

Identifying Latent Structures in Panel Data*

Liangjun Su^a, Zhentao Shi^b, and Peter C. B. Phillips^c

^a *School of Economics, Singapore Management University*

^b *Department of Economics, Chinese University of Hong Kong*

^c *Yale University, University of Auckland*

University of Southampton & Singapore Management University

December 29, 2015

This paper provides a novel mechanism for identifying and estimating latent group structures in panel data using penalized techniques. We consider both linear and nonlinear models where the regression coefficients are heterogeneous across groups but homogeneous within a group and the group membership is unknown. Two approaches are considered – penalized profile likelihood (PPL) estimation for the general nonlinear models without endogenous regressors, and penalized GMM (PGMM) estimation for linear models with endogeneity. In both cases we develop a new variant of Lasso called classifier-Lasso (C-Lasso) that serves to shrink individual coefficients to the unknown group-specific coefficients. C-Lasso achieves simultaneous classification and consistent estimation in a single step and the classification exhibits the desirable property of uniform consistency. For PPL estimation C-Lasso also achieves the oracle property so that group-specific parameter estimators are asymptotically equivalent to infeasible estimators that use individual group identity information. For PGMM estimation the oracle property of C-Lasso is preserved in some special cases. Simulations demonstrate good finite-sample performance of the approach both in classification and estimation. Empirical applications to both linear and nonlinear models are presented.

fi C33, C36, C38, C51

Classification; Cluster analysis; Dynamic panel; Group Lasso; High dimensionality; Nonlinear panel; Oracle property; Panel structure model; Parameter heterogeneity; Penalized least squares; Penalized GMM; Penalized profile likelihood

*The authors thank the Co-editor Elie Tamer and three anonymous referees for many constructive comments on the previous version of the paper. They also thank Stéphane Bonhomme, Xiaohong Chen, Cheng Hsiao, Joon Park, and Yixiao Sun for discussions on the subject matter and comments on the paper. Su acknowledges support from the Singapore Ministry of Education for Academic Research Fund (AcRF) under the Tier-2 grant number MOE2012-T2-2-021. Phillips acknowledges NSF support under Grant Nos. SES-0956687 and SES-1285258. Address Correspondence to: Liangjun Su, School of Economics, Singapore Management University, 90 Stamford Road, Singapore 178903; E-mail: ljsu@smu.edu.sg, Phone: +65 6828 0386.

Panel data are widely used in empirical analysis in many disciplines across the social and medical sciences. Such data usually cover individual units sampled from different backgrounds and with different individual characteristics so that an abiding feature of the data is its heterogeneity, much of which is simply unobserved. Neglecting latent heterogeneity in the data can lead to many difficulties, including inconsistent estimation and misleading inference, as is well explained in the literature (e.g., Hsiao 2014, ch. 6). It is therefore widely acknowledged that an important feature of good empirical modeling is to control for heterogeneity in the data as well as for potential heterogeneity in the response mechanisms that figure within the model. Since heterogeneity is a latent feature of the data and its extent is unknown a priori, respecting the potential influence of heterogeneity on model specification is a serious challenge in empirical research. Even in the simplest linear panel data models the challenge is manifest and clearly stated: do we allow for heterogeneous slope coefficients in regression as well as heterogeneous error variances?

While it may be clearly stated, this challenge to the empirical researcher is by no means easily addressed. While allowing for cross-sectional slope heterogeneity in regression may help to avert misspecification bias, it also sacrifices the power of cross section averaging in the estimation of response patterns that may be common across individuals, or more subtly, certain groups of individuals in the panel. In the absence of prior information on such grouping and with data where every new individual to the panel may bring new idiosyncratic elements to be explained, the challenge is demanding and almost universally relevant.

Traditional panel data models frequently deal with this challenge by avoidance. Complete slope homogeneity is assumed for certain specified common parameters in the panel. Under this assumption, the regression parameters are the same across individuals and unobserved heterogeneity is modeled through individual-specific effects which typically enter the model additively. This approach is an exemplar of a convenient assumption that facilitates estimation and inference. Nevertheless, this assumption has been frequently questioned and rejected in empirical studies; see Hsiao and Tahmiscioglu (1997), Lee, Pesaran, and Smith (1997), Durlauf, Kourtellis, and Minkin (2001), Phillips and Sul (2007a), Browning and Carro (2007, 2010, 2014), and Su and Chen (2013), among others.

Despite general agreement that slope heterogeneity is endemic in empirical work with panels, few methods are available to allow for heterogeneity in the slopes when the extent of the heterogeneity is unknown. Some researchers assume complete slope heterogeneity where regression coefficients are completely different for different individuals; see the survey by Baltagi, Bresson, and Pirotte (2008) and Hsiao and Pesaran (2008). Others consider panel structure models where individuals belong to a number of homogeneous groups within a broadly heterogeneous population, and the regression parameters are the same within each group but differ across groups. Two essential questions remain: how to determine the unknown number of groups (dubbed convergence clubs in the economic growth literature); and how to identify the membership of each individual. These are longstanding questions of statistical classification in panel data. No completely satisfactory solution has yet been found, although various approaches have been adopted in empirical research. For instance, Bester and Hansen (2016) consider a panel structure model where individuals are grouped according to some external classification, geographic location, or observable explanatory variables; Ando and Bai (2014) consider a multifactor asset-pricing model with group-specific pervasive factors where the group membership is known. Here the group structure is assumed to be *completely known* to the researcher, an approach that is common in practical work because of its convenience. In spite of its convenience, this approach to

panel inference is inevitably misleading when the number of groups and individual identities are incorrectly classified.

Several approaches have been proposed to determine an *unknown* group structure in modeling unobserved slope heterogeneity in panels. The first approach applies finite mixture models. For example, Sun (2005) considers a *parametric* finite mixture linear panel data model, and Kasahara and Shimotsu (2009) and Browning and Carro (2011) study identification in discrete choice panel data models for a fixed number of groups using *nonparametric* discrete mixture distributions. The second approach is based on the K-means algorithm in statistical cluster analysis. Lin and Ng (2012) and Sarafidis and Weber (2015) consider linear panel data models where the slope coefficients have latent group structure. They modify the K-means algorithm to estimate the models but do not provide any inference theory. Bonhomme and Manresa (2015, BM hereafter) consider a linear panel data model where the additive fixed effects have group structure and apply the K-means algorithm to estimate the model and study its asymptotic properties. Ando and Bai (2015) extend BM's approach to allow for group structure among the interactive fixed effects. In addition, Phillips and Sul (2007a) develop an algorithm for determining group clusters that relies on the estimation of evaporating trend functions to determine convergence clusters. Hahn and Moon (2010) argue that the group structure has sound foundations in game theory or macroeconomic models where multiplicity of Nash equilibria is expected and they consider nonlinear panel data models where the parameter of interest is common to individuals whereas the fixed effects have finite support.

The present paper proposes a new method for econometric estimation and inference in panel models when the regression parameters are heterogenous across groups, individual group membership is unknown, and classification is to be determined empirically. It is an automated data-determined procedure and does not require the specification of any modeling mechanism for the unknown group structure. The methods proposed here have several novel aspects in relation to earlier research and they contribute to both the Lasso and econometric classification literatures in various ways, which we outline in the following paragraphs.

First, our approach is motivated by a key advantage of Lasso technology in coping with parameter sparsity. This advantage is particularly useful when the set of unknown parameters is potentially large but may also embody certain *sparse* features. In a typical panel structure model, the *effective* number of unknown regression parameters $\{ \beta_i = 1 \}$ is not of order (n) as it would be if these parameters were all incidental, but rather of some order (g_0) where g_0 denotes the number of unknown groups within which the parameters are homogeneous. Hence, in many empirical applications the set of unknown parameters in a panel structure model surely exhibits the desirable sparsity feature, making the use of Lasso technology highly appealing.

Second, the procedures developed in the present paper contribute to the fused Lasso literature in which sparsity arises because some parameters take the same value. The fused Lasso was proposed by Tibshirani, Saunders, Rosset, Zhu, and Knight (2005) and was designed for problems with features that can be ordered

likelihood objective function and when multiple penalty terms are used they enter the objective function *additively*. To achieve simultaneous group classification and estimation in a single step our variant of Lasso involves *additive* penalty terms, each of which takes a *multiplicative* form as a product of g_0 penalty terms. To the best of our knowledge, this paper is the first to propose a mixed additive-multiplicative penalty form that can serve as an engine for simultaneous classification and estimation. The method works by using each of the g_0 penalty terms in the *multiplicative* expression to shrink the individual-level regression parameter vectors to a particular *unknown* group-level parameter vector, thereby producing a joint shrinkage process to unknown quantities. This process is distinct from the prototypical Lasso method that shrinks an individual parameter to the *known* value zero and the group Lasso method that shrinks a parameter vector to a *known* vector of zeros (see Yuan and Lin, 2006). To emphasize its role as a classifier and for future reference, we describe our new Lasso method as the *classifier-Lasso* or *C-Lasso*.

Fourth, we develop a double asymptotic limit theory for the C-Lasso that demonstrates its capacity to achieve simultaneous classification and estimation in a single step. As mentioned in the Abstract, the paper develops two classes of estimators for panel structure models – penalized profile likelihood (PPL) and penalized GMM (PGMM). The former is applicable to both linear and nonlinear panel models without endogeneity and with or without dynamic structures, while the latter is applicable to linear panel models with endogeneity or dynamic structures. Both broaden the scope of applicability of our method as early literature only considers linear panels without endogeneity. In either case, we show uniform classification consistency in the sense that all individuals belonging to a certain group can be classified into the same group correctly uniformly over both individuals and group identities with probability approaching one (w.p.a.1). Conversely, all individuals that are classified into a certain group belong to the same group uniformly over both individuals and group identities w.p.a.1. Such a uniform result allows us to establish an *oracle* property of the PPL estimator that, like the BM K-means estimator, is asymptotically equivalent to the corresponding infeasible estimator of the group-specific parameter that is obtained by knowing all individual group identities. Unfortunately, our PGMM estimator generally does not have the oracle property. But the uniform classification consistency property allows us to develop a limit theory for post-C-Lasso estimators that are obtained by pooling all individuals in an estimated group to estimate the group-specific parameters and these estimators are asymptotically as efficient as the oracle ones in both the PPL and PGMM contexts.

Fifth, C-Lasso enables empirical researchers to study panel structures without *a priori* knowledge of the number of groups, without the need to specify any ancillary regression models to model individual group identities, and with no need to make any distributional assumptions. When the number g_0 of groups is unknown, a BIC-type information criterion is proposed to determine the number of groups for both PPL and PGMM estimation and it is shown that this procedure selects the correct number of groups consistently.

The rest of the paper is organized as follows. We study C-Lasso PPL estimation and inference of panel structure models in Section 2. PGMM estimation and inference is addressed in Section 3. Section 4 reports Monte Carlo simulation findings. Section 5 contains two empirical applications. Section 6 concludes. Proofs of the main results in the paper are given in Appendices A and B. Additional materials may be found in the Supplemental Material.

For any real matrix \mathbf{X} we write the transpose \mathbf{X}' the Frobenius norm $\|\mathbf{X}\|_F$ and the Moore-Penrose inverse as \mathbf{X}^+ . When \mathbf{X} is symmetric, we use $\lambda_{\max}(\mathbf{X})$ and $\lambda_{\min}(\mathbf{X})$ to denote the largest and smallest eigenvalues, respectively. \mathbf{I}_p and $\mathbf{O}_{p \times 1}$ denote the $p \times p$ identity matrix and $p \times 1$ vector of zeros, and $\mathbf{1}\{\cdot\}$ is the indicator function. The operator \xrightarrow{P} denotes convergence in probability, \xrightarrow{D} convergence in distribution,

and plim probability limit. We use $(n, T) \rightarrow \infty$ to signify that n and T pass jointly to infinity.

fi

This section considers panel structure models without endogeneity. It is convenient to assume first that the number of groups is known and later consider the determination of the number of unknown groups.

Given a panel data set $\{(y_{it}, x_{it})\}$ for $i = 1, \dots, n$ and $t = 1, \dots, T$ it is proposed to use fixed effects quasi maximum likelihood to estimate the unknown parameters by solving the minimization problem

$$\min_{\{\beta_i, \mu_i\}} \frac{1}{nT} \sum_{i=1}^n \sum_{t=1}^T \ln f(y_{it}; x_{it}, \beta_i, \mu_i) \quad (2.1)$$

Here $f(y_{it}; x_{it}, \beta_i, \mu_i)$ denotes the logarithm of the pseudo-true conditional density function of y_{it} given x_{it} the history of (y_{it}, x_{it}) , and (β_i, μ_i) , where β_i are scalar individual effects and μ_i are 1×1 vectors of parameters of interest. Traditionally, econometric work has assumed that the β_i are common for all cross sectional units, leading to a homogeneous panel with individual heterogeneity modeled through β_i alone. At the other extreme, the β_i are assumed to differ across individuals and each is estimated at a slow rate without pooling across section. The present paper allows the true values of β_i denoted β_i^0 to follow a group pattern of the general form

$$\beta_i^0 = \sum_{k=1}^{K_0} \mathbf{1}_{k,i} \beta_k^0 \quad (2.2)$$

Here $\beta_j^0 \neq \beta_k^0$ for any $j \neq k$, $\cup_{k=1}^{K_0} \mathbf{1}_{k,i} = \{1, 2, \dots, K_0\}$ and $\mathbf{1}_{k,i} \cap \mathbf{1}_{j,i} = \emptyset$ for any $j \neq k$. Let $n_k = \#\{i : \mathbf{1}_{k,i} = 1\}$ denote the cardinality of the set $\mathbf{1}_{k,i}$. In the economic growth literature (e.g., Phillips and Sul, 2007a), n_k corresponds to the number of convergence clubs and countries (indexed by i) within the same k^{th} club share the same (slope) parameter vector β_k^0 . In the market entry-exit example (e.g., Hahn and Moon, 2010), n_k denotes the number of pure Nash equilibria and markets (indexed by i) selecting the same equilibrium over time exhibit the same parameter vector.

For now, we assume that the number of groups, n_k , is known and fixed but that each individual's group membership is unknown. In addition, following Sun (2005), Lin and Ng (2012), and BM, we implicitly assume that individual group membership does not vary over time. Let $\alpha \equiv (\alpha_1, \dots, \alpha_{K_0})$ and $\beta \equiv (\beta_1, \dots, \beta_N)$. We denote the true values of β_i, μ_i, α and β as $\beta_i^0, \mu_i^0, \alpha^0$ and β^0 respectively. The econometric task is to infer each individual's group identity and to estimate the group-specific parameters β_i^0 . Some examples of models that fall within this framework and the scope of our methodology are as follows.

(Linear panel) The linear panel structure model is generated according to

$$y_{it} = \beta_i^0 x_{it} + \mu_i^0 + \epsilon_{it} \quad (2.3)$$

where x_{it} is a 1×1 vector of exogenous or predetermined variables, μ_i^0 is an individual fixed effect, β_i^0 is a 1×1 vector of slope parameters, and ϵ_{it} is the idiosyncratic error term with mean zero. Gaussian quasi-maximum likelihood estimation (QMLE) of β_i^0 and μ_i^0 is achieved by minimizing (2.1) with $f(y_{it}; x_{it}, \beta_i^0, \mu_i^0) = \frac{1}{\sigma} \exp\left(-\frac{1}{\sigma} (y_{it} - \beta_i^0 x_{it} - \mu_i^0)\right)$ and $\epsilon_{it} = (y_{it} - \beta_i^0 x_{it} - \mu_i^0)'$

(Linear panel with quantile restrictions) Consider the model in (2.3) with the quantile restriction: $\mathbb{P}(it \leq 0 | it, i, i) = \tau$; see, e.g., Kato, Galvo, and Montes-Rojas (2012). We can estimate β_i and α_i by minimizing (2.1) with $\rho_\tau(it; i, i) = \tau(it - i) - \tau(i)it$ where $\tau(\cdot) = \{ - \Phi(-\cdot) \}$ is a smoothed version of the usual check function with Φ being a CDF-type kernel function and h a bandwidth parameter.

(Binary choice panel) The dynamic binary choice panel data model is characterized by $it = \mathbf{1}\{i'it + i - it \geq 0\}$ where it, i and it are as defined in Example 1. In this case, $-\rho_\tau(it; i, i) = it \ln \left(\frac{it - i}{it - i} \right) + (1 - it) \ln \left[1 - \left(\frac{it - i}{it - i} \right) \right]$ where $it = (it, i)'$ and $\Phi(\cdot)$ denotes the conditional CDF (standard logistic or normal) of it given i and the history of (it, it) .

(Tobit panel) The Tobit panel is characterized by $it = \max(0, i'it + i + it)$ where it, i and it are defined as in the above examples. For clarity, assume that it 's are independent and identically distributed (IID) $(0, \frac{2}{\varepsilon})$ given it and the history of (it, it) . In this case, $-\rho_\tau(it; i, i, \frac{2}{\varepsilon}) = \mathbf{1}\{it = 0\} \ln \left(\frac{it - i}{it - i} \right) + \mathbf{1}\{it > 0\} \ln \left[\frac{\phi \left(\frac{it - i}{\varepsilon} \right)}{\Phi \left(\frac{it - i}{\varepsilon} \right)} \right]$ where $it = (it, i)'$ and ϕ denotes the standard normal PDF and CDF, respectively. The presence of the common parameter $\frac{2}{\varepsilon}$ can be addressed by extending the asymptotic analysis below.

f **α** **β**

Following Hahn and Newey (2004) and Hahn and Kuersteiner (2011), the profile log-likelihood function is

$$l_{1,NT}(\beta) = \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \rho_\tau(it; i, i) \quad (2.4)$$

where $\hat{i}_i(i) = \arg \min_{\mu_i} \frac{1}{T} \sum_{t=1}^T \rho_\tau(it; i, i)$. Motivated by the literature on group Lasso (e.g., Yuan and Lin 2006), we propose to estimate β and α by minimizing the following PPL criterion function

$$l_{1NT, \lambda_1}^{(K_0)}(\beta, \alpha) = l_{1,NT}(\beta) + \frac{1}{N} \sum_{k=1}^{K_0} \|\hat{i}_k - \alpha_k\| \quad (2.5)$$

where $\lambda_1 = \lambda_{1NT}$ is a tuning parameter. Minimizing the above criterion function produces *classifier-Lasso* (C-Lasso) estimates $\hat{\beta}$ and $\hat{\alpha}$ of β and α respectively. Let \hat{i}_i and $\hat{\alpha}_k$ denote the i -th and k -th columns of $\hat{\beta}$ and $\hat{\alpha}$, respectively, i.e., $\hat{\alpha} \equiv (\hat{\alpha}_1, \dots, \hat{\alpha}_K)$ and $\hat{\beta} \equiv (\hat{\beta}_1, \dots, \hat{\beta}_N)$.

The penalty term in (2.5) takes a novel mixed *additive-multiplicative* form that does not appear in the literature. Traditional Lasso includes additive penalty terms to an objective function by differentiating zeros from non-zero-valued parameters to select relevant regressors. In contrast, the C-Lasso has K_0 additive terms, each of which takes a multiplicative form as the product of N_0 separate penalties. The multiplicative component is needed because for each i_0 can take any one of the N_0 *unknown* values, $\frac{0}{1}, \dots, \frac{0}{K_0}$. We do not know *a priori* to which point i_0 should shrink, and all N_0 possibilities must be allowed. Each of the N_0 penalty terms in the multiplicative expression permits i_0 to shrink to a particular *unknown* group-level parameter vector α_k . The summation component is needed because we need to pull information from all cross sectional units in order to identify $\{i_0\}$ and $\{\alpha_k\}$ jointly. Our approach differs from the prototypical Lasso method of Tibshirani (1996) that shrinks a parameter to zero as well as the group Lasso method of Yuan and Lin (2006) that shrinks a parameter vector to a zero vector. The main purpose in the latter papers

is to select relevant variables while C-Lasso is designed to determine group membership for each individual. As emphasized in the Introduction, both problems enjoy the same motivation of parameter sparsity despite their different nature. C-Lasso has the additional motivation of classification of unknown parameters into *a priori* unknown groups each with their own *unknown* parameters.

Note that the objective function in (2.5) is not convex in β even though it is (conditionally) convex in k when one fixes j for \neq . The supplement provides an iterative algorithm to obtain the estimates $\hat{\alpha}$ and $\hat{\beta}$.

Let $i(i) \equiv \arg \min_{\mu_i} \Psi_i(i, i)$ where $\Psi_i(i, i) \equiv \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\mu_i(i, i)]$. Note that $i^0 = i(i^0)$. Let $i(i, i) \equiv \mu_i(i, i)$ and $i(i, i) \equiv \beta_i(i, i)$. Let i^{μ_i} and $i^{\mu_i \mu_i}$ denote the first and second derivatives of i with respect to i . Define i^{μ_i} and $i^{\mu_i \mu_i}$ similarly. For notational simplicity, denote $i_t \equiv i(i_t, i_t^0)$ and similarly for $i_t^{\mu_i}$ and $i_t^{\mu_i \mu_i}$. Define

$$\begin{aligned} i_U &\equiv \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\mu_i(i_t)] & i_V &\equiv \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\beta_i(i_t)] & i_{U2} &\equiv \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\mu_i^{\mu_i}(i_t)] & i_{V2} &\equiv \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\beta_i^{\mu_i}(i_t)] \\ U_{it} &\equiv i_t - \frac{i_U}{i_V} i_t & U_{it}^{\beta_i} &\equiv \beta_i(i_t) - \frac{i_U}{i_V} \beta_i'(i_t) & \text{and } U_{it}^{\mu_i} &\equiv \mu_i(i_t) - \frac{i_U}{i_V} \mu_i'(i_t) \end{aligned}$$

Let $\Omega_{iT} \equiv \frac{1}{T} \sum_{t=1}^T \sum_{s=1}^T \mathbb{E}(U_{it} U_{is}')$, $\mathbb{H}_{iT} \equiv \frac{1}{T} \sum_{t=1}^T \mathbb{E}[U_{it}^{\beta_i}]$ and $\mathbb{H}_{kNT} \equiv \frac{1}{N_k} \sum_{i \in G_k^0} \mathbb{H}_{iT}$. Define the two expected Hessian matrices for cross sectional unit i :

$$i_{\mu\mu}(i) \equiv \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\mu_i^{\mu_i}(i_t, i_t)] \quad \text{and} \quad i_{\beta\beta}(i) \equiv \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\beta_i^{\mu_i}(i_t) + \mu_i(i_t) \frac{i(i_t)}{i} \right]$$

where $\beta_i^{\mu_i}(i) = \beta_i^{\mu_i}(i, i(i))$ and similarly for $\mu_i^{\mu_i}(i)$. Let \min_i denote $\min_{1 \leq i \leq N}$ and similarly for \max_i . We make the following assumptions

ASSUMPTION A1. (i) For each $\{i_t : t = 1, 2, \dots\}$ is stationary strong mixing with mixing coefficients $\rho_i(\cdot)$. $\rho_i(\cdot) \equiv \max_i \rho_i(\cdot)$ satisfies $\rho_i(\cdot) \leq \alpha^\tau$ for some $\alpha < 1$ and $\tau \in (0, 1)$. $\{i_t : t = 1, 2, \dots\}$ are independent across

(ii) For each $0 < \min_i [\inf_{(\beta_i, \mu_i): \|(\beta_i, \mu_i) - (\beta_i^0, \mu_i^0)\| > \eta} \Psi_i(i, i) - \Psi_i(i^0, i^0)] > 0$

(iii) Let Θ denote the parameter space for $i = (i, i)'$. Θ is a compact and convex subset of \mathbb{R}^{p+1} such that $i^0 = (i^0, i^0)'$ lies in the interior of Θ for each

(iv) Let $\|\cdot\| \equiv \sum_{j=1}^{p+1} \cdot_j$ and $v(i, i) \equiv |v| (i, i) = (v_1 \cdots v_{p+1})$ where $v = (v_1 \cdots v_{p+1})$ is a vector of nonnegative integers and $v_{(j)}$ denotes the j th element of v . There exists a function $\phi(\cdot)$ such that $\sup_{\theta \in \Theta} \|\phi(v(i, i))\| \leq \phi(i_t) \|v(i, i) - v(i, i)^-\| \leq \phi(i_t) \|v(i, i) - v(i, i)^-\|$ for any $v(i, i)^- \in \Theta$ and $\|\cdot\| \leq 3$ and $\max_i \mathbb{E} \|\phi(i_t)\|^q \leq M$ for some $M < \infty$ and $q \geq 6$

(v) There exists a constant $H > 0$ such that $\min_i \inf_{\beta \in \mathcal{B}} i_{\mu\mu}(i) \geq H$ and $\min_i \min(\beta_i^{\mu_i}(i^0)) \geq H$

(vi) There exists a constant $\alpha < 1$ such that $\min_{1 \leq k < l \leq K_0} \|\frac{0}{k} - \frac{0}{l}\| \geq \alpha$

(vii) i_0 is fixed and $k \rightarrow k \in (0, 1)$ for each $i = 1, \dots, 0$ as $\rightarrow \infty$

ASSUMPTION A2. (i) $\frac{2}{1} (\ln \cdot)^{6+2\nu} \rightarrow \infty$ and $\frac{1}{1} (\ln \cdot)^\nu \rightarrow 0$ for some $\nu > 0$ as $(\cdot) \rightarrow \infty$

(ii) $\frac{1}{2} \frac{1}{1} (\ln \cdot)^9 \rightarrow 0$ and $\frac{2}{2} \frac{1}{1} \rightarrow \infty$ as $(\cdot) \rightarrow \infty$.

ASSUMPTION A3. (i) For each $k = 1, \dots, K_0$, $\Omega_k \equiv \lim_{(N_k, T) \rightarrow \infty} \frac{1}{N_k} \sum_{i \in G_k^0} \Omega_{iT}$ exists and $\Omega_k \neq 0$
(ii) For each $k = 1, \dots, K_0$, $\mathbb{H}_k \equiv \lim_{(N_k, T) \rightarrow \infty} \mathbb{H}_{kNT}$ exists and $\mathbb{H}_k \neq 0$

Assumption A1(i) imposes conditions on $\{\beta_{it}\}$ which are commonly assumed for dynamic nonlinear panel data model; see, e.g., Hahn and Kuersteiner (2011) and Lee and Phillips (2015). With more complicated notation, we can relax the stationarity assumption along the time dimension. A1(ii) imposes an identification condition for the joint identification of (β_i, γ_i) for each i . A1(iii) restricts the parameter space and it is possible to allow Θ to be β_i -dependent. A1(iv) specifies the smoothness and moment conditions on β_i or objects associated with it. A1(v), in conjunction with A1(ii) and (iv), implies that $\min_i [\inf_{\mu_i: |\mu_i - \mu_i(\beta_i)| > \eta} \Psi_i(\beta_i, \gamma_i) - \Psi_i(\beta_i, \gamma_i(\beta_i))] > 0$ and $\min_i [\inf_{\beta_i: \|\beta_i - \beta_i^0\| > \eta} \Psi_i(\beta_i, \gamma_i(\beta_i)) - \Psi_i(\beta_i^0, \gamma_i(\beta_i^0))] > 0$. A1(vi) specifies that the group-specific parameters are separated from each other, similar to the separation requirement in Hahn and Moon (2010). A1(vii) implies that each group has an asymptotically non-negligible membership number of individuals as $N_k \rightarrow \infty$. This assumption can also be relaxed at the cost of more lengthy arguments. Assumption A2(i) imposes conditions on β_i all of which hold if

$$\beta_i \propto \beta_i^{-a} \text{ for any } \beta_i \in (0, 1, 2) \quad (2.6)$$

A2(ii) is needed to ensure some higher order terms are asymptotically negligible. A3 is used to derive the asymptotic bias and variance of the C-Lasso estimator. The theory developed below under these conditions does not require correct specification of the likelihood function and the C-Lasso asymptotics apply under the general QMLE setup.

ff

The following theorem establishes the consistency of the PPL estimates $\{\hat{\beta}_i\}$ and $\{\hat{\gamma}_k\}$

Suppose that Assumption A1 holds and $\beta_i = (1)$. Then (i) $\hat{\beta}_i - \beta_i^0 = O_P(N^{-1/2 + \epsilon})$ for $\epsilon > 0$
(ii) $\frac{1}{N} \sum_{i=1}^N \|\hat{\beta}_i - \beta_i^0\|^2 = O_P(N^{-1})$ and (iii) $(\hat{\beta}_{(1)}, \hat{\gamma}_{(K_0)}) - (\beta_{(1)}^0, \gamma_{(K_0)}^0) = O_P(N^{-1/2})$
where $(\hat{\beta}_{(1)}, \hat{\gamma}_{(K_0)})$ is a suitable permutation of $(\hat{\beta}_1, \hat{\gamma}_{K_0})$

Theorem 2.1(i)-(ii) establish the pointwise and mean-square convergence of $\hat{\beta}_i$. Theorem 2.1(iii) indicates that the group-specific parameters $(\beta_1^0, \gamma_{K_0}^0)$ can be estimated consistently by $(\hat{\beta}_1, \hat{\gamma}_{K_0})$ subject to permutation. As expected and consonant with other Lasso limit theory, the pointwise convergence rate of $\hat{\beta}_i$ depends on the rate at which the tuning parameter λ_1 converges to zero. Somewhat unexpectedly, this requirement is not the case either for mean-square convergence of $\hat{\beta}_i$ or convergence of $\hat{\gamma}_k$. For notational simplicity, hereafter we simply write $\hat{\gamma}_k$ for $\hat{\gamma}_{(k)}$ as the consistent estimator of γ_k^0 , and define

$$\hat{\gamma}_k = \sum_{i \in G_k^0} \hat{\beta}_i : \hat{\beta}_i = \hat{\beta}_k^0 \text{ for } i \in G_k^0 \quad (2.7)$$

fi

Roughly speaking, a classification method is consistent if it classifies each individual to the correct group w.p.a.1. For a rigorous statement of this property we define

$$\hat{\gamma}_{kNT, i} \equiv \sum_{i \in G_k^0} \hat{\beta}_i \mid \beta_i \in G_k^0 \text{ and } \hat{\gamma}_{kNT, i} \equiv \sum_{i \in G_k^0} \hat{\beta}_i \mid \beta_i \in G_k^0 \quad (2.8)$$

where $\hat{0} = 1$ and $\hat{0} = 1$. Let $\hat{k}_{kNT} = \cup_{i \in G_k^0} \hat{k}_{kNT,i}$ and $\hat{k}_{kNT} = \cup_{i \in \hat{G}_k} \hat{k}_{kNT,i}$. \hat{k}_{kNT} and \hat{k}_{kNT} mimic Type I and II errors in statistical tests: \hat{k}_{kNT} denotes the error event of not classifying an element of G_k^0 into the estimated group \hat{k} ; and \hat{k}_{kNT} denotes the error event of classifying an element that does not belong to G_k^0 into the estimated group \hat{k} . Both types of errors must be controlled. We use the following definition.

fi The classification is *individually consistent* if $(\hat{k}_{kNT,i}) \rightarrow 0$ as $(n) \rightarrow \infty \forall i \in G_k^0$ and $i \in \{1, \dots, 0\}$ and $(\hat{k}_{kNT,i}) \rightarrow 0$ as $(n) \rightarrow \infty \forall i \in \hat{k}$ and $i \in \{1, \dots, 0\}$. It is *uniformly consistent* if $(\cup_{k=1}^{K_0} \hat{k}_{kNT}) \rightarrow 0$ and $(\cup_{k=1}^{K_0} \hat{k}_{kNT}) \rightarrow 0$ as $(n) \rightarrow \infty$.

The following theorem establishes uniform consistency for the PPL classifier.

Suppose that Assumptions A1-A2 hold. Then (i) $(\cup_{k=1}^{K_0} \hat{k}_{kNT}) \leq \frac{K_0}{k=1} (\hat{k}_{kNT}) \rightarrow 0$ as $(n) \rightarrow \infty$ and (ii) $(\cup_{k=1}^{K_0} \hat{k}_{kNT}) \leq \frac{K_0}{k=1} (\hat{k}_{kNT}) \rightarrow 0$ as $(n) \rightarrow \infty$.

