, 9\}$ used in the estimation. t-values are reported in parenthesis.

When ignoring the lagged dependent variable coefficient, one finds that in both Table 2 and Table 3 the estimation results for $\hat{\beta}_R$ and the corresponding t-values are quite sensitive to changes in R for very small values of R , but become much more stable as R increases, and actually do not change too much anymore from roughly $R = 2$ onwards. These findings are very well in line with our asymptotic theory, and the dynamic effect of divorce law reform that we find are also similar to the findings in Wolfers (2006) and Kim and Oka (2014). The effect of the law reform on the divorce rates initially increases over time, is certainly significant in year 3-4 after the reform, and declines and becomes insignificant

those paper assume $R = R^0$ known, but we strongly expect that those formulas are robust towards $R > R^0$, as partly justified by Theorem 3.2 above. For the model without lagged dependent variable we also allow for serial correlation in e_{it} when estimating the bias and standard deviation of $\hat{\beta}_R$.

	R = 0	R = 1	R = 2	R = 3	R = 4	R = 5	R = 6	R = 7	R = 8	R = 9
years 1-2	0.023 (0.27)	0.034 (0.54)	0.048 (0.70)	0.102 (1.63)	0.053 (0.86)	0.042 (0.66)	0.088 (1.48)	0.095 (1.57)	0.071 (1.21)	0.107 (1.70)
years 3-4	0.049 (0.58)	0.146* (2.12)	0.155* (2.05)	0.265** (3.51)	0.221** (2.95)	0.186* (2.37)	0.223** (2.81)	0.251** (3.09)	0.210* (2.57)	0.228** (2.70)
years 5-6	-0.055 (-0.51)	0.058 (0.67)	0.045 (0.46)	0.201* (1.97)	0.154 (1.59)	0.106 (1.08)	0.207* (2.22)	0.215* (2.23)	0.175 (1.84)	0.204* (2.13)
years 7-8	-0.024 (-0.18)	0.044 (0.39)	-0.011 (-0.09)	0.192 (1.37)	0.136 (1.03)	0.113 (0.92)	0.190 (1.59)	0.212 (1.78)	0.149 (1.25)	0.159 (1.30)
years 9-10	-0.148 (-0.93)	-0.041 (-0.31)	-0.151 (-0.99)	0.044 (0.27)	-0.023 (-0.15)	-0.050 (-0.35)	0.070 (0.49)	0.093 (0.64)	0.018 (0.13)	0.056 (0.40)
years 11-12	-0.195 (-1.10)	-0.029 (-0.19)	-0.195 (-1.13)	-0.011 (-0.06)	-0.079 (-0.46)	-0.109 (-0.66)	0.045 (0.27)	0.071 (0.42)	0.020 (0.12)	0.030 (0.19)
years 12-14	-0.191 (-0.91)	0.043 (0.23)	-0.183 (-0.92)	-0.043 (-0.21)	-0.135 (-0.70)	-0.159 (-0.85)	0.012 (0.06)	0.032 (0.16)	-0.004 (-0.02)	-0.001 (-0.01)
years 15+	-0.007 (-0.03)	0.284 (1.23)	-0.004 (-0.02)	0.094 (0.41)	-0.005 (-0.02)	-0.019 (-0.09)	0.125 (0.54)	0.152 (0.65)	0.112 (0.50)	0.065 (0.29)

Table 3: Same as Table 2, but without including the lagged dependent variable into the model.

afterwards.³⁹

In contrast, the estimated coefficient on the lagged dependent variable in Table 2 is quite large and highly significant for small values of R , but decreases steadily with R , until it gets close to zero and insignificant for $R \geq 8$. A plausible interpretation of this finding is that the model that includes the lagged dependent variable is misspecified, and that the estimated value of β_0 for small values of R does not correspond to a true state dependence of Y_{it} , but simply reflects the time-serial correlation of the error process being picked up by the autoregressive model. This interpretation also matches the fact that once we include more and more factors into the model we control for more and more serial dependence of the unobserved error term, thus uncovering the true insignificance of β_0 in the estimates for $R \geq 8$.

This empirical example shows that instead of relying on a single estimate \hat{R} for the number of factors and reporting the corresponding $\hat{\beta}_{\hat{R}}$ it can be very informative to calculate $\hat{\beta}_R$ for multiple values of R . Whether the estimated coefficients become stable for sufficiently large R values, as our asymptotic theory suggests, is a useful robustness check for the model. When reporting the final results, then, it is better, within a reasonable range, to choose an R that is too large than one that is too small.

6 Monte Carlo Simulations

In this section we investigate the finite sample properties of $\hat{\beta}_R$ through a small scale Monte Carlo simulation. The model is a static panel model with one regressor ($K = 1$),

³⁹The magnitude of estimates is smaller than those estimated by Wolfers (2006), i.e. controlling for unobserved factors reduced the effect size, as already pointed out by Kim and Oka (2014).

R	N=100							
	T=10		T=30		T=100		T=300	
	Bias	SD	Bias	SD	Bias	SD	Bias	SD
0	0.2286	0.0321	0.2301	0.0167	0.2305	0.0117	0.2305	0.0103
1	0.1061	0.0552	0.1155	0.0296	0.1191	0.0195	0.1200	0.0160
2	-0.0385	0.0342	-0.0166	0.0142	-0.0053	0.0071	-0.0019	0.0040
3	-0.0427	0.0342	-0.0170	0.0142	-0.0053	0.0072	-0.0019	0.0040
4	-0.0450	0.0356	-0.0172	0.0144	-0.0053	0.0072	-0.0019	0.0040
5	-0.0461	0.0370	-0.0175	0.0146	-0.0053	0.0073	-0.0019	0.0041

R	N=300							
	T=10		T=30		T=100		T=300	
	Bias	SD	Bias	SD	Bias	SD	Bias	SD
0	0.2291	0.0298	0.2299	0.0136	0.2306	0.0082	0.2307	0.0065
1	0.1054	0.0500	0.1159	0.0263	0.1193	0.0148	0.1203	0.0105
2	-0.0408	0.0237	-0.0172	0.0085	-0.0054	0.0041	-0.0018	0.0023
3	-0.0442	0.0244	-0.0175	0.0086	-0.0054	0.0041	-0.0018	0.0023
4	-0.0462	0.0258	-0.0179	0.0087	-0.0055	0.0041	-0.0018	0.0023
5	-0.0468	0.0275	-0.0182	0.0088	-0.0055	0.0041	-0.0018	0.0023