Theorem 2.2 implies that all individuals within a group, say G_k^0 can be simultaneously correctly classified into the same group (denoted \hat{k}) w.p.a.1. Conversely, all individuals that are classified into the same group, say \hat{k} simultaneously correctly belong to the same group (G_k^0) w.p.a.1. Let $\hat{0} \equiv \{1, 2, \dots\} \setminus (\cup_{k=1}^{K_0} \hat{k})$ and $\hat{i}_{iNT} \equiv \{i \in \hat{0}\}$. Theorem 2.2(i) implies that $(\cup_{1 \leq i \leq N} \hat{i}_{iNT}) \leq \frac{K_0}{k=1} (\hat{k}_{kNT}) \rightarrow 0$. That is, all individuals can be classified into one of the G_0 groups w.p.a.1. Nevertheless, when n is not large, a small percentage of individuals could be left unclassified if we stick with the classification rule in (2.7). To ensure that all individuals are classified into one of the G_0 groups in finite samples, we can modify the classifier. In particular, we classify $i \in \hat{k}$ if $\hat{i} = \hat{k}$ for some $i = 1, \dots, 0$ and $i \in \hat{l}$ for some $i = 1, \dots, 0$ if $\|\hat{i} - \hat{l}\| = \min\{\|\hat{i} - \hat{1}\|, \dots, \|\hat{i} - \hat{K}_0\|\}$ and $\frac{K_0}{k=1} \mathbf{1}\{\hat{i} = \hat{k}\} = 0$. Since the event $\frac{K_0}{k=1} \mathbf{1}\{\hat{i} = \hat{k}\} = 0$ occurs w.p.a.1 uniformly in n we can ignore it in large samples in subsequent theoretical analysis and restrict our attention to the classification rule in (2.7) to avoid confusion.

Let $\hat{k} \equiv \frac{N}{i=1} \mathbf{1}\{i \in \hat{k}\}$. The following corollary studies the consistency of \hat{k} .

Suppose that Assumptions A1-A2 hold. Then $\hat{k} - k = o_P(1)$ for $n \rightarrow \infty$.

The following theorem reports the oracle property of the Lasso estimator $\{\hat{k}\}$.

Suppose Assumptions A1-A3 hold. Then $\sqrt{k}(\hat{k} - \frac{0}{k}) - \mathbb{H}_{kNT}^{-1} \mathbb{B}_{kNT} \xrightarrow{D} (0, \mathbb{H}_k^{-1} \Omega_k (\mathbb{H}_k^{-1})')$ where $\mathbb{B}_{kNT} = \mathbb{B}_{1kNT} - \mathbb{B}_{2kNT}$, $\mathbb{B}_{1kNT} = \frac{1}{\sqrt{N_k T^3}} \sum_{i \in G_k^0} \frac{-1}{iV} \sum_{s=1}^T \sum_{t=1}^T i_s \mathbb{U}_{it}^{\mu_i}$ and $\mathbb{B}_{2kNT} = \frac{1}{2\sqrt{N_k T}} \sum_{i \in G_k^0} \frac{-2}{iV} (\sum_{i=2}^m \frac{m_{iV2}}{m_{iV}} iU) (\frac{1}{\sqrt{T}} \sum_{t=1}^T it)^2$ for $n \rightarrow \infty$.

\mathbb{B}_{kNT} is written as the difference between two terms that are derived from the first and second order Taylor expansions of the PPL estimating equation, respectively. Comparing the above result with HK, we find that the quantities Ω_k , \mathbb{H}_k and \mathbb{B}_k coincide with the corresponding terms in HK; see the remark after Lemma S1.12 for details. Then we can use the formula in HK to estimate the asymptotic bias and variance with obvious modifications. Alternatively, we can use the jackknife to correct bias; see Hahn and Newey (2004) and Dhaene and Jachmans (2015) for static and dynamic models, respectively.

If group membership is known, the *oracle* estimator of β_k is given by $\hat{G}_k^0 \equiv \arg \min_{\alpha_k} \frac{1}{N_k T} \sum_{t=1}^T \left(\beta_{it}; \beta_k \hat{\beta}_i(\beta_k) \right)$. Then following our asymptotic analysis or that of HK, we can readily show that $\sqrt{\frac{1}{k}} \left(\hat{G}_k^0 - \beta_k^0 \right) - \mathbb{H}_{kNT}^{-1} \mathbb{B}_{kNT} \xrightarrow{D} \left(0 \ \mathbb{H}_k^{-1} \Omega_k (\mathbb{H}_k^{-1})' \right)$ under Assumptions A1 and A3. Theorem 2.4 indicates that the PPL estimator $\hat{\beta}_k$ achieves the same limit distribution as this oracle estimator. In this sense, we say that the PPL estimators $\{\hat{\beta}_k\}$ enjoy the asymptotic oracle property. In addition, given the estimated groups $\hat{\beta}_k$ we can obtain the post-Lasso estimator of β_k by $\hat{G}_k \equiv \arg \min_{\alpha_k} \frac{1}{N_k T} \sum_{t=1}^T \left(\beta_{it}; \beta_k \hat{\beta}_i(\beta_k) \right)$. The following theorem reports the asymptotic distribution of \hat{G}_k .

Suppose Assumptions A1-A3 hold. Then $\sqrt{\frac{1}{k}} \left(\hat{G}_k - \beta_k^0 \right) - \mathbb{H}_{kNT}^{-1} \mathbb{B}_{kNT} \xrightarrow{D} \left(0 \ \mathbb{H}_k^{-1} \Omega_k (\mathbb{H}_k^{-1})' \right)$ for $k = 1, \dots, K$ where \mathbb{B}_{kNT} is as defined in Theorem 2.4.

Theorems 2.4 and 2.5 indicate that $\hat{\beta}_k$ and \hat{G}_k are asymptotically equivalent. In a totally different framework, Belloni and Chernozhukov (2013) study post-Lasso estimators which apply OLS to the model selected by first-step penalized estimators and show that the post-Lasso estimators perform at least as well as Lasso in terms of rate of convergence and have the advantage of smaller bias. Correspondingly, it would be interesting to compare the higher-order asymptotic properties of $\hat{\beta}_k$ and \hat{G}_k in future work.

Note that our asymptotic results are ‘‘pointwise’’ in the sense that the unknown parameters are treated as fixed. The implication is that in finite samples, the distributions of our estimators can be quite different from normal, as discussed in Leeb and Pötscher (2008, 2009). This is a well-known challenge for shrinkage estimators. Despite its importance, developing a thorough theory on uniform inference in this context is beyond the scope of the present work.

In practice, the exact number of groups is typically unknown. We assume that K_0 is bounded from above by a finite integer K_{\max} and study the determination of the number of groups via some information criterion (IC). By minimizing (2.5) with K_0 replaced by K we obtain the C-Lasso estimates $\{\hat{\beta}_i(\beta_k), \hat{\beta}_k(\beta_k)\}$ of $\{\beta_i, \beta_k\}$ where we make the dependence of $\hat{\beta}_i$ and $\hat{\beta}_k$ on (β_k) explicit. As above, we classify individual i into group $\hat{\beta}_k(\beta_k)$ if and only if $\hat{\beta}_i(\beta_k) = \hat{\beta}_k(\beta_k)$, i.e., $\hat{\beta}_k(\beta_k) \equiv \{i \in \{1, 2, \dots, N\} : \hat{\beta}_i(\beta_k) = \hat{\beta}_k(\beta_k)\}$ for $k = 1, \dots, K$. Let $\hat{G}(K) \equiv \{\hat{G}_1(\beta_k), \dots, \hat{G}_K(\beta_k)\}$. The post-Lasso estimator of β_k^0 is denoted as $\hat{G}_{k(K, \lambda_1)}$. We propose to select K to minimize

$$Q_1(K) \equiv \frac{2}{N T} \sum_{k=1}^K \sum_{i \in \hat{G}_k(K, \lambda_1)} \sum_{t=1}^T \left(\beta_{it}; \hat{G}_{k(K, \lambda_1)} \hat{\beta}_i(\hat{G}_{k(K, \lambda_1)}) \right) + \lambda_1 N T \quad (2.9)$$

where $\lambda_1 N T$ is a tuning parameter. Let $\hat{K} \equiv \arg \min_{1 \leq K \leq K_{\max}} Q_1(K)$. See Wang, Li, and Tsai (2007), Liao (2013), and Lu and Su (2016) for the use of a similar IC in various contexts.

Let $\mathcal{G}^{(K)} \equiv \{G_{K,1}, \dots, G_{K,K}\}$ be any K -partition of $\{1, 2, \dots, N\}$ and \mathcal{G}_K a collection of all such partitions. Let $\hat{G}^{(K)} \equiv \frac{2}{N T} \sum_{k=1}^K \sum_{i \in G_{K,k}} \sum_{t=1}^T \left(\beta_{it}; \hat{G}_{K,k} \hat{\beta}_i(\hat{G}_{K,k}) \right)$ where $\hat{G}_{K,k} \equiv \arg \min_{\alpha_k} \frac{1}{N_k T} \sum_{t=1}^T \left(\beta_{it}; \beta_k \hat{\beta}_i(\beta_k) \right)$. We add the following two assumptions.

ASSUMPTION A4. As $(N, T) \rightarrow \infty$, $\min_{1 \leq K < K_0} \inf_{G^{(K)} \in \mathcal{G}_K} \hat{G}^{(K)} \xrightarrow{P} \beta_k^0$ where $\beta_k^0 \equiv \lim_{(N, T) \rightarrow \infty} \frac{2}{N T} \sum_{k=1}^{K_0} \sum_{i \in G_k^0} \sum_{t=1}^T \mathbb{E} \left[\left(\beta_{it}; \beta_k^0 \beta_k^0 \right) \right]$

ASSUMPTION A5. As $(k) \rightarrow \infty$, $\lambda_{1NT} \rightarrow 0$ and $\lambda_{2NT} \rightarrow \infty$.

Assumption A4 is intuitively clear and applies under primitive conditions in a variety of models, such as panel autoregressions. It requires that all under-fitted models yield asymptotic mean square errors that are larger than $\frac{2}{3}$, which is delivered by the true model. A5 reflects the usual conditions for the consistency of model selection: λ_{1NT} cannot shrink to zero either too fast or too slowly.

The following theorem justifies the use of (2.9) as a selector criterion for

Suppose Assumptions A1-A5 hold. Then $\mathbb{P}(\hat{g}_1 = g_0) \rightarrow 1$ as $(k) \rightarrow \infty$

As Theorem 2.6 indicates, as long as λ_1 satisfies Assumption A2(i), we can ensure that the correct number of groups is chosen w.p.a.1. In practice, we can fine-tune this parameter over a finite set, e.g., $\Lambda_1 \equiv \{ \lambda_1 = j^{-1/3} \mid j = 0, \dots, J \}$ for some $J > 0$ and $\lambda_1 \in \Lambda_1$. That is, we pick up $\lambda_1 \in \Lambda_1$ such that $\mathbb{P}(\hat{g}_1 = g_0)$ is minimized. We can show that with such a choice of λ_1 Theorem 2.6 continues to hold. Alternatively, we can consider a data-driven cross-validation procedure.

For the linear model in (2.3) with $\mathbb{E}(it | it, i) = 0$, we have $(it; i, i) = \frac{1}{2} (it - i, it - i)' \hat{\beta}_i(i) = \bar{y}_i - i' \bar{y}_i$ and $\lambda_{1NT}(\beta) = \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T (\tilde{y}_{it} - i' \tilde{y}_{it})^2$ where $\bar{y}_i = \frac{1}{T} \sum_{t=1}^T y_{it}$, $\tilde{y}_{it} = y_{it} - \bar{y}_i$ and \tilde{y}_i and \tilde{y}_{it} are analogously defined. So the PPL problem becomes the penalized least squares (PLS) problem considered in Su, Shi, and Phillips (2014, SSP hereafter). In addition, we can verify that $(i, i) = \mathbb{E}(\bar{y}_i) - i' \mathbb{E}(\bar{y}_i)$, $(it; i, i) = - (it - i, it - i)' i = - (it - i, it - i)' i' (it; i, i) = - (it - i, it - i)' i' \mathbb{U}_{it} = - it [it - \mathbb{E}(\bar{y}_i)]' \mathbb{U}_{it}^{\beta_i} = [it - \mathbb{E}(\bar{y}_i)]' i' \mathbb{U}_{it}^{\mu_i} = it - \mathbb{E}(\bar{y}_i)$, $i' \mu_i(i) = 1$, $i' \beta_i(i) = \frac{1}{T} \sum_{t=1}^T \mathbb{E}\{[it - \mathbb{E}(\bar{y}_i)] [it - \mathbb{E}(\bar{y}_i)]'\}$

$$\Omega_{iT} = \frac{1}{T} \sum_{t=1}^T \sum_{s=1}^T \mathbb{E}\{it is [it - \mathbb{E}(\bar{y}_i)] [is - \mathbb{E}(\bar{y}_i)]'\} \text{ and}$$

$$\mathbb{H}_{iT} = \frac{1}{T} \sum_{t=1}^T \mathbb{E}\{[it - \mathbb{E}(\bar{y}_i)] [it - \mathbb{E}(\bar{y}_i)]'\}$$

With the above calculations, we can readily verify that Assumptions A1(ii), (iv)-(v) and A3 hold under weak conditions. In addition, we can show that

$$\mathbb{B}_{1kNT} = \frac{-1}{\sqrt{k}^3} \sum_{i \in G_k^0} \sum_{t=1}^T \sum_{s=1}^T it [is - \mathbb{E}(\bar{y}_i)] = \mathbb{B}_{1k} + o_P(1) \text{ and } \mathbb{B}_{2kNT}$$

analogously defined as \hat{G}_k . In practice, $\hat{G}_{(K,\lambda_1)}^2$ is frequently replaced by its natural logarithm as in standard BIC to obtain

$$\ln \hat{G}_{(K,\lambda_1)}^2 = \ln \hat{G}_{(K,\lambda_1)}^2 + \ln \hat{G}_{(K,\lambda_1)}^2 + \ln \hat{G}_{(K,\lambda_1)}^2 \quad (2.10)$$

which will be used in our simulations and applications. But because the fixed effects are eliminated in the within-group transformed model, the \sqrt{N} -convergence rates of their estimates won't play a role to ensure the selection consistency of $\hat{\beta}_1$. SSP show that the requirement on $\ln \hat{G}_{(K,\lambda_1)}^2$ can be relaxed with Assumption A5 replaced by:

ASSUMPTION A5*. As $(N, T) \rightarrow \infty$, $\ln \hat{G}_{(K,\lambda_1)}^2 \rightarrow 0$ and $\ln \hat{G}_{(K,\lambda_1)}^2 / \ln \hat{G}_{(K,\lambda_1)}^2 \rightarrow \infty$ where $\ln \hat{G}_{(K,\lambda_1)}^2 = \ln \hat{G}_{(K,\lambda_1)}^2$ if $\ln \hat{G}_{(K,\lambda_1)}^2$ is strictly exogenous and $\ln \hat{G}_{(K,\lambda_1)}^2 = \ln \hat{G}_{(K,\lambda_1)}^2$ otherwise.

In some applications, certain parameters of interest may be common across all individuals whereas others are group-specific. For instance, Pesaran, Shin, and Smith (1999) constrain the long-run coefficients to be identical across individuals while assuming the short-run coefficients to be heterogenous, or in our case, group-specific. Example 4 above is another instance. To keep up with the early notation, we write the negative log-likelihood function as $l(\beta; \alpha, \gamma_i)$ where β is the common parameter and the γ_i have a group structure as before. The negative profile log-likelihood function now becomes $l_{1,NT}(\beta) = \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T l(\beta; \gamma_i, \alpha_i)$ where $\hat{\gamma}_i(\alpha_i) = \arg \min_{\gamma_i} \frac{1}{T} \sum_{t=1}^T l(\beta; \gamma_i, \alpha_i)$. Then we can estimate β and α by minimizing the following PPL criterion function

$$l_{1,NT,\lambda_1}^{(K_0)}(\beta, \alpha) = l_{1,NT}(\beta) + \frac{1}{N} \sum_{i=1}^{K_0} \|\gamma_i - \alpha_i\| \quad (2.11)$$

Our previous analysis can be followed to establish uniform consistency for the classifier and the oracle property for the resulting estimators of the group-specific parameters γ_k and the common parameter β .

When we have time effects $\{\alpha_t\}$ we generally cannot eliminate them through transformation even in a linear panel structure model because of the slope heterogeneity. In this case, we need to estimate $(\alpha_1, \dots, \alpha_T)'$ jointly with β and α in (2.11). A formal asymptotic analysis of this case is left for future work.

This section considers penalized GMM estimation of linear panel structure models when some regressors are lagged dependent variables or endogenous.

$$\alpha = \beta$$

To stay focused, we restrict attention to the linear panel structure model in (2.3).¹ We consider the first differenced system

$$\Delta y_{it} = \alpha_i \Delta y_{it} + \Delta y_{it} \quad (3.1)$$

¹Extension to general nonlinear panel data models with endogeneity and nonadditive fixed effects (e.g., Fernández-Val and Lee 2013) is possible but rigorous analysis raises additional statistical challenges and is left for future research.

where, e.g., $\Delta_{it} = y_{it} - y_{i,t-1}$ for $i = 1, \dots, N$ and $t = 1, \dots, T$ and we assume that y_{i0} and y_{i0} are observed. Let Δ_{it} be a 1×1 vector of instruments for Δ_{it} with $\Delta_{it} \geq 0$. Define $\Delta_i = (\Delta_{i1} \dots \Delta_{iT})'$ with similar definitions for Δ_i and Δ_i . We propose to estimate β and α by minimizing the following penalized GMM (PGMM) criterion function²

$$Q_{2NT, \lambda_2}^{(K_0)}(\beta, \alpha) = Q_{2NT}(\beta) + \frac{2}{\lambda_2} \sum_{i=1}^N \Pi_{k=1}^{K_0} \|\Delta_{it} - \Delta_{it}\| \quad (3.2)$$

where $Q_{2NT}(\beta) = \frac{1}{N} \sum_{i=1}^N \frac{1}{T} \sum_{t=1}^T \Delta_{it} (\Delta_{it} - \Delta_{it})'$ and $\Delta_{iNT} = \frac{1}{T} \sum_{t=1}^T \Delta_{it} (\Delta_{it} - \Delta_{it})'$ is a 1×1 symmetric matrix that is asymptotically nonsingular and $2 = 2NT$ is a tuning parameter. Minimizing (3.2) produces the PGMM estimates $\tilde{\alpha}$ and $\tilde{\beta}$ where $\tilde{\alpha} \equiv (\tilde{\alpha}_1 \dots \tilde{\alpha}_{K_0})$ and $\tilde{\beta} \equiv (\tilde{\beta}_1 \dots \tilde{\beta}_N)$

Let $\tilde{y}_{i,z\Delta x} \equiv \frac{1}{T} \sum_{t=1}^T \Delta_{it} (\Delta_{it})'$ and $\tilde{y}_{i,z\Delta x} \equiv \mathbb{E}[\tilde{y}_{i,z\Delta x}]$. Let $\Delta_{it} \equiv (\Delta_{it} \ (\Delta_{it})' \ \Delta_{it})'$ and $\Delta_{it} \equiv \Delta_{it} (\Delta_{it} - \Delta_{it})'$ and $\Delta_{i,T}(\cdot) \equiv \frac{1}{\sqrt{T}} \sum_{t=1}^T \{\Delta_{it}(\cdot) - \mathbb{E}[\Delta_{it}(\cdot)]\}$. Let \mathcal{B}_i denote the parameter space for Δ_{it} . We make the following assumptions.

- ASSUMPTION B1. (i) $\mathbb{E}[\Delta_{it}^0] = 0$ for each $i = 1, \dots, N$ and $\Delta_{it} = 1$
(ii) $\sup_{\beta \in \mathcal{B}_i} \|\Delta_{i,T}(\cdot)\| = P(1) \frac{1}{N} \sum_{i=1}^N \|\Delta_{i,T}(\cdot)\|^2 = P(1)$ where $\Delta_{it} \in \mathcal{B}_i$ and $(\max_i \|\Delta_{i,T}(\cdot)\| \geq (\ln \cdot)^{3+\nu}) = O(1)$ for any $\nu > 0$ and $\Delta_{it} = 0$
(iii) $(\max_i \|\tilde{y}_{i,z\Delta x} - \tilde{y}_{i,z\Delta x}\| \geq \epsilon) = O(1)$ for any $\epsilon > 0$ and $\liminf_{(N,T) \rightarrow \infty} \min_i \min(\|\tilde{y}_{i,z\Delta x} - \tilde{y}_{i,z\Delta x}\|) = \frac{2}{Q} > 0$
(iv) There exist nonrandom matrices Δ_{it} such that $(\max_i \|\Delta_{iNT} - \Delta_{it}\| \geq \epsilon) = O(1)$ for any $\epsilon > 0$ and $\liminf_{N \rightarrow \infty} \min_i \min(\Delta_{it}) = \underline{W} > 0$
(v) There exists a constant $\alpha > 0$ such that $\min_{1 \leq k < l \leq K_0} \|\Delta_{it}^k - \Delta_{it}^l\| \geq \alpha$
(vi) Δ_{it} is fixed and $\Delta_{it} \rightarrow \Delta_{it} \in (0, 1)$ for each $i = 1, \dots, N$ as $\Delta_{it} \rightarrow \infty$

- ASSUMPTION B2. (i) $\frac{2}{\lambda_2} (\ln \cdot)^{6+2\nu} \rightarrow \infty$ and $\frac{2}{\lambda_2} (\ln \cdot)^\nu \rightarrow 0$ for some $\nu > 0$ as $\Delta_{it} \rightarrow \infty$
(ii) For any given $\epsilon > 0$, $\max_i (\|\Delta_{i,T}(\cdot) - \tilde{y}_{i,z\Delta x}\| \geq \epsilon) \rightarrow 0$ as $\Delta_{it} \rightarrow \infty$

- ASSUMPTION B3. (i) For each $i = 1, \dots, N$, $\Delta_{it}^k \equiv \frac{1}{N_k} \sum_{i \in G_k^0} \Delta_{it}^k \rightarrow \Delta_{it}^k = 0$ as $\Delta_{it} \rightarrow \infty$
(ii) For each $i = 1, \dots, N$, $\frac{1}{\sqrt{N_k T}} \sum_{i \in G_k^0} \Delta_{it}^k \xrightarrow{D} (0, \Delta_{it}^k)$ as $\Delta_{it} \rightarrow \infty$

Assumption B1(i) specifies moment conditions to identify Δ_{it} . B1(ii) is a high level condition. Its first part can be verified by applying Donsker's theorem. For example, if there exists \mathcal{F}_{it} a σ -field, such that $\{\Delta_{it} \mid \mathcal{F}_{it}\}$ is a stationary ergodic adapted mixingale with size -1 (e.g., White 2001, pp. 124-125), and $\text{Var}(\Delta_{i,T}(\cdot)) \rightarrow \Sigma_i \in (0, \infty)$ as $\Delta_{it} \rightarrow \infty$ for some $\Sigma_i > 0$ and any nonrandom $\Delta_{it} \in \mathbb{R}^d$ with $\|\Delta_{it}\| = 1$ then $\Delta_{i,T}(\cdot) \xrightarrow{D} (0, \Sigma_i)$ and the first part of B1(ii) follows. The second and third parts of B1(ii) can be verified by the Markov inequality and the application of Lemma S1.2(iii) in the supplement under strong

²We were unable to establish asymptotic theory for the case where the criterion $Q_{2,NT}$ is replaced by the fully pooled criterion $Q_{2,NT} = \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T z_{it} \Delta y_{it} - \beta_i^0 \Delta x_{it} = W_{NT} \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T z_{it} \Delta y_{it} - \beta_i^0 \Delta x_{it}$, where W_{NT} is asymptotically nonsingular. We also found that Arellano and Bond (1991) GMM estimation is not applicable to handle unobserved slope heterogeneity. Noticing this, Fernández-Val and Lee (2013) used a criterion similar to $Q_{2,NT}$ in the nonlinear panel setup. As we shall see, the use of $Q_{2,NT}$ means that the PGMM estimator generally does not have the oracle property.

mixing conditions. B1(iii) provides a rank condition to identify β_0 . B1(iv) is automatically satisfied for $\Sigma_{NT} = \sigma^2 I$ the \times identity matrix. B1(v)-(vi) and B2(i) parallel A1(vi)-(vii) and A2(i). B2(ii) holds true by Lemma S1.2 in the supplement if $\{(\Delta_{it} - \beta_0) \geq 1\}$ is strong mixing with geometric decay rate and Δ_{it} has six plus moments.

B3(i)-(ii) can be verified under various primitive conditions. For example, if (a) $\mathbb{E} \|\Delta_{it} - \beta_0\|^{2+\sigma} < \infty$ for some $\sigma > 0$ (b) $\{(\Delta_{it} - \beta_0) \geq 1\}$ is strong mixing for each i with mixing coefficients $\alpha_i(\cdot)$ that satisfy $\frac{1}{N_k} \sum_{i \in G_k^0} \sum_{\tau=1}^{\infty} \alpha_i(\tau)^{(2+\sigma)/\sigma} < \infty$ (c) $\{(\Delta_{it} - \beta_0)\}$ is stationary along the time dimension and IID along the individual dimension for all $i \in G_k^0$, and (d) $\beta_0 = \beta_0 \forall i \in G_k^0$ then B3(i) is satisfied with $\beta_k = \{\mathbb{E}[\Delta_{it} - \beta_0]\}' \mathbb{E}[\Delta_{it} - \beta_0] \forall i \in G_k^0$. To verify B3(ii), for simplicity we assume that $\Sigma_{NT} = \sigma^2 I$ and make the following decomposition

$$\begin{aligned}
& \frac{1}{\sqrt{N_k}} \sum_{i \in G_k^0} \sum_{t=1}^T \Delta_{it} \\
&= \frac{1}{N_k^{1/2} T^{3/2}} \sum_{i \in G_k^0} \sum_{s=1}^T \sum_{t=1}^T \mathbb{E}(\Delta_{is} - \beta_0) (\Delta_{it} - \beta_0) + \frac{1}{N_k^{1/2} T^{3/2}} \sum_{i \in G_k^0} \sum_{s=1}^T \sum_{t=1}^T \mathbb{E}(\Delta_{is} - \beta_0) (\Delta_{it} - \beta_0) \\
& \quad + \frac{1}{N_k^{1/2} T^{3/2}} \sum_{i \in G_k^0} \sum_{s=1}^T \sum_{t=1}^T \{[\Delta_{is} - \beta_0 - \mathbb{E}(\Delta_{is} - \beta_0)] (\Delta_{it} - \beta_0) - \mathbb{E}(\Delta_{is} - \beta_0) (\Delta_{it} - \beta_0)\} \\
& \equiv \beta_{kNT} + \gamma_{kNT} + \eta_{kNT} \text{ say,} \tag{3.3}
\end{aligned}$$

where β_{kNT} and γ_{kNT} contribute to the asymptotic bias and variance, respectively, and η_{kNT} is a term that is asymptotically negligible under suitable conditions. Then B3(ii) is satisfied with $\Sigma_{NT} = \sigma^2 I$ if $\beta_{kNT} = \frac{1}{N_k^{1/2} T^{1/2}} \sum_{i \in G_k^0} \sum_{t=1}^T \mathbb{E}(\Delta_{it} - \beta_0) \xrightarrow{D} (0, \sigma^2)$ and $\gamma_{kNT} = o_p(1)$ both of which can be verified by strengthening the conditions given in (a)-(c) above. Note that β_{kNT} signifies the asymptotic bias of $\tilde{\beta}_k$ which may not vanish asymptotically but can be corrected; see Section S2.2 in the supplement.³

We first establish the preliminary consistency rate of $(\tilde{\beta}, \tilde{\alpha})$.

Suppose Assumption B1 holds and $\beta_0 = \beta_0$. Then (i) $\tilde{\beta}_i - \beta_0 = o_p(N^{-1/2} + T^{-1/2})$ for $i = 1, \dots, N$ (ii) $\frac{1}{N} \sum_{i=1}^N \|\tilde{\beta}_i - \beta_0\|^2 = o_p(N^{-1})$ and (iii) $(\tilde{\beta}_{(1)}, \dots, \tilde{\beta}_{(K_0)}) - (\beta_0, \dots, \beta_0) = o_p(N^{-1/2})$ where $(\tilde{\beta}_{(1)}, \dots, \tilde{\beta}_{(K_0)})$ is a suitable permutation of $(\tilde{\beta}_1, \dots, \tilde{\beta}_{K_0})$

Remark 1 applies here with obvious modifications. As before, hereafter we simply write $\tilde{\beta}_k$ for $\tilde{\beta}_{(k)}$ as the consistent estimator of β_0 and define $\tilde{\beta}_k \equiv \{\beta \in \{\beta_0, \beta_0\} : \tilde{\beta}_i = \beta\}$ for $i = 1, \dots, N$

³If Conditions (a)-(b) are satisfied and $E \|z_{it} \Delta \varepsilon_{it}\|^{2+\sigma} < \infty$, by the Davydov inequality, we have $\|B_{kNT}\| \leq \frac{1}{T \sqrt{N_k T}} \sum_{i \in G_k^0} \sum_{t=1}^T \sum_{s=1}^T \|E[\Delta x_{is} z_{is}^0 z_{it} \Delta \varepsilon_{it}]\| = O(N/T^{1/2})$, which is $o_p(1)$ if $T \gg N$ and usually asymptotically non-negligible otherwise.

3.3.2 Classification Consistency

Let $\tilde{H}_{nQW} = \{1 \leq J_n^0\}$ and $\tilde{I}_{nQW} = \{1 \leq J_n^0\}$ for $l = 1, \dots, Q$. Let $\tilde{H}_{nQW} = \{1 \leq J_n^0\}$ and $\tilde{I}_{nQW} = \{1 \leq J_n^0\}$. We establish uniform classification consistency in the next theorem.

Theorem 3.2 Suppose that Assumptions B1-B2 hold. Then (i) $\sum_{n=1}^{\infty} P_{N_0}(\tilde{H}_{nQW} \neq \emptyset) < \infty$ and (ii) $\sum_{n=1}^{\infty} P_{N_0}(\tilde{I}_{nQW} \neq \emptyset) < \infty$ as $(Q) \rightarrow \infty$.

REMARK 8. Remark 2 also holds for the above theorem with obvious modifications. Let $\tilde{J}_0 = \{1 \leq J_n^0\}$ and $\tilde{K}_{lQW} = \{1 \leq J_n^0\}$. Theorem 3.2(i) implies that $\sum_{n=1}^{\infty} P_{N_0}(\tilde{H}_{nQW} \neq \emptyset) < \infty$ meaning that all individuals are classified into one of the N_0 groups w.p.a.1.

Let $\tilde{Q}_n = \sum_{l=1}^Q 1\{1 \leq J_n^0\}$. The following corollary parallels Corollary 2.3.

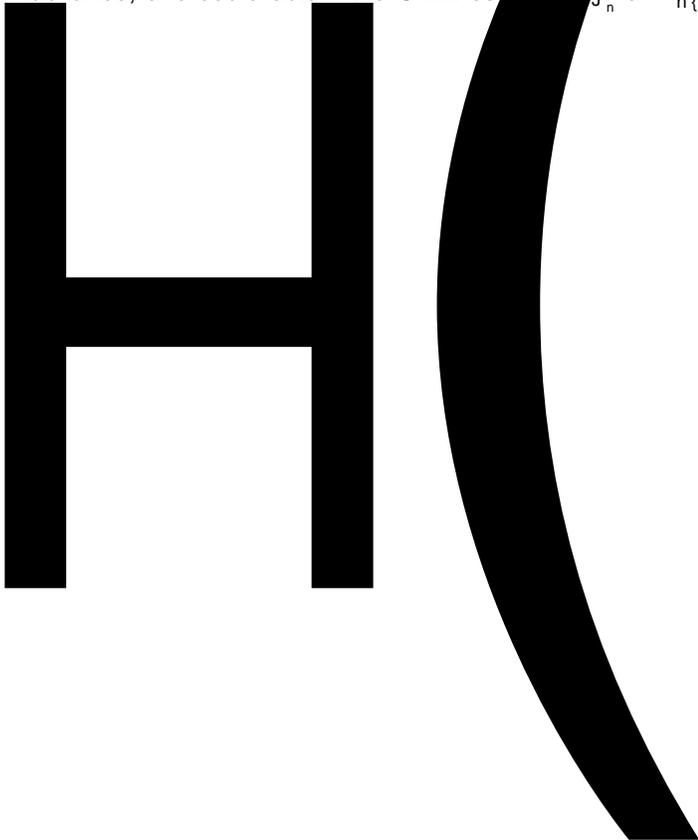
Corollary 3.3 Suppose that Assumptions B1-B2 hold. Then $\tilde{Q}_n \rightarrow Q_n = r_S(1) =$

3.3.3 Improved Convergence and Asymptotic Properties of Post-Lasso

The following theorem establishes the asymptotic distribution of the C-Lasso estimators $\{\tilde{\beta}_n\}$.

Theorem 3.4 Suppose Assumptions B1-B3 hold. Then $\sqrt{n}(\tilde{\beta}_n - \beta) \rightarrow_d N(0, \Sigma)$ for $n \rightarrow \infty$.

REMARK 9. In contrast to the PPL case, the PGMM estimators $\{\tilde{\beta}_n\}$ may fail to possess the oracle property. If the group identities were known in advance, one could obtain the GMM estimator of β .



Suppose Assumptions B1-B4 hold. Then $\sqrt{k}(\tilde{G}_k - \frac{0}{k}) \xrightarrow{D} (0, \Omega_k)$ where $\Omega_k =$

To prove the above theorem, we first apply Theorem 3.2 and show that $\sqrt{k}(\tilde{G}_k - \frac{0}{k}) = \sqrt{k}(G_k^0 - \frac{0}{k}) + o_p(1)$. That is, the post-Lasso GMM estimator \tilde{G}_k is asymptotically equivalent to the oracle estimator G_k^0 . To obtain the most efficient estimator among the class of GMM estimators based on the moment conditions specified in Assumption B1(i), one can set $\frac{(k)}{NT}$ to be a consistent estimator of $\frac{1}{k}$. Alternatively, we can consider Arellano and Bond (1991) GMM estimation based on the estimated groups. The procedure is standard and details are omitted.

If $\frac{(k)}{NT} = \frac{1}{k}$ for each $k \in \mathbb{N}$ in Assumption B3(i) (which is unrealistic before knowing the group identity), and $\frac{(k)}{NT} = 0$ in Assumption B3(ii), then $\frac{(k)}{z\Delta x} = \frac{(k)}{z\Delta x}$, $\frac{(k)}{z\Delta x} = \frac{(k)}{z\Delta x} \Omega_k$ and $\sqrt{k}(\tilde{G}_k - \frac{0}{k}) \xrightarrow{D} (0, \Omega_k)$. Thus in this special case

Table 1: Frequency of selecting $K = 1$ or 5 groups when $K_0 = 3$

N	T			DGP 1				DGP 2				DGP 3		
		1	2	4	5	1	2	4	5	1	2	4	5	
100	15	0	0	0.004	0.002	0	0.232	0.004	0.002					
100	25	0	0	0	0	0	0.016	0	0	0	0.096	0.242	0.016	
100	50	0	0	0	0	0	0	0	0	0	0	0.014	0	
200	15	0	0	0.106	0.004	0	0.022	0.008	0					
200	25	0	0	0	0	0	0	0	0	0	0.106	0.226	0	
200	50	0	0	0	0	0	0	0	0	0	0	1	0	

Next, given the true number of groups, we focus on the classification of individual units and the point estimation of post-Lasso.⁴ Due to space limitation, all tabulated results are produced under $\lambda_j = 0.5$, $\lambda_1 = 1.2$, for the linear models, and $\lambda_1 = 0.05$ for the Probit model. The outcomes are found robust over the specified range of constants. Column 4 of Tables 2 shows the percentage of correct classification of the K_0 units, calculated as $\frac{1}{N} \sum_{k=1}^{K_0} \mathbf{1}_{i \in \hat{G}_k} \{i = k\}$, averaged over the Monte Carlo replications. Columns 5–7 summarize the post-Lasso estimator's root

Table 2: Classification and Point Estimation of α_1

	N	T	% of correct	Post-Lasso			Oracle		
			classification	RMSE	Bias	Coverage	RMSE	Bias	Coverage
DGP 1	100	15	0.8935	0.0594	0.0105	0.8758	0.0463	0.0012	0.9336
	100	25	0.9674	0.0384	0.0018	0.9344	0.0353	0.0001	0.9362
	100	50	0.9964	0.0249	0.0000	0.9528	0.0245	-0.0002	0.9348
	200	15	0.8987	0.0432	0.0077	0.8650	0.0324	-0.0013	0.9410
	200	25	0.9661	0.0272	0.0015	0.9228	0.0250	-0.0006	0.9394
	200	50	0.9966	0.0174	-0.0001	0.9496	0.0171	-0.0002	0.9424
DGP 2	100	15	0.8063	0.0711	-0.0123	0.9562	0.0502	-0.0037	0.9090
	100	25	0.8974	0.0461	-0.0060	0.9760	0.0351	0.0011	0.9336
	100	50	0.9689	0.0278	-0.0011	0.9860	0.0242	-0.0010	0.9320
	200	15	0.8151	0.0557	-0.0159	0.9436	0.0352	-0.0017	0.9308
	200	25	0.9037	0.0328	-0.0047	0.9664	0.0252	-0.0006	0.9442
	200	50	0.9711	0.0193	-0.0014	0.9842	0.0164	0.0000	0.9304
DGP 3	100	25	0.7941	0.1701	0.0805	0.7856	0.1077	0.0114	0.9376
	100	50	0.9456	0.0859	0.0231	0.8970	0.0752	0.0090	0.9504
	200	25	0.8277	0.1325	0.0777	0.7214	0.0821	0.0116	0.9104
	200	50	0.9527	0.0635	0.0223	0.8818	0.0573	0.0121	0.9280

via a dynamic Probit model. Due to space limitation, we only report the estimated coefficients in the main text. Summary statistics, group membership, and additional details of implementation can be found in the Supplementary Material.

f

Understanding the disparate savings behavior across countries is a longstanding research interest in development economics. Theoretical advances and empirical studies have accumulated over many years; see Feldstein (1980), Deaton (1990), Edwards (1996) Bosworth, Collins, and Reinhart (1999), Rodrik (2000), and Li, Zhang, and Zhang (2007), among many others. Empirical research in this area typically employs standard panel data methods to handle heterogeneity or relies on prior information to categorize countries into groups. Classification criteria vary from geographic locations to the notion of developed countries versus developing countries (Loayza, Schmidt-Hebbel and Servén, 2000). This section applies the methodology developed in the present paper to revisit this empirical problem.

Following Edwards (1996), we consider the simple regression model

$$s_{it} = \alpha_1 s_{i,t-1} + \alpha_2 \pi_{it} + \alpha_3 r_{it} + \alpha_4 g_{it} + \alpha_5 \mu_i + \epsilon_{it} \quad (5.1)$$

where s_{it} is the ratio of savings to GDP, π_{it} is the CPI-based inflation rate, r_{it} is the real interest rate, g_{it} is the per capita GDP growth rate, μ_i is a fixed effect, and ϵ_{it} is an idiosyncratic error term. Inflation characterizes the degree of the macroeconomic stability and the real interest rate reflects the price of money. The relationship between the savings rate and GDP growth rate is well documented, with the latter being found to Granger-cause the former (Carroll and Weil, 1994). The first-order lagged savings rate is added to the specification to capture persistence of the savings rate.

Data are obtained from the widely used World Development Indicators, a comprehensive dataset compiled by the World Bank. For many countries the time series of real interest rates are often short in comparison with the other variables. Using the time span 1995–2010, we were able to construct a balanced panel of 56 countries. Substantial heterogeneity across countries was observed in all these major macroeconomic indicators. Evidence of within group homogeneity is therefore particularly important in supporting panel data pooling techniques.

This dynamic panel model can be estimated by either PLS or PGMM. We first try PLS, which has higher correct classification ratio in our simulation when $k = 15$. Following the simulation, λ_{1NT} is set as $\frac{2}{3}(\frac{1}{N})^{-1/2}$, and the IC picks two groups and the tuning parameter constant $\lambda_1 = 1.55$ over all combinations of $k = 1 \dots 5$ and λ_1 in a geometrically increasing sequence of 10 points in $(0.2 \dots 2)$. Based on this choice of tuning parameter, the data determine the group identities. Interestingly, some geographic features remain salient in the classification. For example, we observe a strong collection of Asian countries in Group 1. In particular, except for South Korea and the city state Singapore, Group 1 includes all Eastern Asian and Southeastern Asian countries in our sample, namely, China, Japan, Indonesia, Malaysia, Philippines, and Thailand.

Table 3: PLS and PGMM estimation results

Variables	PLS				PGMM		
	Pooled FE	Group1	Group2	Pooled GMM	Group1	Group2	
Lagged savings	. ***	. ***	. ***	.	.	. **	
Inflation	– .	– . ***	. ***	.	– . **	. ***	
Interest rate	– .	– . ***	. ***	– .	– . **	. *	
GDP growth	. ***	. ***	. **	. ***	. ***	. **	

Note: *** 1% significant, ** 5% significant, * 10% significant

Columns 3–4 in Table 3 report the results for the PLS-based post-Lasso estimation, in comparison with those for the pooled FE estimation in Column 2. The estimates are bias-corrected by the half-panel jackknife (Dhaene and Jochmans, 2015), and the standard errors (in parentheses) are clustered at the country level. Compared with Edwards (1996), the FE results re-confirm the significance of lagged savings and GDP growth rate as well as the insignificance of inflation and interest rates in the determination of savings rate. This result also lends support to the *conventional wisdom*

group, but no such relationship is found in the high-occurrence group.

Table 4: Probit, FE Probit and PPL estimation results

Variables	Probit		FE Probit		Post-Lasso PPL										
					high-occurrence		low-occurrence								
	coef.	s.e.	coef.	s.e.	coef.	s.e.	coef.	s.e.							
Lagged civil war	.	***	0.1156	.	***	0.1140	.	***	0.1363	.	***	0.2707			
GDP per capita growth	-	.	***	0.1155	-	.	***	0.1389	.		0.1193	-	.	***	0.2220
population growth	-	.		0.1107	.		0.1284	-	.		0.1429	.			0.1736

Note: *** 1% significant, ** 5% significant, * 10% significant

We propose a novel and systematic approach to identify and estimate latent group structures in panel data, developing panel penalized profile likelihood (PPL) and panel GMM (PGMM) methods for classification and estimation, and providing asymptotic properties for use in inference. The PPL method enjoys the oracle property but PGMM typically does not. Post-Lasso estimates are also studied and a BIC-type information criterion is proposed to determine the number of groups. These techniques combine to provide a general approach to classifying and estimating panel models with unknown homogeneous groups, heterogeneity across groups, and an unknown number of groups. Simulations show that the approach has good finite sample performance and can be readily implemented in practical work. Two applications reveal the advantages of data-determined identification of latent group structures in empirical panel modeling.

The present work raises interesting issues for further research. First, it may be appealing to consider a more general framework that allows the number (p_0) of groups to grow with the sample size. Close examination of the theory provided in this paper suggests that it is possible to permit p_0 to increase with n but at a very slow rate. Second, both the linear and nonlinear models may be extended to include time effects or interactive fixed effects (IFE). In linear models with IFE but without endogeneity, we remark that the present approach can be used in conjunction with principal component analysis to address cross sectional dependence modeled through IFE. Extension to nonlinear models or to models with endogeneity will raise new statistical and computational challenges. Third, our method can be extended to nonstationary panels where panel unit and cointegrating relationships may possess latent group structures. Some of these topics will be explored in future work.

APPENDIX

(i) Let $l_{NT,i}(\hat{\alpha}_i) = \frac{1}{T} \sum_{t=1}^T l_{it}(\hat{\alpha}_i)$ and $l_{NT,i}^{(K_0)}(\hat{\alpha}_i) = \frac{1}{T} \sum_{t=1}^T l_{it}^{(K_0)}(\hat{\alpha}_i)$. Let $\hat{\alpha}_i = \hat{\alpha}_i - \alpha_i^0$ and $\tilde{\alpha}_i = \hat{\alpha}_i - \alpha_i^0$. Since $\hat{\alpha}_i(\hat{\alpha}_i) = \arg \min_{\mu_i} \frac{1}{T} \sum_{t=1}^T l_{it}(\mu_i)$ we have $\frac{1}{T} \sum_{t=1}^T l_{it}(\hat{\alpha}_i) = 0 \forall \hat{\alpha}_i$. Then by second order Taylor expansion and the envelope theorem, we have

$$\begin{aligned} l_{NT,i}(\hat{\alpha}_i) - l_{NT,i}(\alpha_i^0) &= \frac{1}{T} \sum_{t=1}^T l_{it}(\hat{\alpha}_i) - \frac{1}{T} \sum_{t=1}^T l_{it}(\alpha_i^0) \\ &= \hat{\alpha}_i' \hat{\alpha}_i + \frac{1}{2} \hat{\alpha}_i' \tilde{\alpha}_i \tilde{\alpha}_i' \hat{\alpha}_i \end{aligned} \quad (\text{A.1})$$

where $\tilde{\alpha}_i$

of the property of $\hat{\alpha}$. By (A.1) and the Cauchy-Schwarz inequality

$$\begin{aligned}
 & \left\| \frac{(K_0)}{1NT, \lambda_1} (\beta^0 + -1/2 \mathbf{v} \hat{\alpha}) - \frac{(K_0)}{1NT, \lambda_1} (\beta^0 \alpha^0) \right\|^2 \\
 = & \frac{1}{2} \sum_{i=1}^N \left(\hat{\beta}_i - \beta_i \right)^2 + \frac{\sqrt{N}}{2} \sum_{i=1}^N \left(\hat{\beta}_i - \beta_i \right) + \frac{1}{2} \sum_{i=1}^N \left\| \Pi_{k=1}^{K_0} \hat{\beta}_i - \hat{\beta}_i \right\|^2
 \end{aligned}$$

envelope theorem, the first order condition (with respect to \hat{i}) for the minimization problem in (2.5) yields that

$$\begin{aligned}
\mathbf{0}_{p \times 1} &= \frac{1}{\sqrt{\kappa}} \sum_{t=1}^T \hat{i}(\hat{i}_t; \hat{i}, \hat{i}(\hat{i})) + \sqrt{\kappa} \sum_{j=1}^{K_0} \hat{ij} \prod_{l=1, l \neq j}^{K_0} \|\hat{i} - \hat{l}\| \\
&= \frac{1}{\sqrt{\kappa}} \sum_{t=1}^T \hat{i}_t + \left(\frac{1}{\|\hat{i} - \hat{k}\|} \hat{ki} - \hat{i}\beta\beta \right) \sqrt{\kappa} (\hat{i} - \hat{k}) + \frac{1}{\sqrt{\kappa}} \sum_{t=1}^T \left[\hat{i}(\hat{i}_t; \hat{i}, \hat{i}(\hat{i})) - \hat{i}_t \right] \\
&\quad + \hat{i}\beta\beta \sqrt{\kappa} (\hat{k} - \hat{i}) + \sqrt{\kappa} \sum_{j=1, j \neq k}^{K_0} \hat{ij} \prod_{l=1, l \neq j}^{K_0} \|\hat{i} - \hat{l}\| \\
&\equiv \hat{i}_1 + \hat{i}_2 + \hat{i}_3 + \hat{i}_4 + \hat{i}_5
\end{aligned} \tag{A.9}$$

where $\hat{ij} = \frac{\hat{\beta}_i - \hat{\alpha}_j}{\|\hat{\beta}_i - \hat{\alpha}_j\|}$ if $\|\hat{i} - \hat{j}\| \neq 0$ and $\|\hat{ij}\| \leq 1$ otherwise, the second equality follows from the first order Taylor expansion and rearrangement of terms, $\hat{i}\beta\beta \equiv \hat{i}\beta\beta(\hat{i})$ $\hat{i}\beta\beta(\cdot)$ is defined in (A.2), \hat{i} lies between \hat{i} and \hat{i} elementwise.

Let $\varkappa_{1NT} = (\ln \kappa)^{-1/2} (\ln \kappa)^3 + 1) (\ln \kappa)^\nu$. Let κ denote a generic constant that may vary across lines. By (A.4) and Lemmas S1.6-S1.7 in the Supplement, we can readily show that

$$\left(\max_{\hat{i}} \|\hat{i} - \hat{i}\| \geq \varkappa_{1NT} \right) = (\kappa^{-1}) \text{ for some } \kappa > 0 \tag{A.10}$$

which in conjunction with the proof of Theorem 2.1(iii), implies that

$$\left(\sqrt{\kappa} \|\hat{k} - \hat{i}\| \geq (\ln \kappa)^\nu \right) = (\kappa^{-1}) \text{ and } \left(\max_{\hat{i} \in G_k^0} \|\hat{ki} - \hat{i}\| \geq \frac{0}{\kappa} 2 \right) = (\kappa^{-1}) \tag{A.11}$$

By (A.10)-(A.11), $\left(\max_{\hat{i} \in G_k^0} \|\hat{i}_5\| \geq \sqrt{\kappa} \varkappa_{1NT} \right) = (\kappa^{-1})$. Combining these results with those in Lemmas S1.6(v) and S1.11(i), we have $(\Xi_{kNT}) = 1 - (\kappa^{-1})$ where

$$\begin{aligned}
\Xi_{kNT} &\equiv \left\{ \max_{\hat{i} \in G_k^0} \|\hat{ki} - \hat{i}\| \leq \frac{0}{\kappa} 2 \right\} \cap \left\{ \max_{\hat{i} \in G_k^0} \|\hat{i}\beta\beta - \hat{i}\beta\beta(\hat{i})\| \leq H 2 \right\} \cap \left\{ \max_{\hat{i} \in G_k^0} \|\hat{i}_3\| \leq (\ln \kappa)^{3+\nu} \right\} \\
&\quad \cap \left\{ \max_{\hat{i} \in G_k^0} \|\hat{i}_4\| \leq (\ln \kappa)^\nu \right\} \cap \left\{ \max_{\hat{i} \in G_k^0} \|\hat{i}_5\| \leq \sqrt{\kappa} \varkappa_{1NT} \right\}
\end{aligned}$$

Then conditional on Ξ_{kNT} we have uniformly in $\hat{i} \in G_k^0$

$$\begin{aligned}
\|(\hat{i} - \hat{k})'(\hat{i}_2 + \hat{i}_3 + \hat{i}_4 + \hat{i}_5)\| &\geq \|(\hat{i} - \hat{k})' \hat{i}_2\| - \|(\hat{i} - \hat{k})'(\hat{i}_3 + \hat{i}_4 + \hat{i}_5)\| \\
&\geq \sqrt{\kappa} \|\hat{ki}\| \|\hat{i} - \hat{k}\| - \|\hat{i} - \hat{k}\|^2 2 (\ln \kappa)^{3+\nu} + \sqrt{\kappa} \varkappa_{1NT} \\
&\geq \sqrt{\kappa} \|\hat{i}\beta\beta\| \|\hat{i} - \hat{k}\| 4 \text{ for sufficiently large } (\kappa)
\end{aligned}$$

where the last inequality follows because $\sqrt{\kappa} \|\hat{i}\beta\beta\| \gg 2 (\ln \kappa)^{3+\nu} + \sqrt{\kappa} \varkappa_{1NT}$ by Assumption A2(i). Then for all $\hat{i} \in G_k^0$ we have

$$\begin{aligned}
(\hat{k}_{NT,i}) &= \left(\hat{i} \in G_k^0 \mid \hat{i} \in G_k^0 \right) = \left(-\hat{i}_1 = \hat{i}_2 + \hat{i}_3 + \hat{i}_4 + \hat{i}_5 \right) \\
&\leq \left(\left| (\hat{i} - \hat{k})' \hat{i}_1 \right| \geq \left| (\hat{i} - \hat{k})'(\hat{i}_2 + \hat{i}_3 + \hat{i}_4 + \hat{i}_5) \right| \right) \\
&\leq \left(\|\hat{i}_1\| \geq \sqrt{\kappa} \|\hat{i}\beta\beta\| 4 \Xi_{kNT} \right) + (\Xi_{kNT}^c) \rightarrow 0 \text{ as } (\kappa) \rightarrow \infty
\end{aligned}$$

where Ξ_{kNT}^c denotes the complement of Ξ_{kNT} and the convergence follows by Lemma S1.6(iv) and Assumption A2. Consequently, we conclude that with probability $1 - o_p(1)$ the difference $\|\hat{\mu}_i - \hat{\mu}_k\|$ must reach the point where $\|\hat{\mu}_i - \hat{\mu}_k\|$ is not differentiable with respect to $\hat{\mu}_i$ for any $i \in G_k^0$. That is $(\|\hat{\mu}_i - \hat{\mu}_k\| = 0 \mid i \in G_k^0) = 1 - o_p(1)$.

For uniform consistency, we have: $(\bigcup_{k=1}^{K_0} \hat{\mu}_{kNT}) \leq \max_{k=1}^{K_0} (\hat{\mu}_{kNT}) \leq \max_{k=1}^{K_0} \max_{i \in G_k^0} (\hat{\mu}_{kNT,i})$ and by Lemma S1.6(iv)

$$\begin{aligned} \max_{k=1}^{K_0} \max_{i \in G_k^0} (\hat{\mu}_{kNT,i}) &\leq \max_{k=1}^{K_0} \max_{i \in G_k^0} (\|\hat{\mu}_i\| \geq \sqrt{1 - \frac{1}{4} \Xi_{kNT}}) + (\Xi_{kNT}^c) \\ &\leq \max_{1 \leq i \leq N} \left(\frac{1}{\sqrt{t}} \|\mathbf{1}\| \geq \frac{1}{\sqrt{t}} \frac{K_0 - 1}{4} \right) + (1) = (1) \end{aligned} \quad (\text{A.12})$$

This completes the proof of (i).

(ii) Pretending each individual's membership is random, we have $(i \in G_k^0) = \mu_k \rightarrow \mu_k \in (0, 1)$ for $\mu_k = 1 - \mu_0$ and can interpret previous results as conditional on the group membership assignment. By Bayes theorem,

$$\begin{aligned} (\hat{\mu}_{kNT,i}) &= 1 - (i \in G_k^0 \mid i \in \hat{G}_k) \\ &= \frac{\prod_{l=1, l \neq k}^{K_0} (\mu_l \mid i \in G_l^0) \prod_{l=1}^{K_0} (\mu_l)}{(\mu_k \mid i \in G_k^0) \prod_{l=1}^{K_0} (\mu_l) + \prod_{l=1, l \neq k}^{K_0} (\mu_l \mid i \in G_l^0) \prod_{l=1}^{K_0} (\mu_l)} \end{aligned} \quad (\text{A.13})$$

For the numerator, we have by (A.12)

$$\prod_{l=1, l \neq k}^{K_0} (\mu_l \mid i \in G_l^0) \prod_{l=1}^{K_0} (\mu_l) \leq (1 - \mu_0) \prod_{l=1}^{K_0} (\mu_l) = (1)$$

In addition, noting that $(i \in G_k^0 \mid i \in G_k^0) = 1 - (i \in G_k^0 \mid i \in G_l^0) = 1 - (1)$ uniformly in i and by (i), we have that $(i \in G_k^0 \mid i \in G_k^0) \prod_{l=1}^{K_0} (\mu_l) + \prod_{l=1, l \neq k}^{K_0} (\mu_l \mid i \in G_l^0) \prod_{l=1}^{K_0} (\mu_l) \geq (i \in G_k^0) \frac{1}{2}$ w.p.a.1. It follows that

$$\begin{aligned} (\bigcup_{k=1}^{K_0} \hat{\mu}_{kNT}) &\leq \max_{k=1}^{K_0} \max_{i \in \hat{G}_k} (\hat{\mu}_{kNT,i}) \leq \frac{\prod_{l=1, l \neq k}^{K_0} (\mu_l \mid i \in G_l^0) \prod_{l=1}^{K_0} (\mu_l)}{\min_{1 \leq i \leq N} \min_{1 \leq k \leq K_0} (i \in G_k^0) \frac{1}{2}} \\ &= \frac{(1)}{\min_{1 \leq k \leq K_0} \mu_k \frac{1}{2}} = (1) \quad \blacksquare \end{aligned}$$

. Noting that $\hat{\mu}_k = \frac{1}{N} \sum_{i=1}^N \mathbf{1}\{i \in \hat{G}_k\}$, $\mu_k = \frac{1}{N} \sum_{i=1}^N \mathbf{1}\{i \in G_k^0\}$ and $\mathbf{1}\{i \in \hat{G}_k\} - \mathbf{1}\{i \in G_k^0\} = \mathbf{1}\{i \in \hat{G}_k \setminus G_k^0\} - \mathbf{1}\{i \in G_k^0 \setminus \hat{G}_k\}$ we have $\|\hat{\mu}_k - \mu_k\| = \frac{1}{N} \sum_{i=1}^N |\mathbf{1}\{i \in \hat{G}_k \setminus G_k^0\} - \mathbf{1}\{i \in G_k^0 \setminus \hat{G}_k\}|$. Then by the implication rule and the Markov inequality, for any $\epsilon > 0$

$$\begin{aligned} (\|\hat{\mu}_k - \mu_k\| \geq \epsilon) &\leq \left(\frac{1}{N} \sum_{i=1}^N \mathbf{1}\{i \in \hat{G}_k \setminus G_k^0\} \geq \epsilon \right) + \left(\frac{1}{N} \sum_{i=1}^N \mathbf{1}\{i \in G_k^0 \setminus \hat{G}_k\} \geq \epsilon \right) \\ &= \frac{1}{N} \sum_{i=1}^N (\hat{\mu}_{kNT,i}) + \frac{1}{N} \sum_{i=1}^N (\mu_{kNT,i}) \end{aligned}$$

By (A.12), $\prod_{i=1}^N (\hat{k}_{NT,i}) = \prod_{k=1}^{K_0} \prod_{i \in G_k^0} (\hat{k}_{NT,i}) = (1)$ By the proof of Theorem 2.2(i), $\prod_{i=1}^N (\hat{k}_{NT,i})$
 $= \prod_{k=1}^{K_0} \prod_{i \in \hat{\mathcal{R}}}$

where $\hat{\mathcal{R}}$

We start by proving a useful technical result and then proceed to prove the main results. Let $\bar{y}_{iNT}(\bar{y}_i) \equiv [\frac{1}{T} \sum_{t=1}^T (y_{it} - \bar{y}_i)]'$ $\bar{y}_{iNT} [\frac{1}{T} \sum_{t=1}^T (y_{it} - \bar{y}_i)]$ and $\bar{y}_i(\bar{y}_i) \equiv \{\frac{1}{T} \sum_{t=1}^T \mathbb{E} [$

where $\mathbb{1}_{NT} \equiv \min_{1 \leq i \leq N} \min_{i \in \mathcal{Z}_{\Delta x}} \left(\frac{1}{i} - \frac{1}{i, z_{\Delta x}} \right)$ satisfies that $\liminf_{(N,T) \rightarrow \infty} \mathbb{1}_{NT} \geq \frac{2}{W-Q} > 0$ by Assumptions B1(iii)-(iv). Combining (B.4)-(B.6) yields

$$\mathbb{1}_{NT} \left\| \tilde{\mathbf{y}}_i - \frac{0}{i} \right\|^2 \leq \frac{2}{\underline{c}} \left[2 \frac{0}{i, T} + \frac{1}{\tilde{K}_0} \left(\left\| \tilde{\mathbf{y}}_i - \frac{0}{i} \right\|^2 + 2 \left\| \tilde{\mathbf{y}}_i - \frac{0}{i} \right\|^2 \right) \right]$$

or, $\left(\mathbb{1}_{NT} - \frac{4}{\underline{c}} \frac{1}{\tilde{K}_0} \right) \left\| \tilde{\mathbf{y}}_i - \frac{0}{i} \right\|^2 \leq \frac{2}{\underline{c}} \left(2 \frac{0}{i, T} + \frac{1}{\tilde{K}_0} \left\| \tilde{\mathbf{y}}_i - \frac{0}{i} \right\|^2 \right)$ where $\tilde{K}_0 = K_0(\tilde{\boldsymbol{\alpha}})$. Then

$$\begin{aligned} \left\| \tilde{\mathbf{y}}_i - \frac{0}{i} \right\| &\leq \frac{\frac{2}{\underline{c}} \frac{1}{\tilde{K}_0} + \left[\left(\frac{2}{\underline{c}} \frac{1}{\tilde{K}_0} \right)^2 + \frac{8}{\underline{c}} \left(\mathbb{1}_{NT} - \frac{4}{\underline{c}} \frac{1}{\tilde{K}_0} \right) \left(2 \frac{0}{i, T} + \frac{1}{\tilde{K}_0} \left\| \tilde{\mathbf{y}}_i - \frac{0}{i} \right\|^2 \right) \right]^{1/2}}{2 \left(\mathbb{1}_{NT} - \frac{4}{\underline{c}} \frac{1}{\tilde{K}_0} \right)} \\ &= P \left(-1/2 + \frac{1}{2} \right) \end{aligned} \quad (\text{B.7})$$

As in the proof of Theorem 2.1(ii), we can further demonstrate that $\frac{1}{N} \sum_{i=1}^N \left\| \tilde{\mathbf{y}}_i - \frac{0}{i} \right\|^2 = P \left(-1 \right)$

The proof of (iii) is completely analogous to that of Theorem 2.1(iii), now using the facts that $\left| \mathbb{1}_{NT}(\tilde{\boldsymbol{\beta}} | \boldsymbol{\alpha}) - \mathbb{1}_{NT}(\boldsymbol{\beta}^0 | \boldsymbol{\alpha}^0) \right| = P \left(-1/2 \right)$ and that $0 \geq \mathbb{1}_{NT}(\tilde{\boldsymbol{\beta}} | \tilde{\boldsymbol{\alpha}}) - \mathbb{1}_{NT}(\tilde{\boldsymbol{\beta}} | \boldsymbol{\alpha}^0)$ ■

(i) First, we fix $i \in \{1, \dots, K_0\}$. By the consistency of $\tilde{\mathbf{y}}_k$ and $\tilde{\mathbf{y}}_i$ in Theorem 3.1 and Assumptions B1(v)-(vi), we have $\tilde{\mathbf{y}}_i - \tilde{\mathbf{y}}_l \xrightarrow{P} \frac{0}{k} - \frac{0}{l} \neq 0$ for all $i \in \frac{0}{k}$ and $i \neq l$ and $\tilde{\mathbf{y}}_{ki} \equiv \prod_{l=1, l \neq k}^{K_0} \left\| \tilde{\mathbf{y}}_i - \tilde{\mathbf{y}}_l \right\| \xrightarrow{P} 0 \equiv \prod_{l=1, l \neq k}^{K_0} \left\| \frac{0}{k} - \frac{0}{l} \right\| \geq \frac{K_0 - 1}{\alpha} > 0$ for any $i \in \frac{0}{k}$. Now, suppose that $\left\| \tilde{\mathbf{y}}_i - \tilde{\mathbf{y}}_k \right\| \neq 0$ for some $i \in \frac{0}{k}$. Then the first order condition (with respect to $\tilde{\mathbf{y}}_i$) for the minimization problem in (3.2) implies that

$$\begin{aligned} \mathbf{0}_{p \times 1} &= -2 \frac{1}{i, z_{\Delta x}} \frac{1}{iNT} \frac{1}{\sqrt{t=1}^{\square T}} \left(\Delta_{it} - \tilde{\mathbf{y}}_i' \Delta_{it} \right) + \sqrt{2} \sum_{j=1}^{K_0} \tilde{\mathbf{y}}_{ij} \prod_{l=1, l \neq j}^{K_0} \left\| \tilde{\mathbf{y}}_i - \tilde{\mathbf{y}}_l \right\| \\ &= -2 \frac{1}{i, z_{\Delta x}} \frac{1}{iNT} \frac{1}{\sqrt{t=1}^{\square T}} \left(\Delta_{it} \right) + \left\{ \frac{2 \tilde{\mathbf{y}}_{ki}}{\left\| \tilde{\mathbf{y}}_i - \tilde{\mathbf{y}}_k \right\|} p + 2 \frac{1}{i, z_{\Delta x}} \frac{1}{iNT} \tilde{\mathbf{y}}_{i, z_{\Delta x}} \right\} \sqrt{2} \left(\tilde{\mathbf{y}}_i - \tilde{\mathbf{y}}_k \right) \\ &\quad + 2 \frac{1}{i, z_{\Delta x}} \frac{1}{iNT} \tilde{\mathbf{y}}_{i, z_{\Delta x}} \sqrt{2} \left(\tilde{\mathbf{y}}_k - \frac{0}{k} \right) + \sqrt{2} \sum_{j=1, j \neq k}^{K_0} \tilde{\mathbf{y}}_{ij} \prod_{l=1, l \neq j}^{K_0} \left\| \tilde{\mathbf{y}}_i - \tilde{\mathbf{y}}_l \right\| \\ &\equiv -\tilde{\mathbf{y}}_{i1} + \tilde{\mathbf{y}}_{i2} + \tilde{\mathbf{y}}_{i3} + \tilde{\mathbf{y}}_{i4} \text{ say,} \end{aligned} \quad (\text{B.8})$$

where $\tilde{\mathbf{y}}_{ij} = \frac{\tilde{\beta}_i - \tilde{\alpha}_j}{\left\| \tilde{\beta}_i - \tilde{\alpha}_j \right\|}$ if $\left\| \tilde{\mathbf{y}}_i - \tilde{\mathbf{y}}_j \right\| \neq 0$ and $\left\| \tilde{\mathbf{y}}_{ij} \right\| \leq 1$ if $\left\| \tilde{\mathbf{y}}_i - \tilde{\mathbf{y}}_j \right\| = 0$. Following the proof of Lemma S1.7, we can show that $\left(\max_i \left\| \tilde{\mathbf{y}}_i - \frac{0}{i} \right\| \geq \frac{1}{\kappa_{2NT}} \right) = P \left(-1 \right)$ for any given $\frac{1}{\kappa_{2NT}} > 0$. With this, by (B.7) and Assumptions B2(ii)-(iv), we can readily show that

$$\left(\max_i \left\| \tilde{\mathbf{y}}_i - \frac{0}{i} \right\| \geq \frac{1}{\kappa_{2NT}} \right) = P \left(-1 \right) \text{ for some } \frac{1}{\kappa_{2NT}} > 0 \quad (\text{B.9})$$

where $\kappa_{2NT} = \left(-1/2 (\ln \frac{1}{\underline{c}}) + \frac{1}{2} \right) (\ln \frac{1}{\underline{c}})^\nu$. This, in conjunction with the proof of Theorem 3.1(iii), implies that

$$\left(\sqrt{2} \left\| \tilde{\mathbf{y}}_k - \frac{0}{k} \right\| \geq (\ln \frac{1}{\underline{c}})^\nu \right) = P \left(-1 \right) \text{ and } \left(\max_{i \in \frac{0}{k}} \tilde{\mathbf{y}}_{ki} - \frac{0}{k} \geq \frac{0}{k} \frac{1}{2} \right) = P \left(-1 \right) \quad (\text{B.10})$$