Table 4: For different combinations of sample sizes N and T we report the bias and standard deviation of the estimator $\hat{\beta}_R$, for $R = 0, 1, \dots, 5$, based on simulations with 10,000 repetition of design (6.1), where the true number of factors is $R^0 = 2$.

two factors ($R^0 = 2$), and the following data generating process (DGP):

$$\begin{aligned}
Y_{it} &= \beta^0 X_{it} + \sum_{r=1}^2 \lambda_{ir} f_{tr} + e_{it}, \\
X_{it} &= 1 + \tilde{X}_{it} + \sum_{r=1}^2 (\lambda_{ir} + \chi_{ir})(f_{tr} + f_{t-1,r}), \\
e_{it} &= \frac{1}{\sqrt{2}}(v_{it} + v_{i,t-1}).
\end{aligned} \tag{6.1}$$

The random variables \tilde{X}_{it} , λ_{ir} , f_{tr} , χ_{ir} and v_{it} are mutually independent; with \tilde{X}_{it} and $f_{tr} \sim iid\mathcal{N}(0, 1)$; λ_{ir} and $\chi_{ir} \sim iid\mathcal{N}(1, 1)$; and $v_{it} \sim iid t(5)$, i.e. v_{it} has a Student's t -distribution with 5 degrees of freedom.

Note that this model satisfies Assumptions SF, NC, and LL(i), but not LL(ii). The error term e_{it} is *not* distributed as *iid* normal. The time series of e_{it} follows an MA(1) process with innovations distributed as $t(5)$. The purpose of this design is to demonstrate that the iid normality restriction on e_{it} in Assumption LL(ii) is a technical assumption as mentioned in Section 2 and may be relaxed.

We choose $\beta^0 = 1$, and use 10,000 repetitions in our simulation. The true number of factors is chosen to be $R^0 = 2$. For each draw of Y and X we compute the LS estimator $\hat{\beta}_R$ according to equation (3.1) for different values of R , namely $R \in \{0, 1, 2, 3, 4, 5\}$.

Table 4 reports bias and standard deviation of the estimator $\hat{\beta}_R$ for different combinations of R , N and T . For $R < R^0 = 2$ the model is misspecified and $\hat{\beta}_R$ turns out to be severely biased. There is also bias in $\hat{\beta}_R$ for $R \geq R^0$, due to time-serial correlation of e_{it} . This bias was worked out in Bai (2009a), and bias correction is also discussed there.

R	N=100				T=100				
	2.5%	5%	10%	25%	50%	75%	90%	95%	97.5%
2	-1.95	-1.70	-1.44	-1.01	-0.52	-0.05	0.36	0.61	0.86
3	-1.94	-1.73	-1.47	-1.01	-0.52	-0.05	0.38	0.64	0.87
4	-1.97	-1.73	-1.47	-1.01	-0.52	-0.04	0.39	0.64	0.85
5	-1.97	-1.74	-1.48	-1.02	-0.52	-0.04	0.39	0.64	0.88

R	N=300				T=300				
	2.5%	5%	10%	25%	50%	75%	90%	95%	97.5%
2	-1.91	-1.68	-1.43	-1.00	-0.54	-0.07	0.33	0.57	0.78
3	-1.91	-1.68	-1.44	-1.00	-0.55	-0.08	0.32	0.57	0.78
4	-1.92	-1.68	-1.44	-1.00	-0.55	-0.08	0.33	0.57	0.78
5	-1.91	-1.68	-1.44	-1.00	-0.54	-0.07	0.34	0.58	0.79

Table 5: Quantiles of the distribution of $\sqrt{NT}(\hat{\beta}_R - \beta^0)$ are reported for $N = T = 100$ and $N = T = 300$, with $R = 2, 3, 4, 5$, based on simulations with 10,000 repetition of design (6.1), where the true number of factors is $R^0 = 2$.

R	N=100				N=300			
	T=10	T=30	T=100	T=300	T=10	T=30	T=100	T=300
2	0.252	0.084	0.057	0.051	0.535	0.146	0.055	0.051
3	0.327	0.111	0.062	0.050	0.643	0.209	0.062	0.056
4	0.358	0.141	0.067	0.054	0.672	0.280	0.070	0.057
5	0.349	0.170	0.074	0.056	0.664	0.348	0.078	0.058

Table 6: The empirical size of a t-test with 5% nominal size is reported for different combinations of N , T and R , based on 10,000 repetition of design (6.1). A bias corrected estimator for β is used to calculate the test statistics, and we allow for heteroscedasticity and time-serial correlation when estimating bias and standard deviation. Results for $R = 0, 1$ are not reported since those have size=1 due to misspecification.

We have purposefully chosen a DGP where $\hat{\beta}_R$ exhibits such a bias to illustrate that all features of the asymptotic distribution of $\hat{\beta}_{R^0}$ are replicated by $\hat{\beta}_R$, $R > R^0$, including the bias.

Table 5 reports various quantiles of the distribution of $\sqrt{NT}(\hat{\beta}_R - \beta^0)$ for $N = T = 100$ and $N = T = 300$, and different values of $R \geq R^0$. From these tables, we see that as N, T increases the distribution of $\hat{\beta}_R$ gets closer to that of $\hat{\beta}_{R^0}$.

Table 6 reports the size of a t-test with nominal size equal to 5% for $R \geq R^0$. We use the results in Bai (2009a) to correct for the leading $1/N$ (not actually present in our DGP) and $1/T$ (present in our DGP) biases in $\hat{\beta}_R$ before calculating the t-test statistics, allowing for heteroscedasticity in both panel dimensions and for time-serial correlation when estimating the bias and standard deviation of $\hat{\beta}_R$. The finite sample size distortions are mostly due to residual bias after bias correction, but also partly due to some finite sample downward bias in the standard error estimates. The size distortions increase with R , but for all values of $R \geq R^0$ in Table 6 the size distortions decrease rapidly as T increases.

Monte Carlo Simulation results for an AR(1) model with factors can be found in Section S.7 of the supplementary material. Those additional simulations show that the finite sample properties (e.g. for $T = 30$) of $\hat{\beta}_{R^0}$ and $\hat{\beta}_R$, $R > R^0$, can be quite different, but those differences vanish as T becomes large, as predicted by our asymptotic theory. In general, we always expect some finite sample inefficiency from overestimating the number of factors.

7 Conclusions

We show that under certain assumptions the limiting distribution of the LS estimator of a linear panel regression with interactive fixed effects does not change when we include redundant factors in the estimation. The implication of this is that one can use an upper bound of the number of factors R in the estimation without asymptotic efficiency loss. However, some finite sample efficiency loss from overestimating R is likely, so that R should not be chosen too large in actual applications. We impose *iid* normality of the regression errors to derive the asymptotic result, because we require certain results on the eigenvalues and eigenvectors of random covariance matrices that are only known in that case. We expect that progress in the literature on large dimensional random covariance matrices will allow verification of our high-level assumptions under more general error distributions, and our simulation results suggest that the result also holds for non-normal and correlated errors. We also provide multiple intermediate asymptotic results under more general conditions.

References

- Ahn, S. C. and Horenstein, A. R. (2013). Eigenvalue ratio test for the number of factors. *Econometrica*, 81(3):1203–1227.
- Ahn, S. C., Lee, Y. H., and Schmidt, P. (2001). GMM estimation of linear panel data models with time-varying individual effects. *Journal of Econometrics*, 101(2):219–255.
- Ahn, S. C., Lee, Y. H., and Schmidt, P. (2013). Panel data models with multiple time-varying individual effects. *Journal of Econometrics*, 174(1):1–14.
- Allen, D. W. (1992). Marriage and divorce: Comment. *American Economic Review*, pages 679–685.
- Andrews, D. W. K. (1999). Estimation when a parameter is on a boundary. *Econometrica*, 67(6):1341–1384.
- Bai, J. (2009a). Panel data models with interactive fixed effects. *Econometrica*, 77(4):1229–1279.
- Bai, J. (2009b). Supplement to “Panel data models with interactive fixed effects”: technical details and proofs. *Econometrica Supplemental Material*, 77(4).
- Bai, J. (2013). Likelihood approach to small T dynamic panel models with interactive effects. Manuscript.
- Bai, J. and Ng, S. (2002). Determining the number of factors in approximate factor models. *Econometrica*, 70(1):191–221.
- Bai, Z. (1993). Convergence rate of expected spectral distributions of large random matrices. Part II. Sample covariance matrices. *The Annals of Probability*, 21(2):649–672.
- Bai, Z. (1999). Methodologies in spectral analysis of large dimensional random matrices, a review. *Statistica Sinica*, 9:611–677.
- Bai, Z., Miao, B., and Yao, J. (2004). Convergence rates of spectral distributions of large sample covariance matrices. *SIAM journal on matrix analysis and applications*, 25(1):105–127.
- Bai, Z. D., Silverman, J. W., and Yin, Y. Q. (1988). A note on the largest eigenvalue of a large dimensional sample covariance matrix. *J. Multivar. Anal.*, 26(2):166–168.
- Chudik, A., Pesaran, M. H., and Tosetti, E. (2011). Weak and strong cross-section dependence and estimation of large panels. *The Econometrics Journal*, 14(1):C45–C90.

- Cox, D. D. and Kim, T. Y. (1995). Moment bounds for mixing random variables useful in nonparametric function estimation. *Stochastic processes and their applications*, 56(1):151–158.
- Dhaene, G. and Jochmans, K. (2010). Split-panel jackknife estimation of fixed-effect models. *Unpublished manuscript*.
- Friedberg, L. (1998). Did unilateral divorce raise divorce rates? evidence from panel data. Technical report, JSTOR.
- Geman, S. (1980). A limit theorem for the norm of random matrices. *Annals of Probability*, 8(2):252–261.
- Götze, F. and Tikhomirov, A. (2010). The Rate of Convergence of Spectra of Sample Covariance Matrices. *Theory of Probability and its Applications*, 54:129.
- Gray, J. S. (1998). Divorce-law changes, household bargaining, and married women’s labor supply. *American Economic Review*, pages 628–642.
- Holtz-Eakin, D., Newey, W., and Rosen, H. S. (1988). Estimating vector autoregressions with panel data. *Econometrica*, 56(6):1371–95.
- Johnstone, I. (2001). On the distribution of the largest eigenvalue in principal components analysis. *Annals of Statistics*, 29(2):295–327.
- Kato, T. (1980). *Perturbation Theory for Linear Operators*. Springer-Verlag.
- Kiefer, N. (1980). A time series-cross section model with fixed effects with an intertemporal factor structure. *Unpublished manuscript, Department of Economics, Cornell University*.
- Kim, D. and Oka, T. (2014). Divorce law reforms and divorce rates in the usa: An interactive fixed-effects approach. *Journal of Applied Econometrics*.
- Latala, R. (2005). Some estimates of norms of random matrices. *Proc. Amer. Math. Soc.*, 133:1273–1282.
- Lu, X. and Su, L. (2013). Shrinkage estimation of dynamic panel data models with interactive fixed effects. Technical report, Working paper, Hong Kong University of Science & Technology.
- Marčenko, V. and Pastur, L. (1967). Distribution of eigenvalues for some sets of random matrices. *Sbornik: Mathematics*, 1(4):457–483.
- Moon, H., Shum, M., and Weidner, M. (2014). Interactive fixed effects in the blp random coefficients demand model. *CeMMAP working paper series*.
- Moon, H. and Weidner, M. (2013). Dynamic Linear Panel Regression Models with Interactive Fixed Effects. *CeMMAP working paper series*.
- Nickell, S. (1981). Biases in dynamic models with fixed effects. *Econometrica*, 49(6):1417–1426.
- Onatski, A. (2010). Determining the number of factors from empirical distribution of eigenvalues. *The Review of Economics and Statistics*, 92(4):1004–1016.
- Onatski, A. (2012). Asymptotics of the principal components estimator of large factor models with weakly influential factors. *Journal of Econometrics*, 168(2):244–258.
- Onatski, A. (2013). Asymptotic Analysis of the Squared Estimation Error in Misspecified Factor Models. Manuscript.
- Pesaran, M. H. (2006). Estimation and inference in large heterogeneous panels with a multifactor error structure. *Econometrica*, 74(4):967–1012.
- Peters, H. E. (1986). Marriage and divorce: Informational constraints and private contracting. *American Economic Review*, 76(3):437–54.
- Peters, H. E. (1992). Marriage and divorce: Reply. *American Economic Review*, pages 686–693.