By (B.9)-(B.10), $(\max_{i \in G_k^0} \|\tilde{i}_4\| \geq \sqrt{2\mathcal{X}_{2NT}}) = (\cdot^{-1})$. By Assumptions B1(iii)-(iv), we have $(\max_{i \in G_k^0} \|\tilde{i}'_{i,z\Delta x} - \tilde{i}'_{i,z\Delta x} - \tilde{i}'_{i,z\Delta x}\| \geq \cdot) = (\cdot^{-1})$ for any $\cdot > 0$. This result, in conjunction with (B.10), implies that $(\max_{i \in G_k^0} \|\tilde{i}_3\| \geq (\ln \cdot)^\nu) = (\cdot^{-1})$ for some $\cdot > 0$. It follows that $(\Gamma_{kNT}) = 1 - (\cdot^{-1})$ where

$$\Gamma_{kNT} \equiv \left\{ \max_{i \in G_k^0} \|\tilde{i}_k - \tilde{i}_k\| \leq \frac{0}{k} \right\} \cap \left\{ \max_{i \in G_k^0} \|\tilde{i}_{NT} - \tilde{i}\| \leq \frac{0}{W} \right\} \cap \left\{ \max_{i \in G_k^0} \|\tilde{i}_{i,z\Delta x} - \tilde{i}_{i,z\Delta x}\| \leq \frac{0}{Q} \right\} \\ \cap \left\{ \max_{i \in G_k^0} \|\tilde{i}_3\| \leq (\ln \cdot)^\nu \right\} \cap \left\{ \max_{i \in G_k^0} \|\tilde{i}_4\| \leq \sqrt{2\mathcal{X}_{2NT}} \right\}$$

Then conditional on Γ_{kNT} we have uniformly in $\cdot \in \frac{0}{k}$

$$\begin{aligned} (\tilde{i} - \tilde{k})' (\tilde{i}_2 + \tilde{i}_3 + \tilde{i}_4) &\geq \|\tilde{i} - \tilde{k}\| \|\tilde{i}_2\| - \|\tilde{i} - \tilde{k}\| (\|\tilde{i}_3\| + \|\tilde{i}_4\|) \\ &\geq \sqrt{2} \|\tilde{i} - \tilde{k}\| - \|\tilde{i} - \tilde{k}\| (\ln \cdot)^\nu + \sqrt{2\mathcal{X}_{2NT}} \\ &\geq \sqrt{2} \|\tilde{i} - \tilde{k}\| \frac{0}{k} \geq 4 \text{ for sufficiently large } (\cdot) \end{aligned}$$

because $\sqrt{2} \gg (\ln \cdot)^\nu + \sqrt{2\mathcal{X}_{2NT}}$ by Assumption B2(i). Then by Assumptions B2(i)-(ii)

$$\begin{aligned} (\hat{k}_{NT,i}) &= (\cdot \in \tilde{k} \mid \cdot \in \frac{0}{k}) = (\tilde{i}_1 = \tilde{i}_2 + \tilde{i}_3 + \tilde{i}_4) \\ &\leq \left(\|\tilde{i} - \tilde{k}\| \geq \|\tilde{i}_1\| \geq \|\tilde{i} - \tilde{k}\| (\|\tilde{i}_2 + \tilde{i}_3 + \tilde{i}_4\|) \right) \\ &\leq \left(\|\tilde{i}_1\| \geq \sqrt{2} \frac{0}{k} \geq 4 \Gamma_{kNT} \right) + (\Gamma_{kNT}^c) \rightarrow 0 \text{ as } (\cdot) \rightarrow \infty \end{aligned}$$

It follows that $(\|\tilde{i} - \tilde{k}\| = 0 \mid \cdot \in \frac{0}{k}) \rightarrow 1$ as $(\cdot) \rightarrow \infty$. Now, observe that $(\cup_{k=1}^{K_0} \hat{k}_{NT}) \leq \frac{K_0}{k=1} (\hat{k}_{NT}) \leq \frac{K_0}{k=1} \frac{0}{k} (\hat{k}_{NT,i})$ and by Assumption B2(ii)

$$\begin{aligned} \frac{K_0}{k=1} (\hat{k}_{NT,i}) &\leq \frac{K_0}{k=1} \frac{0}{k} \left(\|\tilde{i}_1\| \geq \sqrt{2} \frac{0}{k} \geq 4 \Gamma_{kNT} \right) + (\Gamma_{kNT}^c) \\ &\leq \max_{1 \leq i \leq N} \left(\|\tilde{i}'_{i,z\Delta x} - \tilde{i}'_{i,z\Delta x} - \tilde{i}'_{i,z\Delta x}\| \geq \frac{0}{k} \geq 4 \Gamma_{kNT} \right) + (1) \\ &\leq \max_{1 \leq i \leq N} \left(\frac{1}{t=1} \|\tilde{i}'_{i,z\Delta x} - \tilde{i}'_{i,z\Delta x} - \tilde{i}'_{i,z\Delta x}\| \geq 2 \frac{K_0-1}{\alpha} (16_{-Q-W}) \right) + (1) = (1) \end{aligned}$$

where we use the fact that $\|\tilde{i}_{i,z\Delta x}\| \|\tilde{i}_{NT}\| \geq (\|\tilde{i}_{i,z\Delta x} - \tilde{i}_{i,z\Delta x} - \tilde{i}_{i,z\Delta x}\|) (\|\tilde{i}\| - \|\tilde{i}_{NT} - \tilde{i}\|) \geq \frac{0}{Q-W} \geq 4$ on the set Γ_{kNT} . Consequently, we have shown (i).

(ii) The proof of (i) is almost identical to that of Theorem 2.2(ii) and is omitted. \blacksquare

The proof follows closely from that of Theorem 2.4. Based on the subdifferential calculus, the KKT conditions for the minimization of (3.2) are that for each $i = 1, \dots, N$ and $j = 1, \dots, K_0$

$$\begin{aligned} \mathbf{0}_{p \times 1} &= -2 \tilde{i}'_{i,z\Delta x} - \tilde{i}'_{i,z\Delta x} - \tilde{i}'_{i,z\Delta x} \frac{1}{t=1} \|\tilde{i}'_{i,z\Delta x} - \tilde{i}'_{i,z\Delta x} - \tilde{i}'_{i,z\Delta x}\| + \frac{2}{j=1} \prod_{l=1, l \neq j}^{K_0} \|\tilde{i} - \tilde{l}\| \text{ and} \\ \mathbf{0}_{p \times 1} &= \frac{1}{i=1} \prod_{l=1, l \neq k}^{K_0} \|\tilde{i} - \tilde{l}\| \end{aligned}$$

where \tilde{y}_{ij} is defined after (B.8). Fix $k \in \{1, \dots, p\}$. As in the proof of Theorem 2.4, we can show that $\frac{2}{NT} \sum_{i \in \tilde{G}_k} \tilde{y}'_{i,z\Delta x} iNT \sum_{t=1}^T \tilde{\Delta}_{it} - \tilde{y}'_k \Delta_{it} + \frac{\lambda_2}{N} \sum_{i \in \tilde{G}_0} \tilde{y}'_{ik} \Pi_{l=1, l \neq k}^{K_0} \|\tilde{y}_i - \tilde{y}_l\| = \mathbf{0}_{p \times 1}$ w.p.a.1. It follows that $\tilde{y}_k = \tilde{y}_{1k} + \tilde{\mathcal{R}}_k$ where $\tilde{y}_{1k} = \left(\frac{1}{N} \sum_{i \in \tilde{G}_k} \tilde{y}'_{i,z\Delta x} iNT \sum_{t=1}^T \tilde{\Delta}_{it} \right)^{-1} \times \frac{1}{NT} \sum_{i \in \tilde{G}_k} \tilde{y}'_{i,z\Delta x} iNT \sum_{t=1}^T \tilde{\Delta}_{it}$ and $\tilde{\mathcal{R}}_k = \left(\frac{1}{N} \sum_{i \in \tilde{G}_k} \tilde{y}'_{i,z\Delta x} iNT \sum_{t=1}^T \tilde{\Delta}_{it} \right)^{-1} \frac{\lambda_2}{2N} \sum_{i \in \tilde{G}_0} \tilde{y}'_{ik} \Pi_{l=1, l \neq k}^{K_0} \|\tilde{y}_i - \tilde{y}_l\|$. By Theorem 3.2, we can readily show that $\left(\sqrt{\frac{\|\tilde{\mathcal{R}}_k\|}{\|\tilde{y}_{1k}\|}} \geq \epsilon \right) = o_p(1)$ for any $\epsilon > 0$, and

$$\frac{1}{k} \sum_{i \in G_k^0} \tilde{y}'_{i,z\Delta x} W_{iNT}$$

- BROWNING, M., AND J. M. CARRO (2007): “Heterogeneity and Microeconometrics Modelling,” In R. Blundell, W. K. Newey and T. Persson (Eds.), *Advances in Economics and Econometrics, Theory and Applications: Ninth World Congress of the Econometric Society*, Volume 3, pp. 45-74. Cambridge University Press, New York.
- BROWNING, M., AND J. M. CARRO (2010): “Heterogeneity in Dynamic Discrete Choice Models,” *Econometrics Journal* 13, 1-39.
- BROWNING, M., AND J. M. CARRO (2014): “Dynamic Binary Outcome Models with Maximal Heterogeneity,” *Journal of Econometrics* 178, 805-823.
- CARROLL, C., AND D.N. WEIL (1994): “Saving and Growth: a Reinterpretation,” *Carnegie-Rochester Conference Series on Public Policy* 40, 133-192
- CHAN, N.H., C.Y. YAU, AND R-M. ZHANG (2014): “Group Lasso for Structural Break Time Series,” *Journal of the American Statistical Association* 109, 590-599.
- COLLIER, P. AND A. HOEFFLER (2004): “Greed and Grievance in Civil Wars,” *Oxford Economic Papers* 56, 563-595.
- DEATON, A. (1990): “Saving in Developing Countries: Theory and Review,” *Proceedings of the World Bank Annual Conference on Development Economics* 61-96, The World Bank, Washington, DC.
- DHAENE, G., AND K. JOCHMANS (2015): “Split-panel Jackknife Estimation of Fixed-effect Models,” *Review of Economic Studies* 82, 991-1030.
- DJANKOV, S. AND M. REYNAL-QUEROL (2010): “Poverty and Civil War: Revisiting the Evidence,” *Review of Economics and Statistics* 92, 1035-1041.
- DURLAUF, S. N., P. A. JOHNSON, AND J. R. TEMPLE (2005): “Growth Econometric,” in P. Aghion and S. Durlauf (eds), *Handbook of Economic Growth*, Vol 1, pp. 555-677. Amsterdam: Elsevier.
- EDWARDS, S (1996): “Why Are Latin America’s Savings Rates So Low? An International Comparative Analysis,” *Journal of Development Economics* 51, 5-44.
- ESTEBAN, J., L. MAYORAL AND D. RAY (2012): “Ethnicity and Conflict: An Empirical Study,” *American Economic Review* 102, 1310-1342.
- FEARON, J.D. AND D. D. LAITIN (2003): “Ethnicity, Insurgency, and Civil War,” *The American Political Science Review* 97, 75-90.
- FELDSTEIN, M. (1980): “International Differences in Social Security and Saving,” *Journal of Public Economics* 14, 225-244.
- FERNÁNDEZ-VAL, I., AND L. LEE (2013): “Panel Data Models with Nonadditive Unobservable Heterogeneity: Estimation and Inference,” *Quantitative Economics* 4, 453-481.
- HAHN, J., AND G. KUERSTEINER (2002): “Asymptotically Unbiased Inference for a Dynamic Panel Model with Fixed Effects When Both N and T are Large,” *Econometrica* 70, 1639-1657.
- HAHN, J., AND G. KUERSTEINER (2011): “Bias Reduction for Dynamic Nonlinear Panel Models with Fixed Effects,” *Econometric Theory* 27, 1152-1191.
- HAHN, J., AND H.R. MOON (2010): “Panel Data Models with Finite Number of Multiple Equilibria,” *Econometric Theory* 26, 863-881.
- HAHN, J., AND W. NEWAY (2004): “Jackknife and Analytical Bias Reduction for Nonlinear Panel Models,” *Econometrica* 72, 1295-1319.
- HARCHAOUI, Z., AND C. LÉVY-LEDUC (2010): “Multiple Change-point Estimation with a Total Variation Penalty,” *Journal of the American Statistical Association* 105, 1481-1493.
- HSIAO, C. (2014): *Analysis of Panel Data*, 3rd edition. Cambridge University Press, New York.

- HSIAO, C., AND H. PESARAN (2008): “Random Coefficient Panel Data Models,” in L. Matyas and P. Sevestre (Eds.), *The Econometrics of Panel Data: Fundamentals and Recent Developments in Theory and Practice*, pp. 187-216, 3rd Edition. Springer-Verlag, Berlin.
- HSIAO, C., AND A. K. TAHMISCIOGLU (1997): “A Panel Analysis of Liquidity Constraints and Firm Investment,” *Journal of the American Statistical Association* 92, 455-465.
- KASAHARA, H., AND K. SHIMOTSU (2009): “Nonparametric Identification of Finite Mixture Models of Dynamic Discrete Choices,” *Econometrica* 77, 135-175.
- KATO, K., A.F. GAVAO, AND G.V. MONTES-ROJAS (2012): “Asymptotics for Panel Quantile Regression Models with Individual Effects,” *Journal of Econometrics* 170, 76-91.
- KIVIET, J. F. (1995): “On Bias, Inconsistency, and Efficiency of Various Estimators in Dynamic Panel Data Models,” *Journal of Econometrics* 68, 53-78.
- LEE, K., M. H. PESARAN, AND R. SMITH (1997): “Growth and Convergence in a Multi-country Empirical Stochastic Growth Model,” *Journal of Applied Econometrics* 12, 357-392.
- LEE, Y. (2012): “Bias in Dynamic Panel Models under Time Series Misspecification,” *Journal of Econometrics* 169, 54-60.
- LEE, Y., AND P.C.B. PHILLIPS (2015): “Model Selection in the Presence of Incidental Parameters,” *Journal of Econometrics* 188, 474-489.
- LEEB, H., AND P.M. PÖTSCHER (2008): “Sparse Estimators and the Oracle Property, or the Return of Hodges’ Estimator,” *Journal of Econometrics* 142, 201-211.
- LEEB, H., AND P.M. PÖTSCHER (2009): “On the Distribution of Penalized Maximum Likelihood Estimators: the LASSO, SCAD, and Thresholding,” *Journal of Multivariate Analysis* 100, 2065-2082.
- LI, H., J. ZHANG, AND J. ZHANG (2007): “Effects of Longevity and Dependency Rates on Saving and Growth,” *Journal of Development Economics* 84, 138-154.
- LIAO, Z. (2013): “Adaptive GMM Shrinkage Estimation with Consistent Moment Selection,” *Econometric Theory* 29, 857-904.
- LIN, C-C., AND S. Ng (2012): “Estimation of Panel Data Models with Parameter Heterogeneity when Group Membership Is Unknown,” *Journal of Econometric Methods* 1, 42-55.
- LOAYZA, N., K. SCHMIDT-HEBBEL, AND L. SERVÉN (2000): “Saving in Developing Countries: An Overview,” *The World Bank Economic Review* 14, 393-414.
- LU, X., AND L. SU (2016): “Shrinkage Estimation of Dynamic Panel Data Models with Interactive Fixed Effects,” *Journal of Econometrics* 190, 148-175.
- MIGUEL, E., S. SATYANATH AND E. SERGENTI (2004): “Economic Shocks and Civil Conflict: an Instrumental Variables Approach,” *Journal of Political Economy* 112, 725-753.
- NUNN, N. AND N. QIAN (2014): “US Food Aid and Civil Conflict,” *American Economic Review* 104, 1630-1666.
- PESARAN, H., Y. SHIN, AND R. SMITH (1999): “Pooled Mean Group Estimation of Dynamic Heterogeneous Panels,” *Journal of the American Statistical Association* 94, 621-634.
- PHILLIPS, P. C. B., AND D. SUL (2007a): “Transition Modeling and Econometric Convergence Tests,” *Econometrica* 75, 1771-1855.
- PHILLIPS, P. C. B., AND D. SUL (2007b): “Bias in Dynamic Panel Estimation with Fixed Effects, Incidental Trends and Cross Section Dependence,” *Journal of Econometrics* 137, 162-188.
- QIAN, J., AND L. SU (2015): “Shrinkage Estimation of Regression Models with Multiple Structural Changes,” *Econometric Theory*, forthcoming.

RODRIK, D (2000): "Saving Transitions," *The World Bank Economic Review* 14, 481-507.

SARAFIDIS, V., AND N. WEBER (2015): "A Partially Heterogenous Framework for Analyzing Panel Data," *Oxford Bulletin of Economics and Statistics* 77, 274-296.

SU, L., AND Q. CHEN (2013): "Testing Homogeneity in Panel Data Models with Interactive Fixed Effects," *Econometric Theory*

Online Supplement to “Identifying Latent Structures in Panel Data”¹

Liangjun Su^a, Zhentao Shi^b, and Peter C. B. Phillips^c

^a*School of Economics, Singapore Management University*

^b*Department of Economics, Chinese University of Hong Kong*

^c*Yale University, University of Auckland,
University of Southampton & Singapore Management University*

This supplement is composed of four parts. Section S1 contains the proofs of some technical lemmas for the proofs of the main results in Section 2. Section S2 gives bias correction formulae in linear panel data models for both PPL and PGMM estimation. Sections S3 and S4 contain some additional simulation and applications results, respectively.

In this appendix, we state and prove some technical lemmas that are used in the proofs of the main results in Section 2. We first state an exponential inequality for strong mixing processes.

Let $\{x_t = 1, 2, \dots\}$ be a zero-mean strong mixing process, not necessarily stationary, with the mixing coefficients satisfying $\alpha(k) \leq \alpha_0 \tau^k$ for some $\alpha_0 > 0$ and $\tau \in (0, 1)$. If $\sup_{1 \leq t \leq T} |x_t| \leq T$ then there exists a constant c_0 depending on α_0 and τ such that for any $\lambda \geq 2$ and $\epsilon > 0$

$$\mathbb{P}\left(\left|\sum_{t=1}^T x_t\right| \geq \lambda\right) \leq \exp\left(-\frac{c_0 \lambda^2}{\frac{2}{\epsilon} + \frac{2}{T} + T(\ln \lambda)^2}\right)$$

where $\frac{2}{\epsilon} = \sup_{t \geq 1} [\text{Var}(x_t) + 2 \sum_{s=t+1}^{\infty} |\text{Cov}(x_t, x_s)|]$

Merlevède, Peilgrad, and Rio (2009, Theorem 2) prove (i) under the condition $\alpha(k) \leq \exp(-2k)$ for some $\epsilon > 0$. If $\alpha_0 = 1$ we can take $\tau = \exp(-2)$ and apply the theorem to obtain the claim. ■

The above lemma is used in the proof of the following lemma.

(i) Let $(x_{it}; y_{it})$ be a \mathbb{R}^{d_ϵ} -valued function indexed by the parameter $\theta \in \Phi$ where Φ is a convex compact set in \mathbb{R}^{p+1} and $\mathbb{E}[(x_{it}; y_{it})] = 0$ for all $\theta \in \Phi$. Assume that there exists a function $\psi(\theta)$ such that $\| (x_{it}; y_{it}) - \psi(\theta) \| \leq \| (x_{it}; y_{it}) \| - \psi(\theta)$ for all $\theta \in \Phi$ and $\sup_{\theta \in \Phi} \| (x_{it}; y_{it}) \| \leq \psi(\theta)$. Assume that $\mathbb{E}|(x_{it}; y_{it})|^q = O(\psi(\theta)^{q/2-1})$ for some $q \geq 6$ such that $\psi(\theta) = O(\psi(\theta)^{q/2-1})$. Let $\{x_i\}$ be a nonstochastic sequence in Φ . Then

(i) $\max_{1 \leq i \leq N} \left\| \frac{1}{\sqrt{T}} \sum_{t=1}^T (x_{it}; y_{it}) \right\| = O_P((\ln N)^3)$

(ii) $\max_{1 \leq i \leq N} \left(\left\| \frac{1}{T} \sum_{t=1}^T (x_{it}; y_{it}) \right\| \geq \epsilon \right) = O(N^{-1})$ for any given $\epsilon > 0$

(iii) $\left(\max_{1 \leq i \leq N} \left\| \frac{1}{T} \sum_{t=1}^T (x_{it}; y_{it}) \right\| \geq \epsilon \right) = O(N^{-1})$ for any given $\epsilon > 0$ if $\psi(\theta) = O(\psi(\theta)^{q/2-1})$

where $\psi(\theta) = N_T$ satisfies $(\ln N)^3 = O(\psi(\theta)^{1/2})$

¹Acknowledgements made in the leading footnote of the main paper apply also to this Online Supplement. In particular, Su acknowledges support from the Singapore Ministry of Education for Academic Research Fund (AcRF) under the Tier-2 grant number MOE2012-T2-2-021. Phillips acknowledges NSF support under Grant Nos. SES-0956687 and SES-1285258. Address Correspondence to: Liangjun Su, School of Economics, Singapore Management University, 90 Stamford Road, Singapore 178903; E-mail: ljsu@smu.edu.sg, Phone: +65 6828 0386.

(i) Let $\mathbf{1}_{it} = \mathbf{1} \{ \|\xi_{it}\| \leq 1 \}$ and $\bar{\mathbf{1}}_{it} = 1 - \mathbf{1}_{it}$. Define

$$\begin{aligned} \mathbf{1}_{it} &= \xi_{it}' \{ \xi_{it} \mathbf{1}_{it} - \mathbb{E}[\xi_{it} \mathbf{1}_{it}] \} \\ \mathbf{2}_{it} &= \xi_{it}' \bar{\mathbf{1}}_{it} \quad \text{and} \quad \mathbf{3}_{it} = -\xi_{it}' \bar{\mathbf{1}}_{it} \end{aligned}$$

Apparently $\mathbf{1}_{it} + \mathbf{2}_{it} + \mathbf{3}_{it} = \xi_{it}'$ as $\mathbb{E}[\xi_{it}] = 0$. We prove the lemma by showing that (i1) $\max_{1 \leq i \leq N} \left\| \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{1}_{it} \right\| = O_p((\ln N)^3)$ (i2) $\max_{1 \leq i \leq N} \left\| \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{2}_{it} \right\| \geq (\ln N)^3$ = (1) for any given $\epsilon > 0$ and (i3) $\max_{1 \leq i \leq N} \left\| \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{3}_{it} \right\| = o_p((\ln N)^3)$

First, we prove (i3). By the Hölder and Markov inequalities

$$\begin{aligned} \max_{1 \leq i \leq N} \left\| \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{3}_{it} \right\| &\leq \frac{1}{2} \max_{1 \leq i \leq N} \max_{1 \leq t \leq T} \|\mathbb{E}[\xi_{it} \bar{\mathbf{1}}_{it}]\| \\ &\leq \frac{1}{2} \max_{1 \leq i \leq N} \max_{1 \leq t \leq T} \mathbb{E} \|\xi_{it}\|^{q/2} \|\bar{\mathbf{1}}_{it}\|^{2/q} \left(\sum_{t=1}^T \|\xi_{it}\|^{q/2} \|\bar{\mathbf{1}}_{it}\|^{2/q} \right)^{1/2} \\ &\leq \frac{1}{2} \max_{1 \leq i \leq N} \max_{1 \leq t \leq T} \|\xi_{it}\|^{q/2} \|\bar{\mathbf{1}}_{it}\|^{2/q} \\ &\leq \frac{1}{2} \max_{1 \leq i \leq N} \max_{1 \leq t \leq T} \|\xi_{it}\|^q \|\bar{\mathbf{1}}_{it}\|^{2/q} \\ &= \frac{1}{2} \max_{1 \leq i \leq N} \max_{1 \leq t \leq T} \|\xi_{it}\|^q \|\bar{\mathbf{1}}_{it}\|^{2/q} = o_p((\ln N)^3) \text{ for any } \epsilon > 0 \end{aligned}$$

where $1_q \equiv \max_{1 \leq i \leq N} \max_{1 \leq t \leq T} \|\xi_{it}\|^{q/2}$ and $2_q \equiv \max_{1 \leq i \leq N} \max_{1 \leq t \leq T} \{\mathbb{E}(\|\xi_{it}\|^q)\}^{(q-2)/q}$

Next, we prove (i2). Noting that $\left\| \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{2}_{it} \right\| \geq (\ln N)^3$ implies that $\max_{1 \leq t \leq T} \|\xi_{it}\| \geq (\ln N)^3$ for some i . By the Boole and Markov inequalities, the dominated convergence theorem, and the stated conditions, we have

$$\begin{aligned} \left[\max_{1 \leq i \leq N} \left\| \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{2}_{it} \right\| \geq (\ln N)^3 \right] &\leq \left[\max_{1 \leq i \leq N} \max_{1 \leq t \leq T} \|\xi_{it}\| \geq (\ln N)^3 \right] \\ &\leq \max_{1 \leq i \leq N} \max_{1 \leq t \leq T} \mathbb{P}(\|\xi_{it}\| \geq (\ln N)^3) \\ &\leq \frac{1}{(\ln N)^{3q}} \max_{1 \leq i \leq N} \max_{1 \leq t \leq T} \mathbb{E} \|\xi_{it}\|^q \\ &= o_p(1) \end{aligned}$$

Now, we prove (i1). We observe that for any $\epsilon > 0$

$$\left[\max_{1 \leq i \leq N} \left\| \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{1}_{it} \right\| \geq (\ln N)^3 \right] \leq \sum_{i=1}^N \left[\left\| \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{1}_{it} \right\| \geq (\ln N)^3 \right]$$

We choose $\epsilon > 0$ and divide Φ into subsets Φ_j , $j = 1, \dots, \epsilon$ such that $\|\xi_{it}\| \leq \sqrt{\epsilon}$ for all $i \in \Phi_j$, where $\epsilon = (\ln N)^{2(p+1)/2}$. Then

$$\begin{aligned} \sum_{i=1}^N \left[\left\| \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{1}_{it} \right\| \geq (\ln N)^3 \right] &\leq \sum_{i=1}^N \left[\sup_{\phi \in \Phi} \left\| \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{1}_{it} \right\| \geq (\ln N)^3 \right] \\ &\leq \sum_{j=1}^{\epsilon} \sum_{i \in \Phi_j} \left[\sup_{\phi \in \Phi_j} \left\| \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{1}_{it} \right\| \geq (\ln N)^3 \right] \end{aligned}$$

Let $\phi_j \in \Phi_j$. Then for any $\phi \in \Phi_j$ we have

$$\begin{aligned} \left\| \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{1}(\phi; \mathbf{it}) \right\| &\leq \left\| \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{1}(\phi; \mathbf{j}) \right\| + \left\| \frac{1}{\sqrt{T}} \sum_{t=1}^T [\mathbf{1}(\phi; \mathbf{j}) - \mathbf{1}(\phi; \mathbf{it})] \right\| \\ &\leq \left\| \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{1}(\phi; \mathbf{j}) \right\| + \frac{2}{\sqrt{T}} \sum_{t=1}^T \|\mathbf{it} - \mathbf{j}\| \\ &\leq \left\| \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{1}(\phi; \mathbf{j}) \right\| + \frac{2}{\sqrt{T}} \left\{ \sum_{t=1}^T \|\mathbf{it} - \mathbb{E}[\mathbf{it}]\| \right\} + \frac{2}{\sqrt{T}} \sum_{t=1}^T \mathbb{E}[\|\mathbf{it}\|] \end{aligned}$$

It follows that

$$\left[\sup_{\phi \in \Phi_j} \left\| \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{1}(\phi; \mathbf{it}) \right\| \geq (\ln 3)^3 \right] \leq \left[\left\| \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{1}(\phi; \mathbf{j}) \right\| \geq (\ln 3)^3 \right] + \left[\frac{2}{\sqrt{T}} \sum_{t=1}^T \|\mathbf{it} - \mathbb{E}[\mathbf{it}]\| \geq (\ln 3)^3 \right]$$

as $\frac{2}{\sqrt{T}} \sum_{t=1}^T \mathbb{E}[\|\mathbf{it}\|] \geq (\ln 3)^3 = 0$. Then

$$\left[\max_{1 \leq i \leq N} \left\| \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{1}(\phi; \mathbf{i}) \right\| \geq (\ln 3)^3 \right] \leq \sum_{i=1}^N \sum_{j=1}^{n_\varepsilon} \left[\left\| \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{1}(\phi; \mathbf{j}) \right\| \geq (\ln 3)^3 \right] + \sum_{i=1}^N \sum_{j=1}^{n_\varepsilon} \left[\frac{2}{\sqrt{T}} \sum_{t=1}^T \|\mathbf{it} - \mathbb{E}[\mathbf{it}]\| \geq (\ln 3)^3 \right]$$

For the first term, we have by Lemma S1.1

$$\begin{aligned} \sum_{i=1}^N \sum_{j=1}^{n_\varepsilon} \left[\left\| \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{1}(\phi; \mathbf{j}) \right\| \geq (\ln 3)^3 \right] &\leq \varepsilon \exp\left(-\frac{2 \ln^6 3}{-2 + 4 \frac{2C}{NT} + \frac{2C}{3} \frac{1}{NT}}\right) \\ &= \exp\left(-\frac{2 \ln^6 3}{-2 + 4 \frac{2C}{NT} + \frac{2C}{3}} + \ln \varepsilon\right) \\ &\rightarrow 0 \text{ for sufficiently large } \varepsilon \end{aligned}$$

Similarly, we can show that $\sum_{i=1}^N \sum_{j=1}^{n_\varepsilon} \left[\frac{2}{\sqrt{T}} \sum_{t=1}^T \|\mathbf{it} - \mathbb{E}[\mathbf{it}]\| \geq (\ln 3)^3 \right] = 0$. Then (i1) follows. This completes the proof of (i).