- Silverstein, J. (1990). Weak convergence of random functions defined by the eigenvectors of sample covariance matrices. *The Annals of Probability*, 18(3):1174–1194.
- Silverstein, J. W. (1989). On the eigenvectors of large dimensional sample covariance matrices. *J. Multivar. Anal.*, 30(1):1–16.
- Soshnikov, A. (2002). A note on universality of the distribution of the largest eigenvalues in certain sample covariance matrices. *Journal of Statistical Physics*, 108(5):1033–1056.
- Stock, J. H. and Watson, M. W. (2002). Forecasting using principal components from a large number of predictors. *Journal of the American Statistical Association*, 97:1167–1179.
- Vershynin, R. (2010). Introduction to the non-asymptotic analysis of random matrices. *arXiv preprint arXiv:1011.3027*.
- Wolfers, J. (2006). Did unilateral divorce laws raise divorce rates? a reconciliation and new results. *American Economic Review*, pages 1802–1820.
- Yin, Y. Q., Bai, Z. D., and Krishnaiah, P. (1988). On the limit of the largest eigenvalue of the large-dimensional sample covariance matrix. *Probability Theory Related Fields*, 78:509–521.
- Zaffaroni, P. (2009). Generalized least squares estimation of panel with common shocks. Manuscript.

A Appendix

A.1 Spectral Norm of Random Matrices

Consider an $N \times T$ matrix u whose entries u_{it} have uniformly bounded second moments. Then we have $\|u\| \leq \|u\|_{HS} = \sqrt{\sum_{i,t} u_{it}^2} = \mathcal{O}_P(\sqrt{NT})$. However, in Assumption LL(*i.b*) and Assumption DX-1(*i*) and Assumption DX-2(*i*) we impose $\|\tilde{X}_k^{\text{str}}\| = \mathcal{O}_P(N^{3/4})$ and $\|\tilde{X}_k\| = \mathcal{O}_P(N^{3/4})$, respectively, as N and T grow at the same rate, and in Assumption SN(*ii*) we impose $\|e\| = \mathcal{O}_P(\sqrt{\max(N, T)})$ under an arbitrary asymptotic $N, T \rightarrow \infty$. Those smaller asymptotic rates for the spectral norms of \tilde{X}_k^{str} , \tilde{X}_k and e can be justified by firstly assuming that the entries of these matrices are mean zero and have certain bounded moments, and secondly imposing weak cross-sectional and time-serial correlation. The purpose of this appendix section is to provide some examples of matrix distributions that make the last statement more precise. We consider the $N \times T$ matrix u , which can represent either e , \tilde{X}_k^{str} or \tilde{X}_k .

Example 1: If we assume that $\mathbb{E}u_{it} = 0$, that $\mathbb{E}u_{it}^4$ is uniformly bounded, and that the u_{it} are independently distributed across i and over t , then the results in Latala (2005) show that $\|u\| = \mathcal{O}_P(\sqrt{\max(N, T)})$.

Example 2: Onatski (2013) provides the following example, which allows for both cross-sectional and time-serial dependence: Let ϵ be an $N \times T$ matrix with mean zero, independent entries that have uniformly bounded fourth moment, let ϵ_t denote the columns of ϵ , and also define past ϵ_t , $t \leq 0$, satisfying the same distributional assumptions. Let $u_t = \sum_{j=0}^m \Psi_{N,j} \epsilon_{t-j}$, where m is a fixed integer, and $\Psi_{N,j}$ are $N \times N$ matrices such that $\max_j \|\Psi_{N,j}\|$ is uniformly bounded. Then, the $N \times T$ matrix u with columns u_t satisfies $\|u\| = \mathcal{O}_P(\sqrt{\max(N, T)})$.

More examples of matrix distributions that satisfy $\|u\| = \mathcal{O}_P(\sqrt{\max(N, T)})$ are discussed in

Onatski (2013) and Moon and Weidner (2013). Theorem 5.48 and Remark 5.49 in Vershynin (2010) can also be used to obtain a slightly weaker bound on $\|u\|$ under very general correlation of u in one of its dimensions.

Note that the random matrix theory literature often only discusses asymptotics where N and T grow at the same rate and shows $\|u\| = \mathcal{O}_P(\sqrt{N})$ under that asymptotic. Those results can easily be extended to more general asymptotics with $N, T \rightarrow \infty$ by considering u as a submatrix of a $\max(N, T) \times \max(N, T)$ matrix u^{big} , and using that $\|u\| \leq \|u^{\text{big}}\|$.

Example 3: The following Lemma provides a justification for the bounds on $\|\tilde{X}_k^{\text{str}}\|$ and $\|\tilde{X}_k\|$, allowing for a quite general type of correlation in both panel dimensions.

Lemma A.1. *Let u be an $N \times T$ matrix with entries u_{it} . Let $\Sigma_{ij} = \frac{1}{T} \sum_{t=1}^T \mathbb{E}(u_{it}u_{jt})$, and let Σ be the $N \times N$ matrix with entries Σ_{ij} . Let $\eta_{ij} = \frac{1}{\sqrt{T}} \sum_{t=1}^T [u_{it}u_{jt} - \mathbb{E}(u_{it}u_{jt})]$, $\Psi_{ij} = \frac{1}{N} \sum_{k=1}^N \mathbb{E}(\eta_{ik}\eta_{jk})$, and $\chi_{ij} = \frac{1}{\sqrt{N}} \sum_{k=1}^N [\eta_{ik}\eta_{jk} - \mathbb{E}(\eta_{ik}\eta_{jk})]$. Consider an asymptotic where $N, T \rightarrow \infty$ such that N/T converges to a finite positive constant, and assume that*

$$(i) \quad \|\Sigma\| = \mathcal{O}(1).$$

$$(ii) \quad \frac{1}{N^2} \sum_{i,j=1}^N \mathbb{E}(\eta_{ij}^2) = \mathcal{O}(1).$$

$$(iii) \quad \frac{1}{N} \sum_{i,j=1}^N \Psi_{ij}^2 = \mathcal{O}(1).$$

$$(iv) \quad \frac{1}{N^2} \sum_{i,j=1}^N \mathbb{E}(\chi_{ij}^2) = \mathcal{O}(1).$$

Then we have $\|u\| = \mathcal{O}_P(N^{5/8})$.

The Lemma does not impose $\mathbb{E}u_{it} = 0$ explicitly, but justification of assumption (i) in the lemma usually requires $\mathbb{E}u_{it} = 0$. The assumptions (ii), (iii) and (iv) in the lemma can e.g. be justified by assuming appropriate mixing conditions in both panel dimensions, see e.g. Cox and Kim (1995) for the time-series case.

As pointed out above, our results in Section 4.2 can be obtained under the weaker condition $\|e\| = o_P(N^{2/3})$, and Lemma A.1 can also be applied with $u = e$ then. In that case, the assumptions in Lemma A.1 are not the same, but are similar to those imposed in Bai (2009a).

A.2 Expansion of Objective Function when $R = R^0$

Here we provide a heuristic derivation of the expansion of $\mathcal{L}_{NT}^0(\beta)$ in Theorem 4.2. We expand the profile objective function $\mathcal{L}_{NT}^0(\beta)$ simultaneously in β and in the spectral norm of e . Let the $K + 1$ expansion parameters be defined by $\epsilon_0 = \|e\|/\sqrt{NT}$ and $\epsilon_k = \beta_k^0 - \beta_k$, $k = 1, \dots, K$, and define the $N \times T$ matrix $X_0 = (\sqrt{NT}/\|e\|)e$. With these definitions we obtain

$$\frac{1}{\sqrt{NT}} (Y - \beta \cdot X) = \frac{1}{\sqrt{NT}} [\lambda^0 f^{0'} + (\beta^0 - \beta) \cdot X + e] = \frac{\lambda^0 f^{0'}}{\sqrt{NT}} + \sum_{k=0}^K \epsilon_k \frac{X_k}{\sqrt{NT}}. \quad (\text{A.1})$$

According to equation (3.3) the profile objective function $\mathcal{L}_{NT}^0(\beta)$ can be written as the sum over the $T - R^0$ smallest eigenvalues of the matrix in (A.1) multiplied by its transposed. We consider $\sum_{k=0}^K \epsilon_k X_k/\sqrt{NT}$ as a small perturbation of the unperturbed matrix $\lambda^0 f^{0'}/\sqrt{NT}$, and

thus expand $\mathcal{L}_{NT}^0(\beta)$ in the perturbation parameters $\epsilon = (\epsilon_0, \dots, \epsilon_K)$ around $\epsilon = 0$, namely

$$\mathcal{L}_{NT}^0(\beta) = \frac{1}{NT} \sum_{g=0}^{\infty} \sum_{k_1, \dots, k_g=0}^K \epsilon_{k_1} \epsilon_{k_2} \dots \epsilon_{k_g} L^{(g)}(\lambda^0, f^0, X_{k_1}, X_{k_2}, \dots, X_{k_g}), \quad (\text{A.2})$$

where $L^{(g)} = L^{(g)}(\lambda^0, f^0, X_{k_1}, X_{k_2}, \dots, X_{k_g})$ are the expansion coefficients.

The unperturbed matrix $\lambda^0 f^{0'} / \sqrt{NT}$ has rank R^0 , so that the $T - R^0$ smallest eigenvalues of the unperturbed $T \times T$ matrix $f^0 \lambda^{0'} \lambda^0 f^{0'} / NT$ are all zero, i.e. $\mathcal{L}_{NT}^0(\beta) = 0$ for $\epsilon = 0$ and thus $L^{(0)}(\lambda^0, f^0) = 0$. Due to Assumption SF the R^0 non-zero eigenvalues of the unperturbed $T \times T$ matrix $f^0 \lambda^{0'} \lambda^0 f^{0'} / NT$ converge to positive constants as $N, T \rightarrow \infty$. This means that the ‘‘separating distance’’ of the $T - R^0$ zero-eigenvalues of the unperturbed $T \times T$ matrix $f^0 \lambda^{0'} \lambda^0 f^{0'} / NT$ converges to a positive constant, i.e. the next largest eigenvalue is well separated. This is exactly the technical condition under which the perturbation theory of linear operators guarantees that the above expansion of \mathcal{L}_{NT}^0 in ϵ exists and is convergent as long as the spectral norm of the perturbation $\sum_{k=0}^K \epsilon_k X_k / \sqrt{NT}$ is smaller than a particular convergence radius $r_0(\lambda^0, f^0)$, which is closely related to the separating distance of the zero-eigenvalues. For details on that see Kato (1980) and Section S.2 of the supplementary appendix, where we define $r_0(\lambda^0, f^0)$ and show that it converges to a positive constant as $N, T \rightarrow \infty$. Note that for the expansion (A.2) it is crucial that we have $R = R^0$, since the perturbation theory of linear operators describes the perturbation of the sum of *all* zero-eigenvalues of the unperturbed matrix $f^0 \lambda^{0'} \lambda^0 f^{0'} / NT$. For $R > R^0$ the sum in $\mathcal{L}_{NT}^R(\beta)$ leaves out the $R - R^0$ largest of these perturbed zero-eigenvalues, which results in a much more complicated mathematical problem, since the structure and ranking among these perturbed zero-eigenvalues needs to be discussed.

The above expansion of $\mathcal{L}_{NT}^0(\beta)$ is applicable whenever the operator norm of the perturbation matrix $\sum_{k=0}^K \epsilon_k X_k / \sqrt{NT}$ is smaller than $r_0(\lambda^0, f^0)$. Since our assumptions guarantee that $\|X_k / \sqrt{NT}\| = \mathcal{O}_P(1)$, for $k = 0, \dots, K$, and $\epsilon_0 = \mathcal{O}_P(\min(N, T)^{-1/2}) = o_P(1)$, we have $\left\| \sum_{k=0}^K \epsilon_k X_k / \sqrt{NT} \right\| = \mathcal{O}_P(\|\beta - \beta^0\|) + o_P(1)$, i.e. the above expansion is always applicable asymptotically within a shrinking neighborhood of β^0 — which is sufficient since we already know that $\hat{\beta}_R$ is consistent for $R \geq R^0$.