(ii) Let $\mathbf{1}_2$ and $\mathbf{3it}$ be as defined in (i). Noting that $\mathbf{3it}$ is nonrandom, it suffices to show that for any given $\varepsilon > 0$ we have (ii1) $\max_{1 \leq i \leq N} \left(\left\| \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{1}(\phi; \mathbf{i}) \right\| \geq \varepsilon \right) = 0$ (ii2) $\max_{1 \leq i \leq N} \left(\left\| \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{2}(\phi; \mathbf{i}) \right\| \geq \varepsilon \right) = 0$ and (ii3) $\max_{1 \leq i \leq N} \left\| \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{3it} \right\| = 0$. Following the analysis of $\frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{3it}$ in (i), we have

$$\max_{1 \leq i \leq N} \left\| \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{3it} \right\| \leq \frac{1}{\sqrt{T}} \sum_{t=1}^T \|\mathbf{3it}\| \leq \frac{1}{\sqrt{T}} \sum_{t=1}^T 3q \|\mathbf{it}\| \leq 3q \sum_{t=1}^T \|\mathbf{it}\| \leq 3q \sum_{t=1}^T t \leq 3q \frac{T(T+1)}{2} \leq 3q \frac{T^2}{2} = 0$$

where we use the fact that $\frac{1}{\sqrt{T}} \gg -1/2(\ln \epsilon)^3$ and $\epsilon \geq 3$ by the stated conditions. Thus, (ii3) follows. Following the analysis of $\frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbb{E}[\| \hat{\mu}_i(t) - \mu_i \|^2]$ in (i2), we have

$$\begin{aligned} \max_{1 \leq i \leq N} \left(\sum_{t=1}^T \mathbb{E}[\| \hat{\mu}_i(t) - \mu_i \|^2] \geq \epsilon \right) &\leq \max_{1 \leq i \leq N} \left(\max_{1 \leq t \leq T} \mathbb{E}[\| \hat{\mu}_i(t) - \mu_i \|^2] \right) \\ &\leq \max_{1 \leq i \leq N} \max_{1 \leq t \leq T} \left(\mathbb{E}[\| \hat{\mu}_i(t) - \mu_i \|^2] \right) \\ &= \epsilon^{1-q/2} = \epsilon \end{aligned}$$

That is, (ii2) follows. For (ii1), the analysis is similar to that of $\max_{1 \leq i \leq N} \sum_{t=1}^T \mathbb{E}[\| \hat{\mu}_i(t) - \mu_i \|^2]$ in (i1) with $(\ln \epsilon)^3$ replaced by $\epsilon^{1/2}$. We now require $\epsilon^{1/2} (\ln \epsilon)^3 \rightarrow \infty$ as $\epsilon \rightarrow \infty$. This completes the proof of (ii).

(iii) Let μ_i and β_{it} be as defined in (i). Noting that β_{it} is nonrandom, it suffices to show that for any given $\epsilon > 0$ we have (iii1) $\mathbb{P}(\max_{1 \leq i \leq N} \sum_{t=1}^T \mathbb{E}[\| \hat{\mu}_i(t) - \mu_i \|^2] \geq \epsilon) = \epsilon$ (iii2) $\mathbb{P}(\max_{1 \leq i \leq N} \sum_{t=1}^T \mathbb{E}[\| \hat{\mu}_i(t) - \mu_i \|^2] \geq \epsilon) = \epsilon$ and (iii3) $\max_{1 \leq i \leq N} \sum_{t=1}^T \mathbb{E}[\| \hat{\mu}_i(t) - \mu_i \|^2] = \epsilon$. Following the analysis of $\frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbb{E}[\| \hat{\mu}_i(t) - \mu_i \|^2]$ in (i), we have

$$\max_{1 \leq i \leq N} \sum_{t=1}^T \mathbb{E}[\| \hat{\mu}_i(t) - \mu_i \|^2] \leq \epsilon^{2q} \epsilon^{(2-q)/2} = \epsilon$$

where we use the fact that $\frac{1}{\sqrt{T}} \gg -1/2(\ln \epsilon)^3$ and $\epsilon \geq 6$ by the stated conditions. Thus, (iii3) follows. Following the analysis of $\frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbb{E}[\| \hat{\mu}_i(t) - \mu_i \|^2]$ in (i2), we have

$$\begin{aligned} \mathbb{P} \left(\max_{1 \leq i \leq N} \sum_{t=1}^T \mathbb{E}[\| \hat{\mu}_i(t) - \mu_i \|^2] \geq \epsilon \right) &\leq \mathbb{P} \left(\max_{1 \leq i \leq N} \max_{1 \leq t \leq T} \mathbb{E}[\| \hat{\mu}_i(t) - \mu_i \|^2] \geq \epsilon \right) \\ &\leq \epsilon^{2q} \max_{1 \leq i \leq N} \max_{1 \leq t \leq T} \left(\mathbb{E}[\| \hat{\mu}_i(t) - \mu_i \|^2] \right) \\ &= \epsilon^{2q} \epsilon^{1-q/2} = \epsilon \end{aligned}$$

That is, (iii2) follows. For (iii1), the analysis is similar to that of $\max_{1 \leq i \leq N} \sum_{t=1}^T \mathbb{E}[\| \hat{\mu}_i(t) - \mu_i \|^2]$ in (i1) with $(\ln \epsilon)^3$ replaced by $\epsilon^{1/2}$. We now require $\epsilon^{1/2} (\ln \epsilon)^3 \rightarrow \infty$ as $\epsilon \rightarrow \infty$. This completes the proof of (iii). ■

Recall that $\hat{\Psi}_i(\beta) = \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\| \hat{\mu}_i(t) - \mu_i \|^2]$ and $\Psi_i(\beta) = \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\| \mu_i - \mu_i \|^2]$. Recall that $\hat{\mu}_i(\beta) = \arg \min_{\mu_i} \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\| \hat{\mu}_i(t) - \mu_i \|^2]$. The following three lemmas study the properties of $\hat{\Psi}_i(\beta)$ and $\hat{\mu}_i(\beta)$.

For any $\epsilon > 0$ we have $\max_{1 \leq i \leq N} \sup_{(\beta, \mu)} |\hat{\Psi}_i(\beta) - \Psi_i(\beta)| \geq \epsilon = \epsilon^{-1}$

The proof is analogous to that of Lemma S1.2(iii). ■

For any $\epsilon > 0$ we have $[\max_{1 \leq i \leq N} |\hat{\mu}_i(\beta) - \mu_i(\beta)| \geq \epsilon] = \epsilon^{-1}$

Let $\epsilon = \min_i [\inf_{\mu_i: |\mu_i - \mu_i(\beta_i)| > \eta} \Psi_i(\beta_i) - \Psi_i(\mu_i(\beta_i))]$. Then $\epsilon > 0$ by Assumptions A1(ii) and (v). Then conditional on the event $\mathcal{E} \equiv \max_{1 \leq i \leq N} \sup_{(\beta, \mu)} |\hat{\Psi}_i(\beta) - \Psi_i(\beta)| \leq \frac{1}{3}$ we have

$$\begin{aligned} \inf_{|\mu_i - \mu_i(\beta_i)| > \eta} \hat{\Psi}_i(\beta_i) &\geq \inf_{|\mu_i - \mu_i(\beta_i)| > \eta} \Psi_i(\beta_i) - \frac{1}{3} \\ &\geq \Psi_i(\mu_i(\beta_i)) + \frac{2}{3} \\ &\geq \hat{\Psi}_i(\mu_i(\beta_i)) + \frac{1}{3} \end{aligned}$$

On the other hand, $\hat{\Psi}_i(\hat{i}_i(i)) \leq \hat{\Psi}_i(i_i(i))$. It follows that $(\max_{1 \leq i \leq N} |\hat{i}_i(i) - i_i(i)| \leq \epsilon) \leq P(\epsilon) = O(\epsilon^{-1})$ by Lemma S1.3. ■

- (i) $\hat{i}_i(i) - i_i(i) = O_P(\epsilon^{-1/2})$ for each
- (ii) $\max_{1 \leq i \leq N} |\hat{i}_i(i) - i_i(i)| = O_P(\epsilon^{-1/2} (\ln N)^3)$
- (iii) $\max_{1 \leq i \leq N} |\Psi_i(\hat{i}_i(i)) - \Psi_i(i_i(i))| = O_P(\epsilon^{-1/2} (\ln N)^3)$
- (iv) $(\max_{1 \leq i \leq N} |\hat{i}_i(i) - i_i(i)| \geq \epsilon^{-1/2} (\ln N)^{3+\nu}) = O(\epsilon^{-1})$ for any $\nu > 0$ and $\epsilon > 0$
- (v) $(\max_{1 \leq i \leq N} |\Psi_i(\hat{i}_i(i)) - \Psi_i(i_i(i))| \geq \epsilon^{-1/2} (\ln N)^{3+\nu}) = O(\epsilon^{-1})$ for any $\nu > 0$ and $\epsilon > 0$

(i)-(ii) Noting that $\hat{i}_i(i) = \arg \min_{\mu_i} \frac{1}{T} \sum_{t=1}^T \ell(\mu_i; i_t(i))$ we have

$$\begin{aligned} 0 &= \frac{1}{T} \sum_{t=1}^T \ell(\mu_i; i_t(\hat{i}_i(i))) \\ &= \frac{1}{T} \sum_{t=1}^T \ell(\mu_i; i_t(i_i(i))) + \frac{1}{T} \sum_{t=1}^T \mu_i(\mu_i; i_t(\tilde{i}_i(i))) [\hat{i}_i(i) - i_i(i)] \end{aligned}$$

where $\tilde{i}_i(i)$ lies between $\hat{i}_i(i)$ and $i_i(i)$ for each i . It follows that

$$\hat{i}_i(i) - i_i(i) = - \left[\frac{1}{T} \sum_{t=1}^T \mu_i(\mu_i; i_t(\tilde{i}_i(i))) \right]^{-1} \frac{1}{T} \sum_{t=1}^T \ell(\mu_i; i_t(i_i(i))) \quad (\text{S1})$$

provided $\frac{1}{T} \sum_{t=1}^T \mu_i(\mu_i; i_t(\tilde{i}_i(i)))$ is asymptotically nonvanishing. Let $i_t(i) = i_t(\mu_i; i_i(i))$. Noting that $\mathbb{E}[i_t(i)] = 0$ and

$$\begin{aligned} \text{Var} \left(\frac{1}{T} \sum_{t=1}^T i_t(i) \right) &= \frac{1}{2} \sum_{t=1}^T \sum_{s=1}^T \text{Cov}(i_t(i), i_s(i)) \\ &\leq 8 \max_{i,t} \{\mathbb{E}|i_t(i)|^q\}^{2/q} \frac{1}{2} \sum_{t=1}^T \sum_{s=1}^T (|t-s|)^{1-2/q} \\ &\leq 8 \max_{i,t} \{\mathbb{E}|i_t(i)|^q\}^{2/q} \frac{1}{\tau-1} (T)^{1-2/q} = O\left(\frac{1}{\tau-1}\right) \end{aligned}$$

by the Davydov inequality (e.g., Corollary A.2 in Hall and Heyde (1980)), we have $\frac{1}{T} \sum_{t=1}^T i_t(i) = O_P(\epsilon^{-1/2})$ by the Chebyshev inequality. In addition, by a simple application of Lemma S1.2(i), we can show that $\max_{1 \leq i \leq N} \left| \frac{1}{T} \sum_{t=1}^T i_t(i) \right| = O_P(\epsilon^{-1/2} (\ln N)^3)$

For $\frac{1}{T} \sum_{t=1}^T \mu_i(\mu_i; i_t(\tilde{i}_i(i)))$ we make the following decomposition:

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \mu_i(\mu_i; i_t(\tilde{i}_i(i))) &= \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\mu_i(\mu_i; i_t(i_i(i)))] \\ &\quad + \frac{1}{T} \sum_{t=1}^T \{ \mu_i(\mu_i; i_t(i_i(i))) - \mathbb{E}[\mu_i(\mu_i; i_t(i_i(i)))] \} \\ &\quad + \frac{1}{T} \sum_{t=1}^T \{ \mu_i(\mu_i; i_t(\tilde{i}_i(i))) - \mu_i(\mu_i; i_t(i_i(i))) \} \end{aligned} \quad (\text{S2})$$

By Assumption A1(v), $\frac{1}{T} \sum_{t=1}^T \mathbb{E} [\mu_i (it; i_i(i))] = i\mu\mu (i) \geq H - 0$ uniformly in i . By a simple application of Lemma S1.2(i), we have

$$\max_{1 \leq i \leq N} \left| \frac{1}{T} \sum_{t=1}^T \{ \mu_i (it; i_i(i)) - \mathbb{E} [\mu_i (it; i_i(i))] \} \right| = P(1)$$

Next, by Assumption A1, and Lemmas S1.2(i) and S1.4, we have

$$\begin{aligned} & \max_{1 \leq i \leq N} \left| \frac{1}{T} \sum_{t=1}^T [\mu_i (it; i_i(i)) - i\mu\mu (i)] \right| \\ & \leq \max_{1 \leq i \leq N} \frac{1}{T} \sum_{t=1}^T | \mu_i (it; i_i(i)) - i\mu\mu (i) | \\ & \leq \left\{ \max_{1 \leq i \leq N} \frac{1}{T} \sum_{t=1}^T \mathbb{E} [| \mu_i (it; i_i(i)) |] + \max_{1 \leq i \leq N} \left| \frac{1}{T} \sum_{t=1}^T \{ \mu_i (it; i_i(i)) - \mathbb{E} [\mu_i (it; i_i(i))] \} \right| \right\} \max_{1 \leq i \leq N} | \hat{i}_i(i) - i_i(i) | \\ & \leq \frac{1}{M} + P(1) P(1) = P(1) \end{aligned} \quad (S3)$$

It follows that $\frac{1}{T} \sum_{t=1}^T \mu_i (it; i_i(i)) = i\mu\mu (i) + P(1)$ uniformly in i and $|\hat{i}_i(i) - i_i(i)| = P(-1/2)$ for each i and $\max_{1 \leq i \leq N} |\hat{i}_i(i) - i_i(i)| = P(-1/2) (\ln N)^3$

(iii) In view of the definition that $\Psi_i(i) = \mathbb{E} [| \mu_i (it; i_i(i)) |]$ we have $\max_{1 \leq i \leq N} |\Psi_i(i) - \hat{\Psi}_i(i)| = \max_{i,t} \mathbb{E} [| \mu_i (it; i_i(i)) |] = P(-1/2) (\ln N)^3$

(iv) We define the following events:

$$\begin{aligned} 1 & \equiv \left\{ \max_{1 \leq i \leq N} | \hat{i}_i(i) - i_i(i) | \leq H \left(6 \frac{1}{M} \right) \right\} \\ 2 & \equiv \left\{ \max_{1 \leq i \leq N} \left| \frac{1}{T} \sum_{t=1}^T \{ \mu_i (it; i_i(i)) - \mathbb{E} [\mu_i (it; i_i(i))] \} \right| \leq \frac{1}{M} \cdot 2 \right\} \\ 3 & \equiv \left\{ \max_{1 \leq i \leq N} \left| \frac{1}{T} \sum_{t=1}^T [\mu_i (it; i_i(i)) - i\mu\mu (i)] \right| \leq H \cdot 4 \right\} \\ 4 & \equiv \left\{ \min_{1 \leq i \leq N} \left| \frac{1}{T} \sum_{t=1}^T \mu_i (it; i_i(i)) \right| \geq H \cdot 2 \right\} \\ 5 & \equiv \left\{ \min_{1 \leq i \leq N} \left| \frac{1}{T} \sum_{t=1}^T \mu_i (it; i_i(i)) \right| \geq H \cdot 4 \right\} \end{aligned}$$

Let \bar{c}_j denote the complement of c_j for $j = 1, 2, 3, 4, 5$. Let $i = \hat{i}_i(i) - i_i(i)$. By Lemmas S1.4 and S1.2(iii), $\bar{c}_1 = \bar{c}_2$ and $\bar{c}_3 = \bar{c}_4$. Then by (S3)

$$\begin{aligned} & \bar{c}_3 \\ & \leq \left(\left\{ \max_{1 \leq i \leq N} \frac{1}{T} \sum_{t=1}^T \mathbb{E} [| \mu_i (it; i_i(i)) |] + \max_{1 \leq i \leq N} \left| \frac{1}{T} \sum_{t=1}^T \{ \mu_i (it; i_i(i)) - \mathbb{E} [\mu_i (it; i_i(i))] \} \right| \right\} \max_{1 \leq i \leq N} | i | \geq H \cdot 4 \right) \\ & \leq \left(\left\{ \max_{1 \leq i \leq N} \frac{1}{T} \sum_{t=1}^T \mathbb{E} [| \mu_i (it; i_i(i)) |] + \max_{1 \leq i \leq N} \left| \frac{1}{T} \sum_{t=1}^T \{ \mu_i (it; i_i(i)) - \mathbb{E} [\mu_i (it; i_i(i))] \} \right| \right\} \max_{1 \leq i \leq N} | i | \geq H \cdot 4 \cdot 2 \right) \\ & \quad + \bar{c}_2 \\ & \leq \left(3 \frac{1}{M} \max_{1 \leq i \leq N} | i | \geq H \cdot 2 \right) + \bar{c}_2 \\ & \leq \bar{c}_1 + \bar{c}_2 = \bar{c}_3 \end{aligned}$$

where $\hat{i}^{(0)}$ lies between $\hat{i}^{(0)}$ and $i^{(0)}$. By Assumptions A1, Lemma S1.5, and the Markov inequality, one can readily show that the first term is $P(-1/2)$ and the second is $P(-1)$. It follows that $\hat{i} - i = P(-1/2)$.

(ii) By a simple application of Lemma S1.2(i), $\max_{1 \leq i \leq N} \|i\| = P(-1/2(\ln)^3)$. Next,

$$\begin{aligned} \max_{1 \leq i \leq N} \|\hat{i} - i\| &\leq \max_{1 \leq i \leq N} \frac{1}{t=1} \sum_{t=1}^T \left(\hat{i}^{(t)} - i^{(t)} \right) \\ &\leq \left\{ \max_{1 \leq i \leq N} \frac{1}{t=1} \sum_{t=1}^T \mathbb{E} \left[\left(\hat{i}^{(t)} - i^{(t)} \right) \right] + \max_{1 \leq i \leq N} \left\{ \frac{1}{t=1} \sum_{t=1}^T \left| \left(\hat{i}^{(t)} - i^{(t)} \right) - \mathbb{E} \left[\left(\hat{i}^{(t)} - i^{(t)} \right) \right] \right| \right\} \right\} \\ &\quad \times \max_{1 \leq i \leq N} \left| \hat{i}^{(0)} - i^{(0)} \right| \\ &= \left\{ (1) + P(1) \right\} P(-1/2(\ln)^3) = P(-1/2(\ln)^3) \end{aligned}$$

(iii) By the Cauchy-Schwarz inequality, $\frac{1}{N} \sum_{i=1}^N \|\hat{i} - i\|^2 \leq \frac{2}{N} \sum_{i=1}^N \|i\|^2 + \frac{2}{N} \sum_{i=1}^N \|\hat{i} - i\|^2$. The first term in $P(-1)$ by the Markov inequality and the calculation in (i). Using the decomposition of $\hat{i} - i$ in (i), we can readily show that the second term is $P(-1)$. Then $\frac{1}{N} \sum_{i=1}^N \|\hat{i} - i\|^2 = P(-1)$.

(iv) The result follows by a simple application of Lemma S1.2(ii) and Assumption A2.

(v) The proof is similar to that of (ii) but we now apply Lemmas S1.2(iii) and S1.5(iv). ■

The next lemma establishes the uniform consistency of \hat{i} .

For any $\eta > 0$ we have $\left(\max_{1 \leq i \leq N} \|\hat{i} - i\| \right) = P(-1)$

Recall that $\frac{(K_0)}{1NT, \lambda_1}(\beta, \alpha) = \frac{1}{1NT}(\beta) + \frac{\lambda_1}{N} \sum_{i=1}^N \Pi_{k=1}^{K_0} \|i - k\|$ where $\frac{1}{1NT}(\beta) = \frac{1}{NT} \sum_{i=1}^N \left(\hat{i}^{(t)} - i^{(t)} \right) = \frac{1}{N} \sum_{i=1}^N \hat{\Psi}_i(\hat{i}^{(t)} - i^{(t)})$. Noting that $(\hat{\beta}, \hat{\alpha}) = \arg \min_{(\beta, \alpha)} \frac{(K_0)}{1NT, \lambda_1}(\beta, \alpha)$ we have $\frac{(K_0)}{1NT, \lambda_1}(\hat{\beta}, \hat{\alpha}) \leq \frac{(K_0)}{1NT, \lambda_1}(\beta^0, \hat{\alpha})$ and

$$\hat{\Psi}_i(\hat{i} - i) + \frac{1}{\Pi_{k=1}^{K_0}} \|\hat{i} - k\| \leq \hat{\Psi}_i(i^{(0)} - i^{(0)}) + \frac{1}{\Pi_{k=1}^{K_0}} \|i^{(0)} - k\| \text{ for } i = 1$$

Let $\mathcal{E}_1 \equiv \left\{ \max_{1 \leq i \leq N} \inf_{\beta_i: \|\beta_i - \beta_i^0\| > \eta} \Psi_i(i^{(0)} - i^{(0)}) - \Psi_i(i^{(0)} - i^{(0)}) \right\}$. Define three events $\mathcal{E}_1 \equiv \left\{ \max_{1 \leq i \leq N} \sup_{(\beta, \mu)} |\hat{\Psi}_i(\hat{i}) - \Psi_i(i^{(0)})| \leq \frac{1}{6} \right\}$ and $\mathcal{E}_2 \equiv \left\{ \max_{1 \leq i \leq N} \sup_{\beta_i} |\Psi_i(i^{(0)} - i^{(0)}) - \Psi_i(i^{(0)} - i^{(0)})| \leq \frac{1}{6} \right\}$ and $\mathcal{E}_3 \equiv \left\{ \frac{1}{N} \max_{\beta_i: \beta_i \in \mathcal{B}} \sum_{k=1}^{K_0} \|i - k\| \leq \frac{1}{6} \right\}$. By Lemmas S1.3, S1.5(v) and Assumption A2(i), $P(\mathcal{E}_1 \cap \mathcal{E}_2 \cap \mathcal{E}_3) \geq 1 - P(\mathcal{E}_1) - P(\mathcal{E}_2) - P(\mathcal{E}_3) = 1 - P(-1)$. Then conditional on $\mathcal{E}_1 \cap \mathcal{E}_2 \cap \mathcal{E}_3$ we have uniformly in

$$\begin{aligned} &\inf_{\beta_i: \|\beta_i - \beta_i^0\| > \eta} \hat{\Psi}_i(\hat{i} - i) + \frac{1}{\Pi_{k=1}^{K_0}} \|\hat{i} - k\| \\ &\geq \inf_{\beta_i: \|\beta_i - \beta_i^0\| > \eta} \Psi_i(i^{(0)} - i^{(0)}) - \frac{1}{6} + 0 \geq \inf_{\beta_i: \|\beta_i - \beta_i^0\| > \eta} \Psi_i(i^{(0)} - i^{(0)}) - \frac{1}{6} - \frac{1}{6} \\ &\geq \Psi_i(i^{(0)} - i^{(0)}) + \frac{1}{6} - \frac{1}{6} \\ &\geq \Psi_i(i^{(0)} - i^{(0)}) - \frac{1}{6} + \frac{1}{6} - \frac{1}{6} \\ &\geq \hat{\Psi}_i(i^{(0)} - i^{(0)}) - \frac{1}{6} - \frac{1}{6} + \frac{1}{6} - \frac{1}{6} \\ &= \hat{\Psi}_i(i^{(0)} - i^{(0)}) + \frac{1}{3} \\ &\geq \hat{\Psi}_i(i^{(0)} - i^{(0)}) + \frac{1}{\Pi_{k=1}^{K_0}} \|i^{(0)} - k\| + \frac{1}{6} \end{aligned}$$

On the other hand, $\hat{\Psi}_i(\hat{i}_i - \hat{i}_i(\hat{i}_i)) + \frac{1}{\sqrt{1 - \Pi_{k=1}^{K_0}}} \|\hat{i}_i - \hat{k}_i\| \leq \hat{\Psi}_i(\hat{i}_i - \hat{i}_i(0)) + \frac{1}{\sqrt{1 - \Pi_{k=1}^{K_0}}} \|\hat{i}_i - \hat{k}_i\|$. It follows that $(\max_{1 \leq i \leq N} \|\hat{i}_i - \hat{i}_i(0)\|) = (\dots)$ ■

To state and prove the next lemma, we follow Hahn and Newey (2004) and introduce some notation. Let F_i and \hat{F}_i denote the cumulative and empirical distribution functions of i_t respectively. Let $\beta_i(\cdot) \equiv \hat{i}_i + \sqrt{\epsilon}(\hat{i}_i - i_i)$ for $\epsilon \in [0, 1/2]$. For fixed i and let $i_i(\cdot) \equiv \arg \min_{\mu_i} (\cdot; i_i) - i_i(\cdot)$ which is the solution to the estimating equation

$$0 = i_i(\cdot; i_i, i_i(i_i, i_i(\cdot))) - i_i(\cdot) \quad (\text{S4})$$

Define $\beta_i(\cdot) = i_i(i_i, i_i(\cdot)) - i_i$. Apparently, $\beta_i(0) = i_i - i_i(-1/2) = \hat{i}_i$

$$\begin{aligned} i_i(i_i) &= i_i(i_i, i_i(0)) \\ \hat{i}_i(i_i) &= i_i(i_i, i_i(-1/2)) \\ \frac{i_i(i_i)}{i_i} &= \frac{i_i(i_i, i_i(0))}{i_i} = \beta_i(0) \quad \text{and} \\ \frac{\hat{i}_i(i_i)}{i_i} &= \frac{i_i(i_i, i_i(-1/2))}{i_i} = \beta_i(-1/2) \end{aligned}$$

We study the properties of $i_i(i_i, i_i(\cdot))$ and $\beta_i(\cdot)$ in the next two lemmas.

- (i) $(\max_{1 \leq i \leq N} \max_{0 \leq \epsilon \leq T^{-1/2}} |i_i(i_i, i_i(\cdot)) - i_i(i_i)| \geq \epsilon) = P(\dots)$ for any $\epsilon > 0$
- (ii) $\max_{1 \leq i \leq N, \max\|\beta_i - \beta_i^0\| = o(1)} |i_i(i_i) - i_i(i_i^0)| = o(1)$
- (iii) $(\max_{1 \leq i \leq N, \max\|\beta_i - \beta_i^0\| = o(1)} |i_i(i_i) - \hat{i}_i(i_i)| \geq \epsilon) = o(1)$ for any $\epsilon > 0$

(i) Let $\epsilon = \min_i [\inf_{\mu_i: |\mu_i - \mu_i(\beta_i)| > \eta} \Psi_i(i_i, i_i) - \Psi_i(i_i, i_i(i_i))]$ > 0 . Noting that

$$(\cdot; i_i, i_i) - i_i(\cdot) = (1 - \sqrt{\epsilon}) \Psi_i(i_i, i_i) + \sqrt{\epsilon} \hat{\Psi}_i(i_i, i_i)$$

we have

$$\begin{aligned} |(\cdot; i_i, i_i) - i_i(\cdot) - \Psi_i(i_i, i_i)| &\leq \sqrt{\epsilon} |\hat{\Psi}_i(i_i, i_i) - \Psi_i(i_i, i_i)| \\ &\leq |\hat{\Psi}_i(i_i, i_i) - \Psi_i(i_i, i_i)| \end{aligned}$$

By Lemma S1.3, we have $[] = o(1)$ where

$$\equiv \left\{ \max_{0 \leq \epsilon \leq T^{-1/2}} \max_{1 \leq i \leq N} |(\cdot; i_i, i_i) - i_i(\cdot) - \Psi_i(i_i, i_i)| \geq \epsilon \right\}$$

Therefore for every $\epsilon \in [0, 1/2]$ and conditional on the event $[]$ we have

$$\begin{aligned} \inf_{\mu_i: |\mu_i - \mu_i(\beta_i)| > \eta} (\cdot; i_i, i_i) - i_i(\cdot) &\geq \inf_{\mu_i: |\mu_i - \mu_i(\beta_i)| > \eta} \Psi_i(i_i, i_i) - \frac{1}{3} \\ &\geq \Psi_i(i_i, i_i(i_i)) + \frac{2}{3} \\ &\geq \Psi_i(i_i, i_i(i_i)) - i_i(\cdot) + \frac{1}{3} \end{aligned}$$

By the triangle inequality,

$$\begin{aligned}
& \left\| \frac{\beta_i(\cdot; i, i(i))}{i} - \frac{\beta_i(\cdot; i, i(i))}{i} \right\| \\
& \leq \left\| \frac{\beta_i(\cdot; i, i(i))}{i} - \frac{\beta_i(\cdot; i, i(i))}{i} \right\| + \left\| \frac{\beta_i(\cdot; i, i(i))}{i} \left[i(i) - i \right] \right\| \\
& = \left\| \frac{\beta_i(\cdot; i, i(i))}{i} - \frac{\beta_i(\cdot; i, i(i))}{i} \right\| + \sqrt{\left\| \frac{\beta_i(\cdot; i, i(i))}{i} \right\| \left(i(i) - i \right)}
\end{aligned}$$

Using Lemma S1.2(iii), we have

$$\left(\max_{1 \leq i \leq N} \max_{0 \leq \epsilon \leq T^{-1/2}} \sqrt{\left\| \frac{\beta_i(\cdot; i, i(i))}{i} \right\| \left(i(i) - i \right)} \geq 4 \right) = (-1)$$

In addition, by Lemma S1.8(i),

$$\begin{aligned}
& \left(\max_{1 \leq i \leq N} \max_{0 \leq \epsilon \leq T^{-1/2}} \left\| \frac{\beta_i(\cdot; i, i(i))}{i} - \frac{\beta_i(\cdot; i, i(i))}{i} \right\| \geq 4 \right) \\
& \leq \left(\max_{1 \leq i \leq N} \left(\cdot \right) i(i) \max_{1 \leq i \leq N} \max_{0 \leq \epsilon \leq T^{-1/2}} |i(i(i)) - i(i)| \geq 4 \right) = (-1)
\end{aligned}$$

Then (S6) follows. Analogously we can prove (S7).