In addition, to guaranteeing converge of the series expansion, the perturbation theory of linear operators also provides explicit formulas for the expansion coefficients $L^{(g)}$, namely for $g = 1, 2, 3$ we have $L^{(1)}(\lambda^0, f^0, X_k) = 0$, $L^{(2)}(\lambda^0, f^0, X_{k_1}, X_{k_2}) = \text{Tr}(M_{\lambda^0} X_{k_1} M_{f^0} X'_{k_2})$, $L^{(3)}(\lambda^0, f^0, X_{k_1}, X_{k_2}, X_{k_3}) = -\frac{1}{3}[\text{Tr}(M_{\lambda^0} X_{k_1} M_f X'_{k_2} \lambda^0 (\lambda^{0'} \lambda^0)^{-1} (f^{0'} f^0)^{-1} f^{0'} X'_{k_3}) + \dots]$, where the dots refer to 5 additional terms obtained from the first one by permutation of k_1, k_2 and k_3 , so that the expression becomes totally symmetric in these indices. A general expression for the coefficients for all orders in g is given in Lemma S.1 in the appendix. One can show that for $g \geq 3$ the coefficients $L^{(g)}$ are bounded as follows

$$\frac{1}{NT} \left| L^{(g)}(\lambda^0, f^0, X_{k_1}, X_{k_2}, \dots, X_{k_g}) \right| \leq a_{NT} (b_{NT})^g \frac{\|X_{k_1}\|}{\sqrt{NT}} \frac{\|X_{k_2}\|}{\sqrt{NT}} \dots \frac{\|X_{k_g}\|}{\sqrt{NT}}, \quad (\text{A.3})$$

where a_{NT} and b_{NT} are functions of λ^0 and f^0 that converge to finite positive constants in probability. This bound on the coefficients $L^{(g)}$ allows us to derive a bound on the remainder term, when

the profile objective expansion is truncated at a particular order. The expansion can be applied under more general asymptotics, but here we only consider the limit $N, T \rightarrow \infty$ with $N/T \rightarrow \kappa^2$, $0 < \kappa < \infty$, i.e. N and T grow at the same rate. Then, apart from the constant $\mathcal{L}_{NT}^0(\beta^0)$, the relevant coefficients of the expansion, which are not treated as part of the remainder term turn out to be $W_{k_1 k_2} = \frac{1}{NT} L^{(2)}(\lambda^0, f^0, X_{k_1}, X_{k_2})$, $C_k^{(1)} = \frac{1}{\sqrt{NT}} L^{(2)}(\lambda^0, f^0, X_k, e) = \frac{1}{\sqrt{NT}} \text{Tr}(M_{\lambda^0} X_k M_{f^0} e')$, and $C_k^{(2)} = \frac{3}{2\sqrt{NT}} L^{(3)}(\lambda^0, f^0, X_k, e, e)$, which corresponds exactly to the definitions in the main text. From the expansion (A.2) and the bound (A.3) we obtain Theorem 4.2. For a more rigorous derivation we refer to Section S.2 in the supplementary appendix.

A.3 $N^{3/4}$ -Convergence Rate of $\widehat{\beta}_R$ for $R > R^0$

The discussion at the end of Section 3 reveals that showing faster than \sqrt{N} convergence of $\widehat{\beta}_R$ is a very important step on the way to the main result. For purely technical reasons we show $N^{3/4}$ -convergence first, but it will usually be the case that if $\widehat{\beta}_R$ is $N^{3/4}$ -consistent, then it is also \sqrt{NT} -consistent as N and T grow at the same rate. We require one of the following two alternative assumptions.

Assumption DX-1 (Decomposition of X_k and Distribution of e_{it} , Version 1).

- (i) For $k = 1, \dots, K$ we have $X_k = \overline{X}_k + \widetilde{X}_k$, where $\text{rank}(\overline{X}_k)$ is bounded as $N, T \rightarrow \infty$, and $\|\overline{X}_k\| = \mathcal{O}_P(\sqrt{NT})$, and $\|\widetilde{X}_k\| = \mathcal{O}_P(N^{3/4})$.
- (ii) Let u be an $N \times T$ matrix whose elements are distributed as i.i.d. $\mathcal{N}(0, 1)$, independent of λ^0, f^0 and $\overline{X}_k, k = 1, \dots, K$, and let one of the following hold
 - (a) either: $e = \Sigma^{1/2} u$, where Σ is an $N \times N$ covariance matrix, independent of u , which satisfies $\|\Sigma\| = \mathcal{O}_P(1)$. In that case, define g to be an $N \times Q$ matrix, independent of u , for some $Q \leq \sum_{k=1}^K \text{rank}(\overline{X}_k)$, such that $g'g = \mathbb{1}_Q$ and $\text{span}(M_{\lambda^0} \overline{X}_k) \subset \text{span}(g)$ for all $k = 1, \dots, K$.⁴⁰
 - (b) or: $e = u \Sigma^{1/2}$, where Σ is a $T \times T$ covariance matrix, independent of u , which satisfies $\|\Sigma\| = \mathcal{O}_P(1)$. In that case, define g to be a $T \times Q$ matrix, independent of u , for some $Q \leq \sum_{k=1}^K \text{rank}(\overline{X}_k)$, such that $g'g = \mathbb{1}_Q$ and $\text{span}(M_{f^0} \overline{X}_k') \subset \text{span}(g)$ for all $k = 1, \dots, K$.

In addition, we assume that there exist a (potentially random) integer sequence $n = n_{NT} > 0$ with $1/n = \mathcal{O}_P(1/N)$ such that $\mu_n(\Sigma) \geq \|g'\Sigma g\|$. Finally, assume that either $R \geq Q$ or that $g'\Sigma g = \|g'\Sigma g\| \mathbb{1}_Q + \mathcal{O}_P(N^{-1/2})$.

Assumption DX-2 (Decomposition of X_k and Distribution of e_{it} , Version 2).

- (i) For $k = 1, \dots, K$ we have $X_k = \overline{X}_k + \widetilde{X}_k$, such that $M_{\lambda^0} \overline{X}_k M_{f^0} = 0$, and $\|\overline{X}_k\| = \mathcal{O}_P(\sqrt{NT})$, and $\|\widetilde{X}_k\| = \mathcal{O}_P(N^{3/4})$.
- (ii) $\|e\| = \mathcal{O}_P(\sqrt{\max(N, T)})$. (same as Assumption SN(ii))

⁴⁰The column space of g thus contains the column space of all $M_{\lambda^0} \overline{X}_k$. $g'g = \mathbb{1}_Q$ is just a normalization.

Theorem A.2. *Let $R > R^0$. Let Assumptions SF, NC and EX hold, and let either Assumption DX-1 or DX-2 be satisfied. Consider $N, T \rightarrow \infty$ with $N/T \rightarrow \kappa^2$, $0 < \kappa < \infty$. Then we have $N^{3/4}(\widehat{\beta}_R - \beta^0) = \mathcal{O}_P(1)$.*

Remarks

- (i) Assumption SN is not explicitly imposed in Theorem A.2, because it is already implied by both Assumption DX-1 and DX-2, see also Lemma A.4 below.
- (ii) The restrictions that Assumption DX-1 imposes on X_k are weaker than those imposed in Assumption LL above. The regressors are decomposed into a low-rank strictly exogenous part \bar{X}_k and a term \tilde{X}_k , which can be both strictly or weakly exogenous. The spectral norm bound $\|\tilde{X}_k\| = \mathcal{O}_P(N^{3/4})$ is satisfied as long as $\tilde{X}_{k,it}$ is mean zero and weakly correlated across i and over t , see Appendix A.1. We can always write $\bar{X}_k = \ell h'$ for some appropriate $\ell \in \mathbb{R}^{N \times \text{rank}(\bar{X}_k)}$ and $h \in \mathbb{R}^{T \times \text{rank}(\bar{X}_k)}$. Thus, the decomposition $X_k = \bar{X}_k + \tilde{X}_k = \ell h' + \tilde{X}_k$ essentially imposes an approximate factor structure on X_k , with factor part \bar{X}_k and idiosyncratic part \tilde{X}_k . In addition to those conditions we need sufficient variation in X_k , as formalized by the non-collinearity Assumption NC.
- (iii) The restrictions that Assumption DX-1 imposes on e are also weaker than those imposed in Assumption LL above. Normality is imposed, but either cross-sectional correlation and heteroscedasticity (case (a)) or time-serial correlation and heteroscedasticity (case (b)), described by Σ , are still allowed. The condition $\|\Sigma\| = \mathcal{O}_P(1)$ requires the correlation of e_{it} to be weak.⁴¹
- (iv) The additional restrictions on Σ in Assumption DX-1 rule out the type of correlation of the low-rank regressor part \bar{X}_k with the second moment structure of e_{it} that was the key feature of the counter example in Section 4.3 above.⁴² Firstly, the condition $\mu_n(\Sigma) \geq \|g'\Sigma g\|$ guarantees that the eigenvectors corresponding to the largest few eigenvectors of Σ (the eigenvectors ν_r of Σ when normalized satisfy $\mu_r(\Sigma) = \nu_r'\Sigma\nu_r$) are not strongly correlated with g (and thus with \bar{X}_k). Secondly, the condition $g'\Sigma g = \|g'\Sigma g\| \mathbb{1}_Q + \mathcal{O}_P(N^{-1/2})$ guarantees that Σ behaves almost as an identity matrix when projected with g , thus not possessing special structure in the “direction of \bar{X}_k ”. Both of these assumption are obviously satisfied when Σ is proportional to the identity matrix.
- (v) Instead of Assumption DX-1 we can also impose Assumption DX-2 to obtain $N^{3/4}$ -consistency in Theorem A.2. The Assumption on e imposed in Assumption DX-2 is the same as in Assumption SN, and as already discussed above, this assumption is quite weak (see also Appendix A.1). However, Assumption DX-2 imposes a much stronger assumption on the regressors by requiring that $M_{\lambda^0} \bar{X}_k M_{f^0} = 0$. This condition implies that $\bar{X}_k = \lambda^0 h' + \ell f^{0'}$ for some $\ell \in \mathbb{R}^{N \times R^0}$ and $h \in \mathbb{R}^{T \times R^0}$, i.e. the factor structure of the regressors is severely restricted. The AR(1) model discussed in Remark (v) of Section 3 does satisfy $M_{\lambda^0} \bar{X}_k = 0$,

⁴¹A sufficient condition for $\|\Sigma\| = \mathcal{O}_P(1)$ is, for example, $\max_i \sum_j |\Sigma_{ij}| = \mathcal{O}_P(1)$, formulated here for case (a). Note that Σ is symmetric.

⁴²However, in the example in Section 4.3 we have both time-serial and cross-sectional correlation in e_{it} , one of which is already ruled out by Assumption DX-1.

and the same is true for a stationary AR(p) model without additional regressors, i.e. for such AR(p) models with factors we obtain $N^{3/4}$ -consistency of $\widehat{\beta}_R$ without imposing strong assumptions (like normality) of e_{it} . Assumption DX-2(i) is furthermore satisfied if $\overline{X}_k = 0$, i.e. if the regressors $X_k = \widetilde{X}_k$ satisfy $\|X_k\| = \mathcal{O}_P(N^{3/4})$, which is true for zero mean weakly correlated processes (see Appendix A.1).

- (vi) Theorem S.5 in the supplementary material provides an alternative $N^{3/4}$ -consistency result, in which Assumptions DX-1 and DX-2 are replaced by a high-level condition, which is more general, but not easy to verify in terms of low-level assumptions.

A.4 Asymptotic Equivalence of $\widehat{\beta}_{R^0}$ and $\widehat{\beta}_R$ for $R > R^0$

Here, we provide high level conditions on the singular values and singular vectors of the error matrix (or equivalently on the eigenvalues and eigenvectors of the corresponding random covariance matrix). Under those assumptions we then establish the main result of the paper that $\widehat{\beta}_{R^0}$ and $\widehat{\beta}_R$ with $R > R^0$ are asymptotically equivalent, that is, $\sqrt{NT}(\widehat{\beta}_R - \widehat{\beta}_{R^0}) = o_P(1)$.

Assumption EV. (Eigenvalues and Eigenvectors of Random Cov. Matrix) *Let the singular value decomposition of $M_{\lambda^0} e M_{f^0}$ be given by $M_{\lambda^0} e M_{f^0} = \sum_{r=1}^Q \sqrt{\rho_r} v_r w_r'$, where $Q = \min(N, T) - R^0$, and $\sqrt{\rho_r}$ are the singular values, and v_r and w_r are normalized N - and T -vectors, respectively.⁴³ Let $\rho_1 \geq \rho_2 \geq \dots \geq \rho_Q \geq 0$. We assume that there exists a constant $c > 0$ and a series of integers $q_{NT} > R - R^0$ with $q_{NT} = o(N^{1/4})$ such that as $N, T \rightarrow \infty$ we have*

$$\begin{aligned}
(i) \quad & \frac{\rho_{R-R^0}}{N} > c, \text{ wpa1.} \\
(ii) \quad & \frac{1}{q_{NT}} \sum_{r=q_{NT}}^Q \frac{1}{\rho_{R-R^0} - \rho_r} = \mathcal{O}_P(1). \\
(iii) \quad & \max_r \|v_r' e P_{f^0}\| = o_P\left(N^{1/4} q_{NT}^{-1}\right), \quad \max_r \|w_r' e' P_{\lambda^0}\| = o_P\left(N^{1/4} q_{NT}^{-1}\right), \\
& \max_r \|v_r' X_k P_{f^0}\| = o_P\left(N q_{NT}^{-1}\right), \quad \max_r \|w_r' X_k' P_{\lambda^0}\| = o_P\left(N q_{NT}^{-1}\right), \\
& \max_{r,s,k} |v_r' X_k w_s| = o_P\left(N^{1/4} q_{NT}^{-1}\right), \quad \text{where } r, s = 1, \dots, Q, \text{ and } k = 1, \dots, K.
\end{aligned}$$