(ii) Recall that

$$\frac{i(i)}{i} = - \frac{\beta_i(\cdot; i, i(i))}{i \mu_i(\cdot; i, i(i))} \quad (S8)$$

To prove (ii), it suffices to show that

$$\max_{1 \leq i \leq N, \max_{\|\beta_i - \beta_i^0\| = o(1)}} \left\| \frac{\beta_i(\cdot; i, i(i))}{i} - \frac{\beta_i(\cdot; 0, i(i^0))}{i} \right\| = (1)$$

and

$$\max_{1 \leq i \leq N, \max_{\|\beta_i - \beta_i^0\| = o(1)}} \left\| \frac{\mu_i(\cdot; i, i(i))}{i} - \frac{\mu_i(\cdot; 0, i(i^0))}{i} \right\| = (1)$$

We only show the first result as the proof of the second one is similar. By Assumption A1(iv) and Lemma S1.8(ii),

$$\begin{aligned}
& \max_{1 \leq i \leq N, \max_{\|\beta_i - \beta_i^0\| = o(1)}} \left\| \frac{\beta_i(\cdot; i, i(i))}{i} - \frac{\beta_i(\cdot; 0, i(i^0))}{i} \right\| \\
& \leq \max_{i,t} \mathbb{E} \left[\left(it \right) \max_{1 \leq i \leq N, \max_{\|\beta_i - \beta_i^0\| = o(1)}} \left\{ \left\| i - \frac{0}{i} \right\| + \left| i(i(i)) - i(i^0) \right| \right\} \right] = (1)
\end{aligned}$$

(iii) By the triangle inequality,

$$\begin{aligned}
\max_{1 \leq i \leq N} \left\| \frac{\hat{i}(i)}{i} - \frac{\hat{i}(i^0)}{i} \right\| & \leq \max_{1 \leq i \leq N} \left\| \frac{\hat{i}(i)}{i} - \frac{i(i)}{i} \right\| + \max_{1 \leq i \leq N} \left\| \frac{\hat{i}(i^0)}{i} - \frac{i(i^0)}{i} \right\| \\
& \quad + \max_{1 \leq i \leq N} \left\| \frac{i(i)}{i} - \frac{i(i^0)}{i} \right\|
\end{aligned}$$

Noting that $\left(\max_{1 \leq i \leq N} \left\| \frac{\partial \hat{\mu}_i(\beta_i)}{\partial \beta_i} - \frac{\partial \mu_i(\beta_i)}{\partial \beta_i} \right\| \geq 3 \right) = (-1)$ by (i) and the last term in the above displayed equation is (1) uniformly in the set $\max_i \left\| i - \frac{0}{i} \right\| = P(1)$ by (ii), we have $\left(\max_{1 \leq i \leq N, \max_{\|\beta_i - \beta_i^0\| = o(1)} \left\| \frac{\partial \hat{\mu}_i(\beta_i)}{\partial \beta_i} - \frac{\partial \mu_i(\beta_i^0)}{\partial \beta_i} \right\| \geq \right) = (-1)$ for any $0 \blacksquare$

Recall from (A.2) that

$$\hat{\beta}_{i\beta}(\mathbf{i}) = \frac{1}{T} \sum_{t=1}^T \left[\beta_i(\mathbf{i}; \hat{\mathbf{i}}_i(\mathbf{i})) + \frac{\mu_i(\mathbf{i}; \hat{\mathbf{i}}_i(\mathbf{i}))}{i} \frac{\hat{\mathbf{i}}_i(\mathbf{i})}{i} \right]$$

Let $\tilde{\beta}_{i\beta}(\mathbf{i}) = \frac{1}{T} \sum_{t=1}^T \left[\beta_i(\mathbf{i}; \mathbf{i}_i(\mathbf{i})) + \frac{\mu_i(\mathbf{i}; \mathbf{i}_i(\mathbf{i}))}{i} \frac{\partial \mu_i(\beta_i)}{\partial \beta'_i} \right]$ Note that $\beta_{i\beta}(\mathbf{i}) = \mathbb{E}[\tilde{\beta}_{i\beta}(\mathbf{i})]$ where $\beta_{i\beta}(\cdot)$ is defined in Section 2.3. The next lemma study the asymptotics of $\hat{\beta}_{i\beta}(\mathbf{i})$

$$(i) \quad \left(\max_{1 \leq i \leq N} \left\| \hat{\beta}_{i\beta}(\tilde{\mathbf{i}}_i) - \beta_{i\beta}(\mathbf{i}_i) \right\| \geq \epsilon \right) = o_p(1)$$

$$(ii) \quad \hat{H} \equiv \min_{1 \leq i \leq N} \left(\hat{\beta}_{i\beta}(\tilde{\mathbf{i}}_i) \right) \geq H - P(1)$$

(i) By the triangle inequality,

$$\begin{aligned} & \max_{1 \leq i \leq N} \left\| \hat{\beta}_{i\beta}(\tilde{\mathbf{i}}_i) - \beta_{i\beta}(\mathbf{i}_i) \right\| \\ & \leq \max_{1 \leq i \leq N} \left\| \hat{\beta}_{i\beta}(\tilde{\mathbf{i}}_i) - \hat{\beta}_{i\beta}(\mathbf{i}_i) \right\| + \max_{1 \leq i \leq N} \left\| \hat{\beta}_{i\beta}(\mathbf{i}_i) - \tilde{\beta}_{i\beta}(\mathbf{i}_i) \right\| + \max_{1 \leq i \leq N} \left\| \tilde{\beta}_{i\beta}(\mathbf{i}_i) - \beta_{i\beta}(\mathbf{i}_i) \right\| \end{aligned}$$

We prove (i) by showing that (i1) $\left(\max_{1 \leq i \leq N} \left\| \hat{\beta}_{i\beta}(\tilde{\mathbf{i}}_i) - \hat{\beta}_{i\beta}(\mathbf{i}_i) \right\| \geq \epsilon \right) = o_p(1)$ (i2) $\left(\max_{1 \leq i \leq N} \left\| \hat{\beta}_{i\beta}(\mathbf{i}_i) - \tilde{\beta}_{i\beta}(\mathbf{i}_i) \right\| \geq \epsilon \right) = o_p(1)$ and (i3) $\left(\max_{1 \leq i \leq N} \left\| \tilde{\beta}_{i\beta}(\mathbf{i}_i) - \beta_{i\beta}(\mathbf{i}_i) \right\| \geq \epsilon \right) = o_p(1)$ For (i1), we make the following decomposition:

$$\begin{aligned} \hat{\beta}_{i\beta}(\tilde{\mathbf{i}}_i) - \hat{\beta}_{i\beta}(\mathbf{i}_i) &= \frac{1}{T} \sum_{t=1}^T \left[\beta_i(\mathbf{i}; \tilde{\mathbf{i}}_i(\tilde{\mathbf{i}}_i)) - \beta_i(\mathbf{i}; \mathbf{i}_i(\mathbf{i}_i)) \right] \\ &+ \frac{1}{T} \sum_{t=1}^T \left[\frac{\mu_i(\mathbf{i}; \tilde{\mathbf{i}}_i(\tilde{\mathbf{i}}_i))}{i} \frac{\hat{\mathbf{i}}_i(\tilde{\mathbf{i}}_i)}{i} - \frac{\mu_i(\mathbf{i}; \mathbf{i}_i(\mathbf{i}_i))}{i} \frac{\hat{\mathbf{i}}_i(\mathbf{i}_i)}{i} \right] \\ &\equiv 11i + 12i \text{ say} \end{aligned}$$

For 11i we have

$$\max_{1 \leq i \leq N} \left\| 11i \right\| \leq \max_{1 \leq i \leq N} \frac{1}{T} \sum_{t=1}^T \left(\left\| \beta_i(\mathbf{i}; \tilde{\mathbf{i}}_i(\tilde{\mathbf{i}}_i)) - \beta_i(\mathbf{i}; \mathbf{i}_i(\mathbf{i}_i)) \right\| + \left\| \hat{\mathbf{i}}_i(\tilde{\mathbf{i}}_i) - \hat{\mathbf{i}}_i(\mathbf{i}_i) \right\| \right)$$

Using the arguments as used in the proof of Lemma S1.5(iv), we can show that

$$\left(\max_{1 \leq i \leq N} \frac{1}{T} \sum_{t=1}^T \left(\left\| \beta_i(\mathbf{i}; \tilde{\mathbf{i}}_i(\tilde{\mathbf{i}}_i)) - \beta_i(\mathbf{i}; \mathbf{i}_i(\mathbf{i}_i)) \right\| \leq 2 \frac{1/q}{M} \right) = 1 - o_p(1) \right)$$

Then by Lemmas S1.7 and S1.8(iii), we can readily show that $\left(\max_{1 \leq i \leq N} \left\| 11i \right\| \geq \epsilon \right) = o_p(1)$ For 12i we make the following decomposition:

$$\begin{aligned} 12i &= \frac{1}{T} \sum_{t=1}^T \left[\frac{\mu_i(\mathbf{i}; \tilde{\mathbf{i}}_i(\tilde{\mathbf{i}}_i))}{i} \frac{\hat{\mathbf{i}}_i(\tilde{\mathbf{i}}_i)}{i} - \frac{\mu_i(\mathbf{i}; \mathbf{i}_i(\mathbf{i}_i))}{i} \frac{\hat{\mathbf{i}}_i(\mathbf{i}_i)}{i} \right] \\ &+ \frac{1}{T} \sum_{t=1}^T \left[\frac{\mu_i(\mathbf{i}; \mathbf{i}_i(\mathbf{i}_i))}{i} \frac{\hat{\mathbf{i}}_i(\tilde{\mathbf{i}}_i)}{i} - \frac{\mu_i(\mathbf{i}; \mathbf{i}_i(\mathbf{i}_i))}{i} \frac{\hat{\mathbf{i}}_i(\mathbf{i}_i)}{i} \right] \\ &\equiv 12i,1 + 12i,2 \text{ say.} \end{aligned}$$

Following the analysis of β_{11i} and applying Lemmas S1.8(i) and (iii) and Lemmas S1.9(i) and (iii), we can readily show that $(\max_{1 \leq i \leq N} \|\beta_{12i,s}\| \geq \delta) = \mathcal{O}(N^{-1})$ for $s = 1, 2$. Then $(\max_{1 \leq i \leq N} \|\beta_{12i}\| \geq \delta) = \mathcal{O}(N^{-1})$. Consequently, we have $(\max_{1 \leq i \leq N} \|\hat{\beta}_{i\beta\beta}(\tilde{\beta}_i) - \hat{\beta}_{i\beta\beta}(\beta_i)\| \geq \delta) = \mathcal{O}(N^{-1})$.

To prove (i2), we make the following decomposition:

$$\begin{aligned} \hat{\beta}_{i\beta\beta}(\beta_i) - \hat{\beta}_{i\beta\beta}(\tilde{\beta}_i) &= \frac{1}{n} \sum_{t=1}^T \left[\beta_i(\beta_i; \beta_i, \beta_i) - \beta_i(\tilde{\beta}_i; \beta_i, \beta_i) \right] \\ &\quad + \frac{1}{n} \sum_{t=1}^T \left[\mu_i(\beta_i; \beta_i, \beta_i) \frac{\tilde{\beta}_i(\beta_i)}{\beta_i} - \mu_i(\tilde{\beta}_i; \beta_i, \beta_i) \frac{\tilde{\beta}_i(\beta_i)}{\beta_i} \right] \\ &\equiv \text{21i} + \text{22i} \end{aligned}$$

Following the analysis of $\max_{1 \leq i \leq N} \|\hat{\beta}_{i\beta\beta}(\tilde{\beta}_i) - \hat{\beta}_{i\beta\beta}(\beta_i)\|$ and using Lemmas S1.2, S1.7, S1.8, and S1.9 and Assumption A1, we can show $(\max_{1 \leq i \leq N} \|\beta_{2si}\| \geq \delta) = \mathcal{O}(N^{-1})$ for $s = 1, 2$. Then (i2) holds.

Next,

$$\begin{aligned} \hat{\beta}_{i\beta\beta}(\beta_i) - \hat{\beta}_{i\beta\beta}(\tilde{\beta}_i) &= \frac{1}{n} \sum_{t=1}^T \left[\beta_i(\beta_i; \beta_i, \beta_i) - \mathbb{E} \left[\beta_i(\beta_i; \beta_i, \beta_i) \right] \right] \\ &\quad + \frac{1}{n} \sum_{t=1}^T \left\{ \mu_i(\beta_i; \beta_i, \beta_i) - \mathbb{E} \left[\mu_i(\beta_i; \beta_i, \beta_i) \right] \right\} \end{aligned}$$

Let ϵ_0 be an arbitrary constant. By Theorem 2.2, $(\|\hat{\mathbf{G}}_{k,1}\| \geq \epsilon_0) \leq P(\hat{\mathbf{G}}_{kNT} \rightarrow 0)$ and $(\|\hat{\mathbf{G}}_{k,2}\| \geq \epsilon_0) \leq P(\hat{\mathbf{G}}_{kNT} \rightarrow 0)$. Thus $\hat{\mathbf{G}}_k = \hat{\mathbf{G}}_k^0 + P(1)$ and it suffices to prove the lemma by showing that (i) $\hat{\mathbf{G}}_k^0 + \mathbb{B}_{kNT} = \frac{1}{\sqrt{N_k T}} \sum_{i \in G_k^0} \sum_{t=1}^T \mathbf{U}_{it} + P(1)$ and (ii) $\frac{1}{\sqrt{N_k T}} \sum_{i \in G_k^0} \sum_{t=1}^T \mathbf{U}_{it} \xrightarrow{D} (0, \Omega_k)$.

We prove $\hat{\mathbf{G}}_k^0 + \mathbb{B}_{kNT} = \frac{1}{\sqrt{N_k T}} \sum_{i \in G_k^0} \sum_{t=1}^T \mathbf{U}_{it} + P(1)$. By second order Taylor expansion,

$$\begin{aligned} \hat{\mathbf{G}}_k^0 &= \frac{1}{\sqrt{k}} \sum_{i \in G_k^0} \sum_{t=1}^T \mathbf{U}_{it} + \frac{1}{\sqrt{k}} \sum_{i \in G_k^0} \sum_{t=1}^T \frac{\mu_i}{it} [\hat{\mathbf{G}}_i(\frac{0}{k}) - \frac{0}{i}] \\ &\quad + \frac{1}{2\sqrt{k}} \sum_{i \in G_k^0} \sum_{t=1}^T \frac{\mu_i \mu_i}{it} [\hat{\mathbf{G}}_i(\frac{0}{k}) - \frac{0}{i}]^2 \\ &\equiv \mathbf{k}_{k,1} + \mathbf{k}_{k,2} + \mathbf{k}_{k,3} \text{ say,} \end{aligned} \tag{S9}$$

where $\hat{\mathbf{G}}_i^*$ lies between $\hat{\mathbf{G}}_i(\frac{0}{k})$ and $\frac{0}{i}$. We will show that $\mathbf{k}_{k,1}$ contributes to the asymptotic variance of $\hat{\mathbf{G}}_k^0$, $\mathbf{k}_{k,3}$ contributes to the asymptotic bias, and $\mathbf{k}_{k,2}$ contributes to both. We analyze $\mathbf{k}_{k,3}$ first. Let $\mathbf{k}_{k,3} = \frac{1}{2\sqrt{N_k T}} \sum_{i \in G_k^0} \sum_{t=1}^T \frac{\mu_i \mu_i}{it} [\hat{\mathbf{G}}_i(\frac{0}{k}) - \frac{0}{i}]^2$. By Assumption A1, the Markov inequality, and Lemma S1.5(ii), we have

$$\begin{aligned} \|\mathbf{k}_{k,3} - \frac{0}{\mathbf{k}_{k,3}}\| &= \frac{1}{2\sqrt{k}} \sum_{i \in G_k^0} \sum_{t=1}^T \left\| \frac{\mu_i \mu_i}{it} [\hat{\mathbf{G}}_i(\frac{0}{k}) - \frac{0}{i}]^2 - \frac{\mu_i \mu_i}{it} [\hat{\mathbf{G}}_i(\frac{0}{k}) - \frac{0}{i}] \right\| \\ &\leq \left\{ \frac{1}{2k} \sum_{i \in G_k^0} \sum_{t=1}^T (\|\hat{\mathbf{G}}_i(\frac{0}{k}) - \frac{0}{i}\|) \right\} \frac{1}{k} \end{aligned}$$

Now, we study $k_{,2}$. By Lemma S1.5(ii), (S1) in its proof, and the fact that $\max_{1 \leq i \leq N} \left| \frac{1}{T} \prod_{t=1}^T \frac{\mu_i}{it} \right| = P(-1/2(\ln \frac{1}{k})^3)$ and $\max_{1 \leq i \leq N} \left| \frac{1}{T} \prod_{t=1}^T \frac{\mu_i}{it} - \frac{1}{iV} \right| = P(-1/2(\ln \frac{1}{k})^3)$ we have

$$\begin{aligned} \hat{\mu}_i \left(\frac{0}{k} \right) - \frac{0}{i} &= - \frac{\frac{1}{T} \prod_{t=1}^T \frac{\mu_i}{it}}{\frac{1}{T} \prod_{t=1}^T \frac{\mu_i}{it} \left(\frac{0}{k} \right) \sim \hat{\mu}_i \left(\frac{0}{k} \right)} = - \frac{\frac{1}{T} \prod_{t=1}^T \frac{\mu_i}{it}}{\frac{1}{T} \prod_{t=1}^T \frac{\mu_i}{it}} + P(-1(\ln \frac{1}{k})^6) \\ &= - \frac{1}{iV} \prod_{t=1}^T \frac{\mu_i}{it} + P(-1(\ln \frac{1}{k})^6) \text{ uniformly in } i \in \frac{0}{k} \end{aligned}$$

But the above expansion is not sufficient to study $k_{,2}$ and we need to get better control on the remainder term. Noting that $\hat{\mu}_i \left(\frac{0}{i} \right) = \arg \min_{\mu_i} \frac{1}{T} \prod_{t=1}^T \left(\frac{\mu_i}{it}; \frac{0}{i} \right)$ we have

$$\begin{aligned} 0 &= \frac{1}{T} \prod_{t=1}^T \frac{\mu_i}{it} \left(\frac{0}{i}; \frac{0}{i} \hat{\mu}_i \left(\frac{0}{i} \right) \right) \\ &= \frac{1}{T} \prod_{t=1}^T \frac{\mu_i}{it} + \frac{1}{T} \prod_{t=1}^T \frac{\mu_i}{it} [\hat{\mu}_i \left(\frac{0}{i} \right) - \frac{0}{i} \left(\frac{0}{i} \right)] + \frac{1}{2} \prod_{t=1}^T \frac{\mu_i \mu_i}{it} \left(\frac{0}{i}; \frac{0}{i} \right) [\hat{\mu}_i \left(\frac{0}{i} \right) - \frac{0}{i} \left(\frac{0}{i} \right)]^2 \end{aligned}$$

where $\tilde{\mu}_i \left(\frac{0}{i} \right)$ lies between $\hat{\mu}_i \left(\frac{0}{i} \right)$ and $\frac{0}{i} \left(\frac{0}{i} \right)$ for each i . It follows that

$$\begin{aligned} &\hat{\mu}_i \left(\frac{0}{i} \right) - \frac{0}{i} \left(\frac{0}{i} \right) \\ &= - \left[\frac{1}{T} \prod_{t=1}^T \frac{\mu_i}{it} \right]^{-1} \left\{ \frac{1}{T} \prod_{t=1}^T \frac{\mu_i}{it} + \frac{1}{2} \prod_{t=1}^T \frac{\mu_i \mu_i}{it} \left(\frac{0}{i}; \frac{0}{i} \right) [\hat{\mu}_i \left(\frac{0}{i} \right) - \frac{0}{i} \left(\frac{0}{i} \right)]^2 \right\} \\ &= - \left[\frac{1}{T} \prod_{t=1}^T \frac{\mu_i}{it} \right]^{-1} \left\{ \frac{1}{T} \prod_{t=1}^T \frac{\mu_i}{it} + \frac{1}{2} \prod_{t=1}^T \frac{\mu_i \mu_i}{it} [\hat{\mu}_i \left(\frac{0}{i} \right) - \frac{0}{i} \left(\frac{0}{i} \right)]^2 \right\} + P(-3(\ln \frac{1}{k})^9) \\ &= - \left[\frac{1}{T} \prod_{t=1}^T \frac{\mu_i}{it} \right]^{-1} \left\{ \frac{1}{T} \prod_{t=1}^T \frac{\mu_i}{it} + \frac{1}{2} \frac{-2}{iV} \prod_{t=1}^T \frac{\mu_i \mu_i}{it} \left(\frac{1}{T} \prod_{t=1}^T \frac{\mu_i}{it} \right)^2 \right\} + P(-3(\ln \frac{1}{k})^9) \\ &= - \left[\frac{1}{T} \prod_{t=1}^T \frac{\mu_i}{it} \right]^{-1} \left\{ \frac{1}{T} \prod_{t=1}^T \frac{\mu_i}{it} + \frac{1}{2} \frac{-2}{iV} \frac{-2}{iV2} \left(\frac{1}{T} \prod_{t=1}^T \frac{\mu_i}{it} \right)^2 \right\} + P(-3(\ln \frac{1}{k})^9) \end{aligned} \quad (\text{S11})$$

where we use the fact $\max_{1 \leq i \leq N} \left| \frac{1}{T} \prod_{t=1}^T \left[\frac{\mu_i \mu_i}{it} \left(\frac{0}{i}; \frac{0}{i} \right) - \frac{\mu_i \mu_i}{it} \right] \right| \leq \max_{1 \leq i \leq N} \frac{1}{T} \prod_{t=1}^T \left(\frac{\mu_i}{it} \right) \times \max_{1 \leq i \leq N} \left| \frac{0}{i} \left(\frac{0}{i} \right) - \frac{0}{i} \left(\frac{0}{i} \right) \right| = P(-1/2(\ln \frac{1}{k})^3)$ and $\max_{1 \leq i \leq N} \left| \frac{1}{T} \prod_{t=1}^T \frac{\mu_i \mu_i}{it} - \frac{1}{iV2} \right| = P(-1/2(\ln \frac{1}{k})^3)$ by Lemma S1.2(i). It follows that

$$\begin{aligned} k_{,2} &= \frac{-1}{\sqrt{k}} \prod_{i \in G_k^0} \prod_{t=1}^T \frac{\mu_i}{it} \left\{ \left[\frac{1}{T} \prod_{t=1}^T \frac{\mu_i}{it} \right]^{-1} \left\{ \frac{1}{T} \prod_{t=1}^T \frac{\mu_i}{it} + \frac{1}{2} \frac{-2}{iV} \frac{-2}{iV2} \left(\frac{1}{T} \prod_{t=1}^T \frac{\mu_i}{it} \right)^2 \right\} + P(-3(\ln \frac{1}{k})^9) \right\} \\ &= - \frac{1}{\sqrt{k}} \prod_{i \in G_k^0} \prod_{t=1}^T \frac{\mu_i}{it} \frac{\frac{1}{T} \prod_{t=1}^T \frac{\mu_i}{it}}{\frac{1}{T} \prod_{t=1}^T \frac{\mu_i}{it}} \\ &\quad - \frac{1}{2\sqrt{k}} \prod_{i \in G_k^0} \prod_{t=1}^T \frac{\mu_i}{it} \frac{-2}{iV} \frac{-2}{iV2} \left(\frac{1}{T} \prod_{t=1}^T \frac{\mu_i}{it} \right)^2 + P(1) \\ &\equiv -k_{,21} - k_{,22} + P(1) \end{aligned} \quad (\text{S12})$$

For $k, 21$ we make the decomposition:

$$\begin{aligned}
 k, 21 &= \frac{1}{\sqrt{k}} \prod_{i \in G_k^0} \prod_{t=1}^T \left(\frac{it - iU}{iV} + \frac{1}{\sqrt{k}} \prod_{i \in G_k^0} \prod_{t=1}^T \frac{\frac{1}{T} \prod_{t=1}^T (\frac{\mu_i}{it} - iU)}{iV} \right) \\
 &+ \frac{1}{\sqrt{k}} \prod_{i \in G_k^0} \prod_{t=1}^T \left(\frac{it - iU}{iV} \left\{ \frac{1}{\frac{1}{T} \prod_{t=1}^T \frac{\mu_i}{it}} - \frac{1}{iV} \right\} \right) \\
 &+ \frac{1}{\sqrt{k}} \prod_{i \in G_k^0} \prod_{t=1}^T \left(\frac{1}{it} \prod_{t=1}^T (\frac{\mu_i}{it} - iU) \left\{ \frac{1}{\frac{1}{T} \prod_{t=1}^T \frac{\mu_i}{it}} - \frac{1}{iV} \right\} \right) \\
 &\stackrel{\text{---}}{=} k, 21a + k, 21b + k, 21c + k, 21d \tag{S13}
 \end{aligned}$$

Apparently, $k, 21b = \frac{1}{\sqrt{N_k T^3}} \prod_{i \in G_k^0} \prod_{t=1}^T \prod_{s=1}^T \left(\frac{\mu_i}{it} - \mathbb{E} \left(\frac{\mu_i}{it} \right) \right)$. For $k, 21c$ we can use the fact that $\max_{1 \leq i \leq N} \left| \frac{1}{T} \prod_{t=1}^T \frac{\mu_i}{it} - iV \right| = P^{-1/2} (\ln)^3$ and $\max_{1 \leq i \leq N} \left| \frac{1}{T} \prod_{t=1}^T \frac{\mu_i}{it} - iV \right| = P^{-1/2} (\ln)^3$ to show that

$$k, 21c \leq \frac{1}{\sqrt{N_k T^3}} \prod_{i \in G_k^0} \prod_{t=1}^T \prod_{s=1}^T \left(\frac{\mu_i}{it} - \mathbb{E} \left(\frac{\mu_i}{it} \right) \right) \leq P^{-1/2} (\ln)^3$$

□

Combining (S12)-(S17) yields

$$k,2 = -\frac{1}{\sqrt{k}} \prod_{i \in G_k^0} \prod_{t=1}^T \frac{it}{iV} - \left\{ \frac{1}{\sqrt{k^3}} \prod_{i \in G_k^0} \prod_{t=1}^T \prod_{s=1}^T \frac{-1}{iV} \right\}$$

second derivatives of U_i with respect to \hat{G}_k . Let $U_{it}^{\mu_i} = U_i^{\mu_i}(it; \begin{smallmatrix} 0 \\ i \end{smallmatrix} \begin{smallmatrix} 0 \\ i \end{smallmatrix})$ and $U_{it}^{\mu_i \mu_i} = U_i^{\mu_i \mu_i}(it; \begin{smallmatrix} 0 \\ i \end{smallmatrix} \begin{smallmatrix} 0 \\ i \end{smallmatrix})$. Following HK, the asymptotic bias term of \hat{G}_k takes the form:

$$\mathbb{B}_{kNT}^{HK} = \frac{1}{\sqrt{k}} \sum_{i \in G_k^0} \left[\frac{-1}{iV} \frac{1}{\sqrt{k}} \sum_{t=1}^T \right] \left[\frac{1}{\sqrt{k}} \sum_{t=1}^T \left(U_{it}^{\mu_i} - \frac{iU_2}{2iV} \right) \right]$$

where $iU_2 \equiv \frac{1}{T} \sum_{t=1}^T \mathbb{E}(U_{it}^{\mu_i \mu_i})$. Note that

$$\begin{aligned} \mathbb{B}_{kNT}^{HK} &= \frac{1}{\sqrt{k}} \sum_{i \in G_k^0} \frac{-1}{iV} \sum_{s=1}^T \sum_{t=1}^T i_s U_{it}^{\mu_i} - \frac{1}{2\sqrt{k}} \sum_{i \in G_k^0} \frac{-2}{iV} iU_2 \left(\frac{1}{\sqrt{k}} \sum_{t=1}^T \right)^2 \\ &\equiv \mathbb{B}_{1kNT}^{HK} - \mathbb{B}_{2kNT}^{HK} \text{ say.} \end{aligned}$$

Let \mathbb{B}_{1kNT} and \mathbb{B}_{2kNT} be as defined in Theorem 2.4. Apparently, $\mathbb{B}_{1kNT}^{HK} = \mathbb{B}_{1kNT}$. Noting that $iU_2 = \frac{1}{T} \sum_{t=1}^T \mathbb{E}(U_{it}^{\mu_i \mu_i} - \frac{m_{iU}}{m_{iV}} U_{it}^{\mu_i \mu_i}) = iU_2 - \frac{m_{iU}}{m_{iV}} iV_2$ with $iU_2 = \frac{1}{T} \sum_{t=1}^T \mathbb{E}(U_{it}^{\mu_i \mu_i})$ we have

$$\mathbb{B}_{2kNT}^{HK} = \frac{1}{2\sqrt{k}} \sum_{i \in G_k^0} \frac{-2}{iV} \left(iU_2 - \frac{iU}{iV} iV_2 \right) \left(\frac{1}{\sqrt{k}} \sum_{t=1}^T \right)^2 = \mathbb{B}_{2kNT}$$

It follows that $\mathbb{B}_{kNT}^{HK} = \mathbb{B}_{kNT}$

Let $\hat{G}_k \equiv \frac{1}{N_k T} \sum_{i \in \hat{G}_k} \sum_{t=1}^T \left[\beta_i(it; \hat{G}_k) + \mu_i(it; \hat{G}_k) \frac{\partial \hat{\mu}_i(\hat{G}_k)}{\partial \alpha_k} \right]$ and \hat{G}_k^0 lying between \hat{G}_k and \hat{G}_k^0 elementwise. Then $\hat{G}_k = \mathbb{H}_{kNT} + P(\hat{G}_k)$ where $\hat{G}_k = \min(1, \frac{1}{N_k T})$

As in the proof of Lemma S1.12, we can readily show that $\hat{G}_k = \hat{G}_k^0 + P(1)$ where $\hat{G}_k^0 = \frac{1}{N_k T} \sum_{i \in G_k^0} \sum_{t=1}^T \left[\beta_i(it; \hat{G}_k) + \mu_i(it; \hat{G}_k) \frac{\partial \hat{\mu}_i(\hat{G}_k)}{\partial \alpha_k} \right]$ following decomposition

$$\hat{G}_k^0 \equiv \frac{1}{k} \sum_{i \in G_k^0} \sum_{t=1}^T \left[\beta_i(it; \begin{smallmatrix} 0 \\ k \end{smallmatrix} \begin{smallmatrix} 0 \\ k \end{smallmatrix}) + \mu_i(it; \begin{smallmatrix} 0 \\ k \end{smallmatrix} \begin{smallmatrix} 0 \\ k \end{smallmatrix}) \frac{\partial \hat{\mu}_i(\hat{G}_k)}{\partial \alpha_k} \right]$$

It follows that $\hat{\cdot}_{(k)} = \mathbb{H}_{kNT} + P(\cdot)$

As before, we classify $i \in \hat{\mathcal{C}}_k(\epsilon)$ if $\|\hat{\alpha}_i - \alpha_k\| = 0$ for $i = 1, \dots, k$ and $i \in \hat{\mathcal{C}}_0(\epsilon)$ otherwise. Suppose that $i \in \hat{\mathcal{C}}_k(\epsilon)$ for $i \in \{0, 1, \dots, k\}$. Then by the pointwise consistency of $\hat{\alpha}_i$, we know that the probability limit of $\hat{\mathcal{C}}_k$ must be given by one of the columns in $\alpha^0 = \begin{pmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{pmatrix}$ and it converges in probability to the true value at the rate $O_p(n^{-1/2})$. Apparently, if \mathcal{C} contains k elements with probability limit given by α_k , we can derive that $\|\hat{\mathcal{C}}_k - \alpha_k\| = O_p(\min(n^{-1/(2n_k)}, n^{-1/2}))$ for $i = 1, \dots, k$. Without loss of generality, assume that if

By Lemma S1.12 and the fact that $\mathbb{B}_k = P^{(\quad)^{1/2}}$ we can show that ${}_{1NT,2} = P^{(\frac{-2}{NT})}$. Let ${}_{1NT,3} \equiv \frac{1}{NT} \prod_{k=1}^{K_0} \prod_{i \in G_k^0} \prod_{t=1}^T i^{(it; \frac{0}{k} \hat{i}^{(0)})} \frac{\partial \hat{\mu}_i(\alpha_k^0)}{\partial \alpha_k}$. Then

$$\begin{aligned} {}_{1NT,3} &= \frac{1}{\prod_{i=1}^N \prod_{t=1}^T} \prod_{i=1}^N \prod_{t=1}^T i^{(\frac{0}{i} \hat{i}^{(0)})} + \frac{1}{\prod_{i=1}^N \prod_{t=1}^T} \prod_{i=1}^N \prod_{t=1}^T i^{(\frac{0}{i} \hat{i}^{(0)})} \left(\frac{\hat{i}^{(0)}}{i} - \frac{i^{(0)}}{i} \right) \\ &\quad + \frac{1}{\prod_{i=1}^N \prod_{t=1}^T} \prod_{i=1}^N \prod_{t=1}^T i^{(\frac{0}{i} \hat{i}^{(0)})} \frac{\mu_i(\hat{i}^{(0)} - 0)}{it} \frac{\hat{i}^{(0)}}{k} + P^{(-1)} \\ &\equiv {}_{1NT,31} + {}_{1NT,32} + {}_{1NT,33} + P^{(-1)} \end{aligned}$$

By the Chebyshev and Davydov inequalities, we can readily show that ${}_{1NT,31} = P^{(\quad)^{-1/2}}$. By (S5),

$$\begin{aligned} \frac{\hat{i}^{(0)}}{i} - \frac{i^{(0)}}{i} &= \frac{i^{(0)} i^{(-1/2)}}{i} - \frac{i^{(0)} i^{(0)}}{i} \\ &= \frac{\prod_{i=1}^N \beta_i(\cdot; \frac{0}{i} \hat{i}^{(0)})}{\prod_{i=1}^N \mu_i(\cdot; \frac{0}{i} \hat{i}^{(0)})} \frac{i}{i} - \frac{\prod_{i=1}^N \beta_i(\cdot; \frac{0}{i} \hat{i}^{(0)})}{\prod_{i=1}^N \mu_i(\cdot; \frac{0}{i} \hat{i}^{(0)})} \frac{\hat{i}^{(0)}}{i} \\ &= \frac{iV - \hat{i}V}{iV - \hat{i}V} = \frac{iV \hat{i}V - \hat{i}V iV}{iV \hat{i}V} \\ &= \frac{iV(\hat{i}V - iV) + (iV - \hat{i}V) iV}{iV \hat{i}V} \end{aligned} \tag{S20}$$

where $iV \equiv \prod_{i=1}^N \mu_i(\cdot; \frac{0}{i} \hat{i}^{(0)})$, $\hat{i}V \equiv \prod_{i=1}^N \beta_i(\cdot; \frac{0}{i} \hat{i}^{(0)})$, $\hat{i}V \hat{i}V \equiv \prod_{i=1}^N \mu_i(\cdot; \frac{0}{i} \hat{i}^{(0)}) \hat{i}^{(0)}$ and recall $iV \equiv \prod_{i=1}^N \mu_i(\cdot; \frac{0}{i} \hat{i}^{(0)})$. Then by (S11) and Lemma S1.2(i), we can show that

$${}_{1NT,32} = \frac{1}{\prod_{i=1}^N \prod_{t=1}^T} \prod_{i=1}^N \prod_{t=1}^T i^{(\frac{-2}{iV} iV(\hat{i}V - iV))} - \frac{1}{\prod_{i=1}^N \prod_{t=1}^T} \prod_{i=1}^N \prod_{t=1}^T i^{(\frac{-1}{iV}(\hat{i}V - iV))} + P^{(\quad)}$$

For ${}_{1NT,33}$ using (S11), (S20), and Lemma S1.2(i), we can readily show that

$$\begin{aligned} {}_{1NT,33} &= \frac{1}{i} \prod_{i=1}^N \prod_{t=1}^T \frac{\mu_i \left(\binom{0}{i} - \binom{0}{i} \right)}{k} + P \left(\binom{-1}{NT} \right) \\ &= \frac{1}{i} \prod_{i=1}^N \prod_{t=1}^T \frac{i \binom{0}{k}}{k} + P \left(\binom{-1}{NT} \right) \\ &= P \left(\binom{-1}{NT} \right) + P \left(\binom{-1}{NT} \right) = P \left(\binom{-1}{NT} \right) \end{aligned}$$

Then ${}_{1NT,3} = P \left(\binom{-1}{NT} \right)$ and ${}_{1NT,3} = P \left(\binom{-2}{NT} \right)$

By Taylor expansion,

$$\begin{aligned} {}_{1NT,1} - \frac{-2}{G^0} &= \frac{2}{i} \prod_{i=1}^N \prod_{t=1}^T \left(\binom{0}{i} \binom{0}{i} \right) - \frac{-2}{G^0} \\ &= \frac{2}{i} \prod_{i=1}^N \prod_{t=1}^T \left[\binom{0}{i} \binom{0}{i} \left[\binom{0}{i} - \binom{0}{i} \right] + \frac{2}{i} \prod_{i=1}^N \prod_{t=1}^T \mu_i \left(\binom{0}{i} \binom{0}{i} \right) \left[\binom{0}{i} - \binom{0}{i} \right]^2 \right] \\ &\equiv {}_{1NT,11} + {}_{1NT,12} \end{aligned}$$

Using (S11), we can readily show that ${}_{1NT,11} = P \left(\binom{-1}{NT} \right)$ and ${}_{1NT,12} = P \left(\binom{-1}{NT} \right)$. Then ${}_{1NT} = \frac{-2}{G^0} + P \left(\binom{-1}{NT} \right)$. It follows that $\hat{G}_{(K,\lambda_1)}^2 - \frac{-2}{G^0} = P \left(\binom{-1}{NT} \right)$ for each $0 \leq \leq \max$ ■

Other choices of kernels are possible. So the bias-corrected PLS C-Lasso estimator is given by

$$\hat{\beta}_k^{(c)} = \hat{\beta}_k - \frac{1}{\hat{\alpha}_k} \mathbb{H}_{kNT}^{-1} \hat{\mathbb{B}}_{1kNT}$$

Similarly, we can obtain the bias-corrected estimator for the post-Lasso estimator $\hat{\beta}_{\hat{G}_k}$

Let $\mathbf{i} \equiv (\mathbf{i}_1 \quad \mathbf{i}_T)'$ and $\mathbf{i} \equiv (\mathbf{i}_1 \quad \mathbf{i}_T)'$. Let $\|\cdot\|_a = \{\mathbb{E} \|\cdot\|^a\}^{1/a}$ for any $a \geq 1$. Let κ denote a generic positive constant that does not depend on k and T . We add the following assumption.