Theorem A.3. *Let $R > R^0$. Let Assumptions SF, NC, EX, and EV hold, and let either Assumption DX-1 or DX-2 hold, and assume that $C^{(1)} = \mathcal{O}_P(1)$. In the limit $N, T \rightarrow \infty$ with $N/T \rightarrow \kappa^2$, $0 < \kappa < \infty$, we then have*

$$\sqrt{NT} \left(\widehat{\beta}_R - \beta^0 \right) = \sqrt{NT} \left(\widehat{\beta}_{R^0} - \beta^0 \right) + o_P(1) = \mathcal{O}_P(1).$$

Remarks

- (i) Theorem A.3 also holds if we replace the Assumptions EX, DX-1, DX-2 by any other condition that guarantees that Assumption SN holds and that $N^{3/4} \left(\widehat{\beta}_R - \beta^0 \right) = \mathcal{O}_P(1)$.

⁴³Thus, w_r is the normalized eigenvector corresponding to the eigenvalue ρ_r of $M_{f^0} e' M_{\lambda^0} e M_{f^0}$, while v_r is the normalized eigenvector corresponding to the eigenvalue ρ_r of $M_{\lambda^0} e M_{f^0} e' M_{\lambda^0}$. We use a convention where eigenvalues with non-trivial multiplicity appear multiple times in the list of eigenvalues ρ_r , but under standard distributional assumptions on e all eigenvalues are simple with probability one anyways.

- (ii) Consider Assumption EV(*iii*). Since v_r and w_r are the normalized singular vectors of $M_{\lambda^0} e M_{f^0}$ we expect them to be essentially uncorrelated with X_k and $e P_{f^0}$, and therefore we expect $v_r' X_k w_s = \mathcal{O}_P(1)$, $\|v_r' e P_{f^0}\| = \mathcal{O}_P(1)$, $\|w_r' e' P_{\lambda^0}\| = \mathcal{O}_P(1)$. We also expect $\|v_r' X_k P_{f^0}\| = \mathcal{O}_P(\sqrt{T})$ and $\|w_r' X_k' P_{\lambda^0}\| = \mathcal{O}_P(\sqrt{N})$, which is different to the analogous expressions with e , since X_k can be correlated with f^0 and λ^0 . The key to making this discussion rigorous is a good knowledge of the properties of the eigenvectors v_r and w_r . If the entries e_{it} are *iid* normal, then the distribution of v_r and w_r can be characterized as follows: Let \tilde{v} be an N -vector with *iid* $\mathcal{N}(0, 1)$ entries and let \tilde{w} be an T -vector with *iid* $\mathcal{N}(0, 1)$ entries. Then we have $v_r =_d \|M_{\lambda^0} \tilde{v}\|^{-1} M_{\lambda^0} \tilde{v}$ and $w_r =_d \|M_{f^0} \tilde{w}\|^{-1} M_{f^0} \tilde{w}$, see also Lemma S.13 in the supplementary material. Here, $=_d$ refers to “equal in distribution”. Thus, if $R^0 = 0$, then v_r and w_r are distributed as *iid* $\mathcal{N}(0, 1)$ vectors, normalized to satisfy $\|v_r\| = \|w_r\| = 1$. This follows from the rotational invariance of the distribution of e when e_{it} is *iid* normally distributed. Using this characterization of v_r and w_r one can formally show that Assumption EV(*iii*) holds, see Lemma A.4 below. The conjecture in the random matrix theory literature is that the limiting distribution of the eigenvectors of a random covariance matrix is “distribution free”, i.e. is independent of the particular distribution of e_{it} (see, e.g., Silverstein (1990), Bai (1999)). However, we are not aware of a formulation and corresponding proof of this conjecture that is sufficient for our purposes, which is one reason why we have to impose *iid* normality of e_{it} .
- (iii) Assumption EV(*ii*) imposes a condition on the eigenvalues ρ_r of the random covariance matrix $M_{f^0} e' M_{\lambda^0} e M_{f^0}$. Eigenvalues are studied more intensely than eigenvectors in the random matrix theory literature, and it is well-known that the properly normalized empirical distribution of the eigenvalues (the so called empirical spectral distribution) of an *iid* sample covariance matrix converges to the Marčenko-Pastur-law (Marčenko and Pastur (1967)) for asymptotics where N and T grow at the same rate. This means that the sum over the function of the eigenvalues ρ_s in Assumption EV(*ii*) can be approximated by an integral over the Marčenko-Pastur limiting spectral distribution. To bound the asymptotic error of this approximation one needs to know the convergence rate of the empirical spectral distribution to its limit law, which is an ongoing research subject in the literature, e.g. Bai (1993), Bai, Miao and Yao (2004), Götze and Tikhomirov (2010). This literature usually considers either *iid* or *iid* normal distributions of e_{it} .
- (iv) For random covariance matrices from *iid* normal errors, it is known from Johnstone (2001) and Soshnikov (2002) that the properly normalized few largest eigenvalues converge to the Tracy-Widom law.⁴⁴ This result can be used to verify Assumption EV(*i*) in the case of *iid* normal e_{it} .
- (v) Details on how to derive Theorem A.3 are given in Section S.4 of the supplementary material.

The following Lemma provides the connection between Theorem A.3 and our main result Theorem 3.1. The proof is given in the supplementary material.

⁴⁴To our knowledge this result is not established for error distributions that are not normal. Soshnikov (2002) has a result under non-normality but only for asymptotics with $N/T \rightarrow 1$.

Lemma A.4. *Let Assumption LL hold, let $R^0 = \text{rank}(\lambda^0) = \text{rank}(f^0)$, and consider a limit $N, T \rightarrow \infty$ with $N/T \rightarrow \kappa^2$, $0 < \kappa < \infty$. Then Assumptions SN, EX, DX-1 and EV are satisfied, and we have $C^{(1)} = \mathcal{O}_P(1)$.*