ASSUMPTION D1. (i) For each $k = 1, \dots, K$, $\{(\mathbf{i}_t \quad \mathbf{i}_t) : t = 1, 2, \dots, T\}$ is strong mixing with mixing coefficients $\{\alpha_i(\cdot)\}$ such that $\alpha_i(\cdot) \leq \alpha_{\alpha, i} \tau$ for some $\alpha_{\alpha, i} < \infty$ and $\tau \in (0, 1)$. $\frac{1}{N_k} \sum_{i \in G_k^0} \alpha_{\alpha, i}^{(2q-1)/q} = O(1)$

(ii) $(\mathbf{i}_t \quad \mathbf{i}_t)$ are independent across $t \in G_k^0$ where $\mathbf{i}_t = \mathbf{1}$ or $\mathbf{0}$

(iii) $\max_{i,t} \mathbb{E} \|\mathbf{i}_t\|^{4q} < \infty$ and $\max_{i,t} \mathbb{E} \|\mathbf{i}_t\|^{4q} < \infty$ for some $q \geq 1$

(iv) As $(\frac{1}{N_k} \sum_{i \in G_k^0} \alpha_{\alpha, i}^{(2q-1)/q}) \rightarrow 0$, $T \rightarrow \infty$, $\frac{2}{T} \rightarrow 0$, $\frac{2}{T} k^{-3} \rightarrow 0$ and $k^{-1/2} \frac{1}{2} \sum_{i \in G_k^0} \alpha_i(\frac{2q-1}{2q}) \rightarrow 0$ for each $k = 1, \dots, K$

Assumption D1(i) assumes the usual mixing condition. D1(ii) assumes cross sectional independence to simplify the proof which can be relaxed at the cost of lengthy arguments. D1(iii) assumes moment conditions. The last condition in D1(iv) can be easily ensured under D1(i) because for any $T \gg -\frac{2q}{(2q-1)\ln q} \ln(k^{-1/2} \frac{1}{2})$ (e.g., $T = (\ln(k^{-1/2} \frac{1}{2}))^{1+\epsilon}$ for some $\epsilon > 0$), we have

$$\begin{aligned} k^{-1/2} \frac{1}{2} \sum_{i \in G_k^0} \alpha_i(\frac{2q-1}{2q}) &\leq \left(\frac{1}{k} \sum_{i \in G_k^0} \alpha_{\alpha, i}^{(2q-1)/(2q)} \right)^{1/2} \frac{1}{2} M_T^{(2q-1)/(2q)} \\ &= O(1) \exp\left(\ln(k^{-1/2} \frac{1}{2}) + \frac{(2-q)T}{2} \ln \right) \rightarrow 0 \end{aligned}$$

The first three requirements in D1(iv) can be easily satisfied too. For example, if $k^{-a} < \frac{1}{T}$ for some $a > 3$ it suffices to set $T = k^{1/b}$ for some $b = \max\{2, 2(3-a)\}$

Suppose Assumption D1 holds. Then $\mathbb{H}_{kNT}^{-1} \hat{\mathbb{B}}_{1kNT} - \mathbb{H}_{kNT}^{-1} \mathbb{B}_{1kNT} = o_p(1)$

Noting that $\mathbb{H}_{kNT}^{-1} \hat{\mathbb{B}}_{1kNT} - \mathbb{H}_{kNT}^{-1} \mathbb{B}_{1kNT} = (\mathbb{H}_{kNT}^{-1} - \mathbb{H}_k^{-1}) \mathbb{B}_{1kNT} + (\mathbb{H}_{kNT}^{-1} - \mathbb{H}_{kNT}^{-1}) (\hat{\mathbb{B}}_{1kNT} - \mathbb{B}_{1kNT}) + \mathbb{H}_{kNT}^{-1} (\hat{\mathbb{B}}_{1kNT} - \mathbb{B}_{1kNT})$, $\mathbb{H}_{kNT}^{-1} = O(1)$ and $\mathbb{B}_{1kNT} = o_p(k^{-1})$ it suffices to show that (i) $\mathbb{H}_{kNT} - \mathbb{H}_k = o_p(N_T)$ and (ii) $\hat{\mathbb{B}}_{1kNT} - \mathbb{B}_{1kNT} = o_p(1)$ where $N_T = \min(1, \frac{1}{k})$

We first prove (i). Let $\bar{\mathbb{H}}_{kNT} \equiv \frac{1}{N_k T} \sum_{i \in G_k^0} \sum_{t=1}^T \tilde{\mathbf{i}}_t \tilde{\mathbf{i}}_t'$. It suffices to show that (i1) $\mathbb{H}_{kNT} - \bar{\mathbb{H}}_{kNT} = o_p(N_T)$ and (i2) $\bar{\mathbb{H}}_{kNT} - \mathbb{H}_k = o_p(N_T)$. Note that

$$\begin{aligned} \mathbb{H}_{kNT} - \bar{\mathbb{H}}_{kNT} &= \frac{1}{k} \sum_{i \in \hat{G}_k} \sum_{t=1}^T \mathbf{i}_t \mathbf{i}_t' - \frac{1}{k} \sum_{i \in G_k^0} \sum_{t=1}^T \tilde{\mathbf{i}}_t \tilde{\mathbf{i}}_t' \\ &= \frac{1}{k} \left(\sum_{i \in \hat{G}_k} \mathbf{i}_t \mathbf{i}_t' - \sum_{i \in G_k^0} \mathbf{i}_t \mathbf{i}_t' \right) \sum_{t=1}^T \tilde{\mathbf{i}}_t \tilde{\mathbf{i}}_t' + \frac{k - \hat{k}}{k} \sum_{i \in G_k^0} \sum_{t=1}^T \tilde{\mathbf{i}}_t \tilde{\mathbf{i}}_t' \\ &\equiv k_{k,1} + k_{k,2} \end{aligned}$$

By Corollary 2.3, we can readily show that $k_{k,2} = o_p(k^{-1}) = o_p(N_T)$. For any $\epsilon > 0$ we have by the proof of Theorem 2.2, $(\|\mathbf{i}_{k,1}\| \geq N_T) \leq (|\hat{k}_{kNT} - k_{kNT}| \geq \epsilon) + (|\hat{k}_{kNT} - k_{kNT}| \geq \epsilon)$. It follows that $\mathbb{H}_{kNT} - \bar{\mathbb{H}}_{kNT} = o_p(N_T)$

$P(NT)$ Now,

$$\begin{aligned}
 \bar{\mathbb{H}}_{kNT} - \mathbb{H}_{kNT} &= \frac{1}{k} \sum_{i \in G_k^0} \sum_{t=1}^T \{ \tilde{it} \tilde{it}' - \mathbb{E}\{ [it - \mathbb{E}(\cdot)] [it - \mathbb{E}(\cdot)]' \} \} \\
 &= \frac{1}{k} \sum_{i \in G_k^0} \sum_{t=1}^T \{ [it - \mathbb{E}(\cdot)] [it - \mathbb{E}(\cdot)] - \mathbb{E}\{ [it - \mathbb{E}(\cdot)] [it - \mathbb{E}(\cdot)]' \} \} \\
 &\quad + \frac{1}{k} \sum_{i \in G_k^0} \sum_{t=1}^T \{ \tilde{it} \tilde{it}' - [it - \mathbb{E}(\cdot)] [it - \mathbb{E}(\cdot)]' \}
 \end{aligned}$$

For the first term, we can apply Lemma A.2(ii) in Gao (2007) and show that it is $O_p(k^{-1})$. For the second term, we can apply the Davydov inequality directly to show that it is bounded from above by

$$\frac{2}{k} \sum_{i \in G_k^0} \left(8 \sum_{s=1}^q \left\| \frac{1}{\sqrt{k}} \sum_{t=1}^T \varepsilon_{it} \right\|_{4q} \right)^2 = O_p(k^{-1})$$

It follows that $\mathbb{B}_{1kNT} - \widehat{\mathbb{B}}_{1kNT} = O_p(k^{-1/2}) = o_p(1)$.

We now show (ii2), we first make the following decomposition:

$$\begin{aligned} \mathbb{B}_{1kNT} - \widehat{\mathbb{B}}_{1kNT} &= \frac{1}{k} \sum_{i \in \widehat{G}_k} \sum_{s=1}^q \sum_{t=1}^T M_T(\varepsilon_{it}) \varepsilon_{it} - \frac{1}{k} \sum_{i \in G_k^0} \sum_{s=1}^q \sum_{t=1}^T \mathbb{E}(\varepsilon_{it} \varepsilon_{is}) \\ &= \frac{1}{k} \sum_{i \in G_k^0} \sum_{s=1}^q \sum_{t=1}^T M_T(\varepsilon_{it}) \varepsilon_{it} - \frac{1}{k} \sum_{i \in G_k^0} \sum_{s=1}^q \sum_{t=1}^T \mathbb{E}(\varepsilon_{it} \varepsilon_{is}) + o_p(1) \\ &= \frac{1}{k} \sum_{i \in G_k^0} \sum_{s=1}^q \sum_{t=1}^T M_T(\varepsilon_{it}) (\varepsilon_{it} - \mathbb{E}(\varepsilon_{it})) \\ &\quad + \frac{1}{k} \sum_{i \in G_k^0} \sum_{s=1}^q \sum_{t=1}^T M_T(\varepsilon_{it}) [\varepsilon_{is} \varepsilon_{it} - \mathbb{E}(\varepsilon_{is} \varepsilon_{it})] \\ &\quad + \frac{1}{k} \sum_{i \in G_k^0} \sum_{s=1}^q \sum_{t=1}^T M_T(\varepsilon_{it}) \mathbb{E}(\varepsilon_{is} \varepsilon_{it}) \\ &\quad + \frac{1}{k} \sum_{i \in G_k^0} \sum_{s=1}^q \sum_{t=1}^T [1 - M_T(\varepsilon_{it})] \mathbb{E}(\varepsilon_{is} \varepsilon_{it}) + o_p(1) \\ &\equiv \widehat{k}_{NT,1} + \widehat{k}_{NT,2} + \widehat{k}_{NT,3} + \widehat{k}_{NT,4} + o_p(1) \end{aligned}$$

where the $o_p(1)$ term arises due to the replacement of \widehat{G}_k by G_k^0 and this can be easily justified by using the uniform classification consistency result and arguments as used in the proof of Theorem 2.5. We prove (ii) by demonstrating that $\widehat{k}_{NT,s} = o_p(1)$ for $s = 1, 2, 3$ and 4.

We first study $\widehat{k}_{NT,1}$. Noting that $\varepsilon_{it} = \varepsilon_{it} - \frac{1}{T} \sum_{t=1}^T \varepsilon_{it} + \frac{1}{T} \sum_{t=1}^T \varepsilon_{it}$ and $\varepsilon_{it} = \frac{1}{T} \sum_{t=1}^T \varepsilon_{it} + \varepsilon_{it} - \frac{1}{T} \sum_{t=1}^T \varepsilon_{it}$ for $i \in G_k^0$ we have that for $i \in G_k^0$

$$\varepsilon_{it} - \mathbb{E}(\varepsilon_{it}) = \varepsilon_{it} - \frac{1}{T} \sum_{t=1}^T \varepsilon_{it} + \frac{1}{T} \sum_{t=1}^T \varepsilon_{it} - \mathbb{E}(\varepsilon_{it}) = \frac{1}{T} \sum_{t=1}^T (\varepsilon_{it} - \mathbb{E}(\varepsilon_{it})) - \mathbb{E}(\varepsilon_{it})$$

where $\mathbb{E}(\varepsilon_{it}) = \frac{1}{T} \sum_{t=1}^T \mathbb{E}(\varepsilon_{it})$. Then

$$\begin{aligned} \widehat{k}_{NT,1} &= \frac{1}{k} \sum_{i \in G_k^0} \sum_{s=1}^q \sum_{t=1}^T M_T(\varepsilon_{it}) \frac{1}{T} \sum_{t=1}^T (\varepsilon_{it} - \mathbb{E}(\varepsilon_{it})) - \frac{1}{k} \sum_{i \in G_k^0} \sum_{s=1}^q \sum_{t=1}^T M_T(\varepsilon_{it}) \mathbb{E}(\varepsilon_{it}) \\ &\equiv \widehat{k}_{NT,1}(1) - \widehat{k}_{NT,1}(2) \quad \text{say.} \end{aligned}$$

In view of the fact that $\hat{G}_k - \frac{0}{k} = P((k)^{-1/2} + (-1))$ and $\hat{k} = k(1 + P(1))$ we have

$$\begin{aligned}
\| \bar{k}_{NT,1}(1) \| &= \frac{1}{k^{1/2} 3/2} \left\| \sum_{i \in G_k^0} \sum_{s=1}^T \sum_{t=1}^T M_T(\cdot) \text{is} \sim \text{it} \left(\frac{0}{k} - \hat{G}_k \right) \right\| \\
&\leq \frac{k^{1/2}}{k^{1/2} 3/2} \left\| \frac{0}{k} - \hat{G}_k \right\| \frac{1}{k} \sum_{i \in G_k^0} \sum_{|s-t| \leq M_T} \left\| \text{is} \sim \text{it} \right\| \\
&= \frac{1/2}{k} \frac{1/2}{1/2} P((k)^{-1/2} + (-1)) P(T) \\
&= P\left(1 + \frac{1/2}{k} (-1/2)\right) P(T) = P(1)
\end{aligned}$$

where we use the fact that $\frac{1}{N_k T^2} \sum_{i \in G_k^0} \sum_{|s-t| \leq M_T} \left\| \text{is} \sim \text{it} \right\| = P(T)$ by moment calculation and the Markov inequality. Let $\bar{k}_{NT,1}(2) \equiv \frac{1}{N_k^{1/2} T^{3/2}} \sum_{i \in G_k^0} \sum_{s=1}^T \sum_{t=1}^T M_T(\cdot) \text{'is} \sim \text{ir}$, where $\text{'is} \sim \text{ir}$ is any $\times 1$ nonrandom vector such that $\| \text{'is} \sim \text{ir} \| = 1$. Then by Assumptions D1(i), (iii) and (iv),

$$\begin{aligned}
|\mathbb{E}[\bar{k}_{NT,1}(2)]| &\leq \frac{1}{k^{1/2} 5/2} \sum_{i \in G_k^0} \sum_{s=1}^T \sum_{t=1}^T \sum_{r=1}^T M_T(\cdot) |\mathbb{E}(\text{'is} \sim \text{ir})| \\
&\leq \frac{8}{k^{1/2} 5/2} \sum_{i \in G_k^0} \sum_{s=1}^T \sum_{t=1}^T \sum_{r=1}^T M_T(\cdot) \|\text{'is}\|_{4q} \|\text{ir}\|_{4q} i(|-|)^{(2q-1)/(2q)} \\
&\leq \frac{1/2}{k^{3/2}} \left\{ \frac{1}{k} \sum_{i \in G_k^0} \binom{2q-1}{\alpha, i} \right\} \left\{ \frac{1}{k} \sum_{t,s,r: |s-t| \leq M_T} \binom{2q-1}{|r-s|} \right\} \\
&= \frac{1/2}{k} \frac{-3/2}{(1)} (T) = \left(\frac{1/2}{T} \frac{-3/2}{k} \right) = (1)
\end{aligned}$$

Similarly, by Assumptions D1(i)-(iv),

$$\begin{aligned}
\text{Var}(\bar{k}_{NT,1}(2)) &= \frac{1}{k} \frac{1}{5} \sum_{i \in G_k^0} \text{Var} \left(\sum_{s=1}^T \sum_{t=1}^T \sum_{r=1}^T M_T(\cdot) \text{'is} \sim \text{ir} \right) \\
&\leq \frac{1}{k} \frac{1}{5} \sum_{i \in G_k^0} \mathbb{E} \left[\left(\sum_{s=1}^T \sum_{t=1}^T \sum_{r=1}^T M_T(\cdot) \text{'is} \sim \text{ir} \right)^2 \right] \\
&= \frac{1}{k} \frac{1}{5} \sum_{i \in G_k^0} \sum_{1 \leq t_1, t_2, \dots, t_6 \leq T} M_T(1, 2) M_T(4, 5) \mathbb{E}(\text{'it}_2 \text{it}_3 \text{'it}_5 \text{it}_6) \\
&\leq \frac{1}{k} \frac{1}{5} \sum_{i \in G_k^0} \sum_{\substack{1 \leq t_1, t_2, \dots, t_6 \leq T \\ |t_1 - t_2| \leq M_T, |t_4 - t_5| \leq M_T}} |\mathbb{E}(\text{'it}_2 \text{it}_3 \text{'it}_5 \text{it}_6)| \\
&= \left(\frac{2}{T} \right) = (1)
\end{aligned}$$

Consequently, $\bar{k}_{NT,1}(2) = P(1)$. This, in conjunction with Corollary 2.3, implies that $\bar{k}_{NT,1}(2) = P(1)$ as $\text{'is} \sim \text{ir}$ is arbitrary. Thus we have shown that $\hat{k}_{NT,1} = P(1)$

For $\hat{k}_{NT,2}$ note that $\hat{k}_{NT,2} = \bar{k}_{NT,2} \frac{1/2}{k} \hat{1/2} = \bar{k}_{NT,2} (1 + P(1))$ where $\bar{k}_{NT,2} = \frac{1}{N_k^{1/2} T^{3/2}} \sum_{i \in G_k^0} \sum_{s=1}^T \sum_{t=1}^T M_T(\cdot) [\text{is} \sim \text{it} - \mathbb{E}(\text{is} \sim \text{it})]$. By construction $\mathbb{E}(\bar{k}_{NT,2}) = 0$. By Assumptions D1(ii)-(iii) and

Jensen inequality,

$$\begin{aligned}
\text{Var}(\hat{\mu}_{kNT,2}^{(c)}) &= \frac{1}{k} \sum_{i \in G_k^0} \text{Var} \left[\sum_{s=1}^T \sum_{t=1}^T M_T(\hat{\mu}_{kNT,2}^{(c)})' [is \ it - \mathbb{E}(is \ \Delta \ it)] \right] \\
&\leq \frac{1}{k} \sum_{i \in G_k^0} \sum_{s=1}^T \sum_{t=1}^T \sum_{r=1}^T \sum_{l=1}^T M_T(\hat{\mu}_{kNT,2}^{(c)})' M_T(\hat{\mu}_{kNT,2}^{(c)}) \mathbb{E}(\hat{\mu}_{kNT,2}^{(c)'} is \ it \ ir \ il) \\
&\leq \frac{1}{k} \sum_{i \in G_k^0} \sum_{|s-t| \leq M_T} \sum_{|r-l| \leq M_T} |\mathbb{E}(\hat{\mu}_{kNT,2}^{(c)'} is \ it \ ir \ il)| = O\left(\frac{2}{T}\right) = O(1)
\end{aligned}$$

where the last equality follows from the fact that $\|\mathbb{E}(\hat{\mu}_{kNT,2}^{(c)'} is \ it \ ir \ il)\| \leq \max_{i,s,t} \|is \ it\|_2^2 \leq \max_{i,t} \|it\|_4^2 \times \max_{i,t} \|it\|_4^2 \propto$ by Assumption D1(iii). Then $\hat{\mu}_{kNT,2}^{(c)} = O_P(1)$ by the Chebyshev inequality and thus $\hat{\mu}_{kNT,2}^{(c)} = O_P(1)$

By Corollary 2.3 and the Davydov inequality,

$$\begin{aligned}
\|\hat{\mu}_{kNT,3}^{(c)}\| &= \frac{\|\hat{\mu}_{kNT,3}^{(c)} - \hat{\mu}_{kNT,2}^{(c)}\|}{3/2 \left(\frac{1}{k}^{-1/2} + \frac{1}{k}^{-1/2} \right)} \sum_{i \in G_k^0} \sum_{s=1}^T \sum_{t=1}^T M_T(\hat{\mu}_{kNT,2}^{(c)}) \mathbb{E}(is \ it) \\
&\leq \frac{\|\hat{\mu}_{kNT,3}^{(c)} - \hat{\mu}_{kNT,2}^{(c)}\|}{1/2 \left(\frac{1}{k}^{-1/2} + \frac{1}{k}^{-1/2} \right)} \left\{ \frac{1}{k} \sum_{i \in G_k^0} \sum_{|s-t| \leq M_T} \|\mathbb{E}(is \ it)\| \right\} \\
&= O_P\left(\frac{1}{k}^{-1/2} \frac{1}{k}^{-1/2}\right) = O_P(1)
\end{aligned}$$

By Assumptions D1(i)-(iv) and the Davydov inequality,

$$\begin{aligned}
\|\hat{\mu}_{kNT,4}^{(c)}\| &= \frac{1}{k} \sum_{i \in G_k^0} \sum_{s=1}^T \sum_{t=1}^T [1 - M_T(\hat{\mu}_{kNT,2}^{(c)})] \mathbb{E}(is \ it) \\
&= \sum_{i \in G_k^0} \sum_{s=1}^T \sum_{t=1}^T \frac{1}{k} [1 - M_T(\hat{\mu}_{kNT,2}^{(c)})] \mathbb{E}(is \ it) \\
&\leq \frac{8}{k} \sum_{i \in G_k^0} \sum_{|s-t| > M_T} i (|s-t|)^{(2q-1)/(2q)} \|is\|_{4q} \|it\|_{4q} \\
&\leq \frac{1}{k} \sum_{i \in G_k^0} i (M_T)^{(2q-1)/(2q)} = O(1)
\end{aligned}$$

This completes the proof of the proposition. \blacksquare

With the above result in hand, we can readily show that

$$\begin{aligned}
\frac{1}{k} (\hat{\mu}_{kNT}^{(c)} - \hat{\mu}_{kNT}^{(c)}) &= \frac{1}{k} (\hat{\mu}_{kNT}^{(c)} - \hat{\mu}_{kNT}^{(c)}) - \mathbb{H}_{kNT}^{-1} \mathbb{B}_{1kNT} + \left(\frac{1}{k} \hat{\mu}_{kNT}^{(c)} \right)^{1/2} \mathbb{H}_{kNT}^{-1} \mathbb{B}_{1kNT} - \mathbb{H}_{kNT}^{-1} \hat{\mathbb{B}}_{1kNT} \\
&\quad + \left[1 - \left(\frac{1}{k} \hat{\mu}_{kNT}^{(c)} \right)^{1/2} \right] \mathbb{H}_{kNT}^{-1} \mathbb{B}_{1kNT} \\
&= \frac{1}{k} (\hat{\mu}_{kNT}^{(c)} - \hat{\mu}_{kNT}^{(c)}) - \mathbb{H}_{kNT}^{-1} \mathbb{B}_{1kNT} + O_P(1) + O_P\left(\frac{1}{k}\right) \left(\frac{1}{k} \right)^{1/2} \\
&= \frac{1}{k} (\hat{\mu}_{kNT}^{(c)} - \hat{\mu}_{kNT}^{(c)}) - \mathbb{H}_{kNT}^{-1} \mathbb{B}_{1kNT} + O_P(1)
\end{aligned}$$

That is, $\sqrt{k} (\hat{\mu}_{kNT}^{(c)} - \hat{\mu}_{kNT}^{(c)})$ has the desired limiting distribution centered on the origin.

Bias correction for the PGMM C-Lasso estimator in dynamic panel data models can be done analogously. For simplicity we focus on the case where $\beta_{iNT} = \beta$ for all i . Recall from Theorem 3.4 and the remark regarding Assumption B3(iii) (see (3.3) in particular) that

$$\frac{1}{k} (\tilde{\beta}_k - \beta) - \frac{1}{k} \beta_{kNT} \xrightarrow{D} (0, \beta_k^{-1}) \text{ for } k \rightarrow \infty$$

where $\tilde{\beta}_k \equiv \frac{1}{N_k} \sum_{i \in G_k^0} \tilde{\beta}'_{i,z\Delta x} \tilde{\beta}_{i,z\Delta x}$ and $\beta_{kNT} = \frac{1}{N_k^{1/2} T^{3/2}} \sum_{i \in G_k^0} \sum_{s=1}^T \sum_{t=1}^T \mathbb{E}(\Delta_{is} \tilde{\beta}'_{is} \tilde{\beta}_{it} \Delta_{it})$. Based on (3.3), in order to verify Assumption B3(iii) we also need to show

$$\beta_{kNT} = \frac{1}{k^{1/2} T^{1/2}} \sum_{i \in G_k^0} \sum_{t=1}^T \tilde{\beta}'_{i,z\Delta x} \tilde{\beta}_{it} \Delta_{it} \xrightarrow{D} (0, \beta_k) \text{ and} \quad (\text{S1})$$

$$\beta_{kNT} = \frac{1}{k^{1/2} T^{3/2}} \sum_{i \in G_k^0} \sum_{s=1}^T \sum_{t=1}^T \{[\Delta_{is} \tilde{\beta}'_{is} - \mathbb{E}(\Delta_{is} \tilde{\beta}'_{is})] \tilde{\beta}_{it} \Delta_{it} - \mathbb{E}(\Delta_{is} \tilde{\beta}'_{is} \tilde{\beta}_{it} \Delta_{it})\} = o_p(1) \quad (\text{S2})$$

The first part is assured by a version of the CLT. Below we first propose an estimate of the bias $\beta_k^{-1} \beta_{kNT}$ and then demonstrate (S2).

To correct the bias, we propose to obtain consistent estimates of $\tilde{\beta}_k$ and β_{kNT} respectively by

$$\tilde{\beta}_k = \frac{1}{k} \sum_{i \in \tilde{G}_k} \tilde{\beta}'_{i,z\Delta x} \tilde{\beta}_{i,z\Delta x} \text{ and } \tilde{\beta}_{kNT} = \frac{1}{k^{1/2} T^{3/2}} \sum_{i \in \tilde{G}_k} \sum_{s=1}^T \sum_{t=1}^T M_T(\cdot) \Delta_{is} \tilde{\beta}'_{is} \tilde{\beta}_{it} \Delta_{it}$$

where $\tilde{\Delta}_{it} = \Delta_{it} - \tilde{\beta}'_{i,z\Delta x} \tilde{\beta}_{i,z\Delta x}$ for all $i \in \tilde{G}_k$. $M_T(\cdot)$ is as defined above: $M_T(\cdot) = \frac{1}{M_T} (|\cdot|)$ and $\frac{1}{M_T}(\cdot)$ denotes the Bartlett kernel: $\frac{1}{M_T}(\cdot) = (1 - |\cdot|/T) \mathbf{1}\{|\cdot| \leq T\}$. Note that we also allow dynamic misspecification here. If one is sure that the model is dynamically correctly specified in the sense that $\mathbb{E}(\Delta_{it} | \mathcal{F}_{i,t-1}) = 0$ where $\mathcal{F}_{i,t-1} = (\Delta_{i,t-1}, \Delta_{i,t-1}, \tilde{\beta}_{i,t-1}; \Delta_{i,t-2}, \Delta_{i,t-2}, \tilde{\beta}_{i,t-2}; \dots)$ one can use the one-sided kernel: $M_T(\cdot) = \frac{1}{M_T}(\cdot)$ where $\frac{1}{M_T}(\cdot) = (1 - \frac{\cdot}{T}) \mathbf{1}\{0 \leq \cdot \leq T\}$. The bias-corrected C-Lasso estimator of β_k would be

$$\tilde{\beta}_k^{(c)} = \tilde{\beta}_k - \frac{1}{k} \beta_k^{-1} \tilde{\beta}_{kNT}$$

Note that Theorem 3.4 indicates that there is no need to consider bias correction for the post Lasso estimator $\tilde{\beta}_{\tilde{G}_k}$.

Let $\tilde{\beta}_i \equiv (\tilde{\beta}_{i1}, \dots, \tilde{\beta}_{iT})'$ and $\tilde{\beta}_i \equiv (\tilde{\beta}_{i1}, \dots, \tilde{\beta}_{iT})'$. We add the following assumption.

ASSUMPTION D2. (i) For each $k = 1, \dots, K$ $\{(\Delta_{it} \tilde{\beta}'_{it} \tilde{\beta}_{it} \Delta_{it}) : i = 1, 2, \dots\}$ is strong mixing with mixing coefficients $\{\alpha_i(\cdot)\}$. In addition, $\alpha_i(\cdot) \leq \alpha_{\tilde{G}_k} \tau$ for some $\alpha_{\tilde{G}_k} < \infty$ and $\tau \in (0, 1)$ where $\frac{1}{N_k} \sum_{i \in G_k^0} \alpha_{\tilde{G}_k}^{(2q-1)/(2q)}$

$$= (1) \text{ and } \frac{1}{N_k} \sum_{i \in G_k^0} \alpha_{\tilde{G}_k}^{(q-1)/q} = (1)$$

(ii) $(\tilde{\beta}_i, \tilde{\beta}_i)$ are independent across $i \in \tilde{G}_k$ where $k = 1, \dots, K$

(iii) $\max_{i,t} \mathbb{E} \|\Delta_{it} \tilde{\beta}'_{it}\|^{4q} < \infty$ and $\max_{i,t} \mathbb{E} \|\tilde{\beta}_{it} \Delta_{it}\|^{4q} < \infty$ for some $q > 1$

(iv) As $(N_k, T) \rightarrow \infty$ $\frac{2}{T} \rightarrow 0$ and $\frac{1}{k^{1/2} T^{1/2}} \sum_{i \in G_k^0} \alpha_i(\frac{2}{T})^{(2q-1)/(2q)} \rightarrow 0$ for each $k = 1, \dots, K$

Assumptions D2(i)-(iv) parallel D1(i)-(iv). The major difference is that we do not need $\frac{2}{T} \rightarrow 0$ in D2(iv) but require $\frac{2}{T} \rightarrow 0$ in D2(iii).

³Observe that $\alpha_k - \alpha_k^0 = o_p(N_k T^{-1/2} T^{-1})$ and $\alpha_{\tilde{G}_k} - \alpha_k^0 = o_p(N_k T^{-1/2})$. We recommend using the post-Lasso estimator $\tilde{\beta}_{\tilde{G}_k}$.

Suppose that the conditions of Theorem 3.4 hold. Suppose Assumption D2 holds. Then

$$\tilde{k}_{kNT}^{-1} \tilde{\Sigma}_{kNT}^{-1} = k_{kNT}^{-1} \Sigma_{kNT}^{-1} + o_p(1)$$

1 + $P(1)$ by Corollary 3.3, we have

$$\begin{aligned} \|k_{NT,1}\| &= \frac{1}{k^{1/2} 3/2} \left\| \sum_{i \in G_k^0} \sum_{s=1}^T \sum_{t=1}^T M_T(\cdot) \Delta_{is} \Delta_{is}' \Delta_{it} (\Delta_{it})' \left(\frac{0}{k} - \tilde{G}_k \right) \right\| \\ &\leq \left(\frac{1}{k} \right)^{1/2} \left\| \frac{0}{k} - \tilde{G}_k \right\| \frac{1}{k} \frac{1}{2} \sum_{i \in G_k^0} \sum_{|s-t| \leq M_T} \|\Delta_{is} \Delta_{is}' \Delta_{it} (\Delta_{it})'\| \\ &= P(1) k_{NT,1} \end{aligned}$$

where $k_{NT,1} = \frac{1}{N_k T^2} \sum_{i \in G_k^0} \sum_{|s-t| \leq M_T} \|\Delta_{is} \Delta_{is}' \Delta_{it} (\Delta_{it})'\|$. By the Markov inequality, $k_{NT,1} = P(T)$. It follows that $\|k_{NT,1}\| = P(T) = P(1)$ under Assumption D2(iv).

For $k_{NT,2}$ note that $k_{NT,2} = k_{NT,2} \frac{1}{k} \frac{1}{2} \sum_{i \in G_k^0} \sum_{|s-t| \leq M_T} \|\Delta_{is} \Delta_{is}' \Delta_{it} (\Delta_{it})'\| = k_{NT,2} (1 + P(1))$ where

$$k_{NT,2} = \frac{1}{k^{1/2} 3/2} \sum_{i \in G_k^0} \sum_{s=1}^T \sum_{t=1}^T M_T(\cdot) [\Delta_{is} \Delta_{is}' \Delta_{it} \Delta_{it} - \mathbb{E}(\Delta_{is} \Delta_{is}' \Delta_{it} \Delta_{it})]$$

Let \mathbf{v} be any $\times 1$ nonrandom vector such that $\|\mathbf{v}\| = 1$. Then $\mathbb{E}(\mathbf{v}' k_{NT,2}) = 0$. By Assumptions D2(ii)-(iv) and Jensen inequality,

$$\begin{aligned} &\text{Var}(\mathbf{v}' k_{NT,2}) \\ &= \frac{1}{k^3} \sum_{i \in G_k^0} \text{Var} \left[\sum_{s=1}^T \sum_{t=1}^T M_T(\cdot) \mathbf{v}' \{\Delta_{is} \Delta_{is}' \Delta_{it} \Delta_{it} - \mathbb{E}(\Delta_{is} \Delta_{is}' \Delta_{it} \Delta_{it})\} \right] \\ &\leq \frac{1}{k^3} \sum_{i \in G_k^0} \sum_{s=1}^T \sum_{t=1}^T \sum_{r=1}^T \sum_{l=1}^T M_T(\cdot) M_T(\cdot) \mathbf{v}' \mathbb{E}[\Delta_{is} \Delta_{is}' \Delta_{it} \Delta_{it} \Delta_{il} \Delta_{il}' \Delta_{ir} \Delta_{ir}] \\ &\leq \frac{1}{k^3} \sum_{i \in G_k^0} \sum_{|s-t| \leq M_T} \sum_{|r-l| \leq M_T} \|\mathbb{E}[\mathbf{v}' \Delta_{is} \Delta_{is}' \Delta_{it} \Delta_{it} \Delta_{il} \Delta_{il}' \Delta_{ir} \Delta_{ir}]\| \\ &= \left(\frac{2}{T} \right) = (1) \end{aligned}$$

where the last equality follows from the fact that $\|\mathbb{E}[\mathbf{v}' \Delta_{is} \Delta_{is}' \Delta_{it} \Delta_{it} \Delta_{il} \Delta_{il}' \Delta_{ir} \Delta_{ir}]\| \leq \max_{i,s} \mathbb{E} \|\Delta_{is} \Delta_{is}'\|^4 \frac{1}{2} \times \max_{i,t} \mathbb{E} \|\Delta_{it} \Delta_{it}'\|^4 \frac{1}{2} \infty$ by Assumption D2(iii). It follows that $k_{NT,2} = P(1)$

By Corollary 3.3 and the Davydov inequality,

$$\begin{aligned} \|k_{NT,3}\| &= \frac{\left\| \frac{-1}{k} - \frac{\tilde{-1}}{k} \right\|}{3/2 \left(\frac{-1/2}{k} + \frac{\tilde{-1/2}}{k} \right)} \left\| \sum_{i \in G_k^0} \sum_{s=1}^T \sum_{t=1}^T M_T(\cdot) \mathbb{E}(\Delta_{is} \Delta_{is}' \Delta_{it} \Delta_{it}) \right\| \\ &\leq \frac{\left\| \frac{\tilde{-1}}{k} - \frac{-1}{k} \right\|}{1/2 \left(\frac{-1/2}{k} + \frac{\tilde{-1/2}}{k} \right)} \left\{ \frac{1}{k} \sum_{i \in G_k^0} \sum_{|s-t| \leq M_T} \|\mathbb{E}(\Delta_{is} \Delta_{is}' \Delta_{it} \Delta_{it})\| \right\} \\ &= P\left(\frac{-1/2}{k} - 1/2 \right) (1) = P(1) \end{aligned}$$

By Assumptions D2(i)-(iii) and the Davydov inequality,

$$\begin{aligned}
\|k_{NT,4}\| &= \left\| \frac{1}{k} \sum_{i \in G_k^0} \sum_{s=1}^T \sum_{t=1}^T [1 - M_T(\cdot)] \mathbb{E}(\Delta_{is} \Delta_{it}) \right\| \\
&\leq \frac{8}{k} \sum_{i \in G_k^0} \sum_{|s-t| > M_T} i(|s-t|)^{(2q-1)/(2q)} \|\Delta_{is} \Delta_{it}\|_{4q} \\
&\leq \frac{8}{k} \sum_{i \in G_k^0} i(T)^{(2q-1)/(2q)} = (1)
\end{aligned}$$

This completes the proof of the proposition. ■

With the above result in hand, we can readily show that

$$\begin{aligned}
\sqrt{k}(\tilde{c}_k - c_0) &= \sqrt{k}(\tilde{c}_k - c_0) - \frac{1}{k} k_{NT} + \left(\tilde{c}_k\right)^{1/2} \frac{1}{k} k_{NT} - \frac{1}{k} k_{NT} \\
&\quad + \left[1 - \left(\tilde{c}_k\right)^{1/2}\right] \frac{1}{k} k_{NT} \\
&= \sqrt{k}(\tilde{c}_k - c_0) - \frac{1}{k} k_{NT} + P(1) + P\left(\tilde{c}_k^{-1}\right) \left(\tilde{c}_k\right)^{1/2} \\
&= \sqrt{k}(\tilde{c}_k - c_0) - \frac{1}{k} k_{NT} + P(1)
\end{aligned}$$

That is, $\sqrt{k}(\tilde{c}_k - c_0)$ has the desired limiting distribution centered on the origin.

Now, we demonstrate (S2). Let $\Delta_{is} = \Delta_{is} - \mathbb{E}(\Delta_{is})$ and $\Delta_{it} = \Delta_{it} - \mathbb{E}(\Delta_{it}) = 0$ and $\mathbb{E}(\Delta_{it}) = 0$ we have

$$\begin{aligned}
k_{NT} &= \frac{1}{k} \sum_{i \in G_k^0} \sum_{s=1}^T \sum_{t=1}^T [\Delta_{is} \Delta_{it} - \mathbb{E}(\Delta_{is} \Delta_{it})] \\
&= \frac{1}{k} \sum_{i \in G_k^0} \sum_{t=1}^T [\Delta_{it} \Delta_{it} - \mathbb{E}(\Delta_{it} \Delta_{it})] + \frac{1}{k} \sum_{i \in G_k^0} \sum_{1 \leq s < t \leq T} [\Delta_{is} \Delta_{it} - \mathbb{E}(\Delta_{is} \Delta_{it})] \\
&\quad + \frac{1}{k} \sum_{i \in G_k^0} \sum_{1 \leq t < s \leq T} [\Delta_{is} \Delta_{it} - \mathbb{E}(\Delta_{is} \Delta_{it})] \\
&\equiv k_{NT,1} + k_{NT,2} + k_{NT,3} \text{ say.}
\end{aligned}$$

It is trivial to show that $k_{NT,1} = P(-1)$ by the Chebyshev and Davydov inequalities. For $k_{NT,2}$ we have $\mathbb{E}(k_{NT,2}) = 0$ by construction, and by Assumption D2(ii) and Jensen inequality

$$\mathbb{E}(k_{NT,2}^2) = \frac{1}{k} \sum_{i \in G_k^0} \text{Var} \left(\sum_{1 \leq t_1 < t_2 \leq T} [\Delta_{it_1} \Delta_{it_2}] \right) \quad \square \square$$

kNT, α

and obtain the updated estimate $(\hat{\beta}_1^{(r,1)})$ of (β_1) . Next we choose (β_2) to minimize

$${}_{1NT, \lambda_1}^{(r,2,K)}(\beta_2) = {}_{1NT}(\beta) + \frac{1}{\lambda_1} \sum_{i=1}^N \left\| \hat{\alpha}_i - \hat{\alpha}_1^{(r,1)} \right\| + \sum_{k \neq 1, 2}^K \left\| \hat{\alpha}_i - \hat{\alpha}_k^{(r-1)} \right\|$$

and obtain the updated estimate $(\hat{\beta}_2^{(r,2)})$ of (β_2) . Repeat this procedure until we choose (β_K) to minimize

$${}_{1NT, \lambda_1}^{(r,K,K)}(\beta_K) = {}_{1NT}(\beta) + \frac{1}{\lambda_1} \sum_{i=1}^N \left\| \hat{\alpha}_i - \hat{\alpha}_1^{(r,K)} \right\| + \sum_{k=1}^{K-1} \left\| \hat{\alpha}_i - \hat{\alpha}_k^{(r,k)} \right\|$$

and obtain the updated estimate $(\hat{\beta}_K^{(r,K)})$ of (β_K) . Let $\hat{\alpha}^{(r)} = (\hat{\alpha}_1^{(r)}, \dots, \hat{\alpha}_K^{(r)})$ and $\hat{Q}_{1NT}^{(r,K)} = \sum_{k=1}^K {}_{1NT, \lambda_1}^{(r,k,K)}(\hat{\beta}_k^{(r,k)}, \hat{\alpha}_k^{(r)})$. Update the iteration index from r to $r+1$.

Repeat Step 2 until a convergence criterion is achieved, e.g., when

$$\left| \hat{Q}_{1NT}^{(r-1,K)} - \hat{Q}_{1NT}^{(r,K)} \right| < \text{tol} \quad \text{and} \quad \frac{\sum_{k=1}^K \left\| \hat{\alpha}_k^{(r)} - \hat{\alpha}_k^{(r-1)} \right\|^2}{\sum_{k=1}^K \left\| \hat{\alpha}_k^{(r-1)} \right\|^2 + 0.0001} < \text{tol}$$

where tol is some prescribed tolerance level (e.g., 0.0001). Define the final iterative estimate of α as $\hat{\alpha} = (\hat{\alpha}_1^{(R)}, \dots, \hat{\alpha}_K^{(R)})$ for a sufficiently large r such that the convergence criterion is met. Intuitively, individual i is classified to group k if $\hat{\alpha}_i^{(R,k)} = \hat{\alpha}_k^{(R)}$; otherwise, $\hat{\alpha}_i$ is assigned to be the $\hat{\alpha}_k^{(R)}$ that is closest to some $\hat{\alpha}_i^{(R,l)}$, $l = 1, \dots, K$. In either case, we can write the individual estimate as $\hat{\alpha}_i = \hat{\alpha}_{k^*}^{(R)}$, where $k^* = \text{argmin}_{k \in \{1, \dots, K\}} \left\| \hat{\alpha}_i^{(R, l^*(k))} - \hat{\alpha}_k^{(R)} \right\|$ and $l^*(i) = \text{argmin}_{l \in \{1, \dots, K\}} \left\| \hat{\alpha}_i^{(R,l)} - \hat{\alpha}_k^{(R)} \right\|$.

The optimization of ${}_{1NT, \lambda_1}^{(r,k,K)}(\beta_k)$ is conducted on the $(K+1)$ -dimensional parameter space for (β_k) . When β_k is non-trivial, this is a high-dimensional optimization problem. Obviously, in the penalty term $\sum_{i=1}^N \dots$ and $\sum_{k=1}^K$ are jointly convex, given $\sum_{l=k+1}^K \left\| \hat{\alpha}_i^{(r-1)} - \hat{\alpha}_l^{(r-1)} \right\|$ and $\sum_{l=1}^{k-1} \left\| \hat{\alpha}_i^{(r)} - \hat{\alpha}_l^{(r)} \right\|$ for each $i = 1, \dots, N$. If ${}_{1NT}(\beta)$ is convex in β , then ${}_{1NT, \lambda_1}^{(r,k,K)}(\beta_k)$ as the summation of ${}_{1NT}(\beta)$ and the penalty, is also convex in (β_k) . Convexity can substantially reduce the computational burden of high-dimensional optimization.

A convex ${}_{1NT}(\beta)$ is common in panel data models. Convexity apparently holds in the linear models in Examples 1 and 2. It also holds in the nonlinear models in Example 3 with $\text{logit}(\cdot)$ as the standard logistic or normal CDF, and in Example 4 after re-parameterizing the original parameter (β_i, β_i^2) into $(\beta_i = \beta_i^2, \beta_i^2 = \beta_i^2, \beta_i^3 = 1 - \beta_i^2)$. We utilize the convexity throughout our numerical works.

Given the convexity in each substep (3), the proposed algorithm consists of a sequence of convex problems implemented in an iterative manner. In particular, the only difference between the standard Lasso and a single substep of PPL is that Lasso shrinks the coefficients to a known center (zero), while the center of PPL is determined in the convex programming. Thus a PPL iteration has the same computational complexity as Lasso, which is $O(N^3)$ in our context of panel linear regression (Efron, Hastie and Johnstone, 2004, p.443). The computational cost of a single iteration is minimal.

Since the additive-multiplicative penalty is not jointly convex in all the parameters (β, α) , we can take advantage of convexity in each substep for (β_k) but not simultaneously for (β, α) . As a consequence of

In this section we carry out two more simulation exercises, one using PGMM to estimate a static panel model with endogenous regressors as in DGP 4 below, and the other using PLS to estimate the linear panel AR(1) in DGP2.

(Linear static panel with endogeneity.) We maintain the linear panel structure model with two explanatory variables as in DGP 1, but the first regressors is endogenous as it is generated from the following underlying reduced-form equation: $y_{it1} = 0.2 \frac{0}{i} + 0.5 \text{it1} + 0.5 \text{it2} + 0.5 \text{it}$ where it1 and it2 are two excluded instrumental variables, and the reduced-form error it and the structural-equation idiosyncratic shock it follow a bivariate normal distribution:

$$\begin{pmatrix} \text{it} \\ \text{it} \end{pmatrix} \sim \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & 0.3 \\ 0.3 & 1 \end{pmatrix} \right)$$

The second regressor it2 is exogenous, and $(\text{it2}, \text{it1}, \text{it2}) \sim \text{IID}(0, 3)$ is independent of (it, it) . All variables are independent across i and t . The econometrician observes $(\text{it}, \text{it1}, \text{it2}, \text{it1}, \text{it2})$. The true coefficients of the three groups are $(0.2, 1, 8)$, $(1, 1)$ and $(1, 8, 0, 2)$, respectively.

We report the statistics in Tables S1 and S2, which correspond to Tables 1 and 2, respectively, in the main text. The choice of tuning parameters are exactly the same as described in Section 4. When we compare PLS estimation with PGMM in DGP 2, we find that the PLS works better in determining the correct number of groups and in classifying the individual units. The 95% coverage probabilities are comparable to those of PGMM when $N = 50$, but are lower than PGMM when N is small. Similar to PPL in DGP 3, the lower coverage probabilities is caused by the bias. The analytical bias correction removes the bias asymptotically, but the effect is limited when N is small, as is shown in the oracle. The post-Lasso has larger coverage probability than the oracle, as the estimated standard deviation is inflated by a few misclassified units.

Table S1: Frequency of selecting $K = 1$ and 5 groups when $K_0 = 3$

N	T	DGP 4				DGP 2 (PLS)			
		1	2	4	5	1	2	4	5
100	15	0	0.022	0.076	0	0	0.106	0	0
100	25	0	0	0.028	0.006	0	0	0	0
100	50	0	0	0.004	0	0	0	0	0
200	15	0	0	0.058	0.002	0	0	0	0
200	25	0	0	0.046	0.004	0	0	0	0
200	50	0	0	0.006	0	0	0	0	0

Table S3 reports the RMSE and bias of α_1 from post-Lasso and C-Lasso under the true K_0 and the IC-determined \hat{K} (or \tilde{K} for PGMM). These estimates are bias corrected whenever necessary in the DGPs. For example, the RMSE of PPL under K_0 is calculated as $\left(\frac{1}{S} \sum_{s=1}^S \sum_{k=1}^{K_0} \frac{\hat{N}_k^{(s)}}{N} \left(\hat{\alpha}_{k,1}^{(s)} - \alpha_{k,1} \right)^2 \right)^{1/2}$ where $\hat{\alpha}_{k,1}$ and

Table S2: Classification and point estimation of α_1 in additional simulations

	N	T	% of correct classification	Post-Lasso			Oracle		
				RMSE	Bias	Coverage	RMSE	Bias	Coverage
DGP 4	100	15	0.8287	0.1583	0.0462	0.7850	0.0806	0.0018	0.9344
	100	25	0.9281	0.0883	0.0195	0.8880	0.0617	0.0009	0.9380
	100	50	0.9885	0.0517	0.0075	0.9406	0.0437	-0.0012	0.9422
	200	15	0.8378	0.1155	0.0484	0.7860	0.0577	-0.0016	0.9454
	200	25	0.9320	0.0643	0.0199	0.8742	0.0436	0.0001	0.9506
	200	50	0.9881	0.0364	0.0074	0.9356	0.0311	-0.0005	0.9450
DGP 2 (PLS)	100	15	0.8907	0.0413	0.0061	0.9148	0.0352	0.0041	0.8524
	100	25	0.9511	0.0261	0.0041	0.9710	0.0241	0.0028	0.9076
	100	50	0.9908	0.0160	0.0015	0.9908	0.0156	0.0013	0.9334
	200	15	0.8949	0.0294	0.0064	0.9154	0.0253	0.0052	0.8576
	200	25	0.9520	0.0188	0.0037	0.9714	0.0178	0.0036	0.8808
	200	50	0.9912	0.0113	0.0017	0.9934	0.0111	0.0015	0.9282

of post-Lasso, although C-Lasso appears to have larger RMSE in the PGMM estimation of DGP 2, where it does not enjoy the oracle property.

When $\gamma \neq 0$, we generalize the definition of the set of true group-specific parameters. For $\gamma = 0$, we shrink $\alpha_1^0 = \begin{pmatrix} 0 \\ 1,1 \end{pmatrix} \begin{pmatrix} 0 \\ K_{0,1} \end{pmatrix}$ into a $K_{0,1}$ -element subset $\alpha_1^0(\gamma)$. For $\gamma \neq 0$, we augment α_1^0 by adding $K_{0,1} - K_{0,1}$ elements choosing from $\begin{pmatrix} 0 \\ k,1 \end{pmatrix} \begin{pmatrix} 0 \\ K_{0,1} \end{pmatrix}$ so that the resulting $\alpha_1^0(\gamma)$ contains α_1^0 . Elements are eliminated or concatenated in each replication to fit $\hat{\alpha}(\hat{\gamma}^{(s)})$. In this scenario, the RMSE is calculated as $\left(\frac{1}{S} \sum_{s=1}^S \sum_{k=1}^K \frac{\hat{N}_k^{(s)}}{N} \left(\hat{\alpha}_{k,1}(\hat{\gamma}^{(s)}) - \alpha_{k,1}(\gamma) \right)^2 \right)^{1/2}$. According to the simulation, the effect of not knowing $\gamma = 0$ is noticeable when $K = 15$ in the linear models and $K = 25$ in the nonlinear model, but it does not necessarily enlarges the RMSE, for the estimator under $\gamma = 0$ is also noisy when K is small. The discrepancy of the RMSE and bias between $\gamma = 0$ and $\hat{\gamma}$ (or $\tilde{\gamma}$) quickly vanishes when N grows.

f

All data are downloaded from the World Bank.⁴ We extract all countries with all the variables in (5.1) available. Using the time span 1995–2010, we were able to construct a balanced panel of 57 countries. We remove one outlier Bulgaria, whose 1997 economic collapse produced hyperinflation in the CPI that significantly distorted the overall mean and the standard deviation. In total we collect 56 countries. The summary statistics are shown in Table S4.

In the implementation, we scale-normalize all the variables for each individual unit to guarantee that the coefficients are comparable. Moreover, in PGMM we use Δ_{t-2} and a constant as two excluded IVs. Although the constant is uncorrelated with the endogenous variable, adding it here stabilizes the post-Lasso estimation in finite samples.

Table S5 displays the group membership. The country names in bold are the 47 coincidences of PLS and PGMM classification out of the total 56 countries.

⁴<http://data.worldbank.org/data-catalog/world-development-indicators>.

Table S3: Estimation of α_1 by post-lasso and C-Lasso under α_0 and $\hat{\alpha}$ or $\tilde{\alpha}$

	N	T	Post-Lasso				C-Lasso				Oracle	
			K	K_0	K	K	K	K_0	K	K	RMSE	Bias
			RMSE	Bias	RMSE	Bias	RMSE	Bias	RMSE	Bias		
DGP 1	100	15	0.0596	0.0108	0.0829	0.0092	0.0619	0.0133	0.0839	0.0120	0.0463	0.0012
	100	25	0.0385	0.0019	0.0385	0.0019	0.0396	0.0040	0.0396	0.0040	0.0353	0.0001
	100	50	0.0249	0.0000	0.0249	0.0000	0.0255	0.0011	0.0255	0.0011	0.0245	-0.0002
	200	15	0.0434	0.0079	0.1373	0.0081	0.0457	0.0107	0.1353	0.0114	0.0324	-0.0013
	200	25	0.0273	0.0015	0.0273	0.0015	0.0280	0.0040	0.0280	0.0040	0.0250	-0.0006
	200	50	0.0174	-0.0001	0.0174	-0.0001	0.0181	0.0011	0.0181	0.0011	0.0171	-0.0002
DGP 2	100	15	0.0848	-0.0090	0.0787	-0.0016	0.1311	-0.0372	0.1188	-0.0250	0.0502	-0.0037
(PGMM)	100	25	0.0556	-0.0055	0.0561	-0.0051	0.1042	-0.0267	0.1045	-0.0255	0.0351	0.0011
	100	50	0.0278	-0.0012	0.0278	-0.0012	0.0418	-0.0130	0.0418	-0.0130	0.0242	-0.0010
	200	15	0.0712	-0.0141	0.0743	-0.0145	0.1491	-0.0399	0.1483	-0.0383	0.0352	-0.0017
	200	25	0.0333	-0.0051	0.0333	-0.0051	0.0932	-0.0284	0.0932	-0.0284	0.0252	-0.0006
	200	50	0.0193	-0.0014	0.0193	-0.0014	0.0277	-0.0134	0.0277	-0.0134	0.0164	0.0000
DGP 3	100	25	0.1722	0.0587	0.1516	0.0727	0.2154	0.0615	0.1641	0.0688	0.1077	0.0114
	100	50	0.0853	0.0379	0.0878	0.0383	0.1178	0.0487	0.1191	0.0489	0.0752	0.0090
	200	25	0.1342	0.0483	0.1401	0.0649	0.1826	0.0487	0.1441	0.0573	0.0821	0.0116
	200	50	0.0632	0.0264	0.0632	0.0264	0.0948	0.0372	0.0948	0.0372	0.0573	0.0121
DGP 4	100	15	0.1691	0.0487	0.1803	0.0376	0.2148	0.1087	0.2102	0.0941	0.0806	0.0018
	100	25	0.0724	0.0189	0.1217	0.0207	0.0882	0.0523	0.1323	0.0539	0.0617	0.0009
	100	50	0.0450	0.0031	0.0645	0.0042	0.0532	0.0204	0.0707	0.0215	0.0437	-0.0012
	200	15	0.1271	0.0512	0.1348	0.0466	0.1777	0.1128	0.1793	0.1074	0.0577	-0.0016
	200	25	0.0513	0.0153	0.1392	0.0235	0.0720	0.0498	0.1485	0.0577	0.0436	0.0001
	200	50	0.0314	0.0036	0.0549	0.0049	0.0399	0.0221	0.0602	0.0234	0.0311	-0.0005
DGP 2	100	15	0.0482	0.0081	0.0487	0.0065	0.0747	0.0297	0.0715	0.0254	0.0352	0.0041
(PLS)	100	25	0.0263	0.0043	0.0263	0.0043	0.0418	0.0189	0.0418	0.0189	0.0241	0.0028
	100	50	0.0160	0.0016	0.0160	0.0016	0.0218	0.0085	0.0218	0.0085	0.0156	0.0013
	200	15	0.0295	0.0064	0.0295	0.0064	0.0567	0.0293	0.0567	0.0293	0.0253	0.0052
	200	25	0.0188	0.0037	0.0188	0.0037	0.0307	0.0174	0.0307	0.0174	0.0178	0.0036
	200	50	0.0113	0.0017	0.0113	0.0017	0.0171	0.0084	0.0171	0.0084	0.0111	0.0015

Table S6: Summary statistics for the civil war dataset

	mean	median	s.e.	min	max
Civil war incidence	0.352	0	0.478	0	1
GDP per capita growth	0.020	0.024	0.040	-0.811	0.306
Population growth	0.012	0.015	0.076	-0.507	0.661

“low-occurrence” groups with results as follows.

(23 countries): Guatemala, Peru, Argentina, Mali, Senegal, Chad, Congo (Dem.), Congo (Rep.), Somalia, Morocco, Sudan, Turkey, Iraq, Lebanon, Afghanistan, China, Pakistan, Sri Lanka, Nepal, Cambodia, Laos, Philippines, Indonesia

(15 countries): Haiti, Dominican, El Salvador, Nicaragua, UK, Yugoslavia, Cyprus, Russia, Liberia, Nigeria, Central African Republic, Ethiopia, South Africa, Iran, Jordan

In this section, we use the data provided by Bonhomme and Manresa (2015) to revisit the link between income growth and democracy across countries. Following BM’s Equation (22), we specify a linear dynamic model, where the dependent variable is a country’s democracy index (measured by Freedom House indicator between 0 (the lowest) and 1 (the highest)), and the explanatory variables are the first-order lagged democracy index and the income (measured by the logarithm of GDP per capita).

The dataset contains a balanced panel of 84 countries and 8 periods at a five year interval over 1965–2000. We use PLS to estimate the model in this short panel. Many developed countries, such as the United States or United Kingdom, kept their democracy index at the highest level throughout the time. Due to the lack of within-group variation in these countries, we scale normalize each variable by its pooled standard deviation. This standardization makes sure that the parameter α_{it-1} can still be interpreted as the autoregressive coefficient, and the magnitude is comparable with the income coefficient.

Table S7: Summary statistics for the democracy dataset

	mean	median	s.e.	min	max
Democracy index	0.5760	0.6667	0.3712	0	1
GDP per capita (in logarithm)	8.2981	8.3039	1.0685	6.0937	10.4450

Following practice in the simulation, the IC with $\lambda_{NT} = \frac{2}{3}(\lambda_1)^{-1/2}$ picks out $\lambda_1 = 3$ and $\lambda_1 = 1/20$ in all combinations of $\lambda_1 = 1/5$ and λ_1 in a geometrically increasing sequence of 10 points in $(0, 2, \dots, 2)$. Under $\lambda_1 = 3$ and $\lambda_1 = 1/20$, C-Lasso categorizes the 84 countries into the following groups:

(30 countries): Belgium, Bolivia, Brazil, Canada, Dominican, Ecuador, El Salvador, Finland, Guatemala, Guinea, Iceland, Indonesia, Italy, Japan, Jordan, Luxembourg, Mali, Morocco, Nepal, Panama, Peru, Philippines, Portugal, Romania, South Africa, Thailand, Turkey, United Kingdom, Uruguay, Venezuela

(36 countries): Algeria, Argentina, Australia, Austria, Barbados, Burkina Faso, Burundi, Cameroon, Chile, China, Colombia, Costa Rica, Cote d’Ivoire, Denmark, Egypt, France, Gabon, Ghana, Greece, India, Iran, Israel, Jamaica, Kenya, Malawi, Malaysia, Mexico, Nigeria, Norway, Paraguay, Rwanda, Spain, Sweden, Togo, Trinidad and Tobago, United States

(18 countries): Benin, Chad, Congo (Rep.), Honduras, Ireland, Korea (Rep.), Madagascar, Netherlands, New Zealand, Nicaragua, Niger, Sri Lanka, Switzerland, Syrian, Tanzania, Tunisia, Uganda, Zambia

The post-Lasso and pooled FE estimates are shown in Table S8. We focus on the coefficient for income. The common FE coefficient is positive and significant. The positive effect is echoed by Groups 1 and 2, but contrasts with Group 3, which consists mainly of low-income and low-democracy nations combined with a few selected OECD countries. OECD countries such as Ireland, Netherlands, New Zealand and Switzerland maintained their democracy index at 1 throughout the sample period. The lack of variation in the dependent variable makes them uninformative about the income coefficient.

Table S8: PLS estimation results

	Pooled FE		PLS					
	coef.	s.e.	Group 1		Group 2		Group 3	
			coef.	s.e.	coef.	s.e.	coef.	s.e.
Lagged democracy	. ***	0.0491	. ***	0.0643	. ***	0.0733	- .	0.0521
Income	. ***	0.0489	. ***	0.0930	. ***	0.0448	- . ***	0.0860

Note: ***1% significant, ** 5% significant, * 10% significant

- BLATTMAN C. AND E. MIGUEL (2010): "Civil Wars," *Journal of Economic Literature* 48, 3-57.
- BONHOMME, S., AND E. MANRESA (2015): "Grouped Patterns of Heterogeneity in Panel Data," *Econometrica* 83, 1147-1184.
- COLLIER, P. AND A. HOFFLER (2004): "Greed and Grievance in Civil Wars," *Oxford Economic Papers* 56, 563-595.
- DJANKOV, S. AND M. REYNAL-QUEROL (2010): "Poverty and Civil War: Revisiting the Evidence," *Review of Economics and Statistics* 92, 1035-1041.
- EFRON, B., T. HASTIE, I. JOHNSTONE and R. TIBSHIRANI (2004), "Least Angle Regression," *Annals of Statistics* 32, 407-499.
- FEARON, J.D. AND D. D. LAITIN (2003): "Ethnicity, Insurgency, and Civil War," *The American Political Science Review* 97, 75-90.
- HAHN, J., AND G. KUERSTEINER (2011): "Bias Reduction for Dynamic Nonlinear Panel Models with Fixed Effects," *Econometric Theory* 27, 1152-1191.
- HAHN, J., AND W. NEWAY (2004): "Jackknife and Analytical Bias Reduction for Nonlinear Panel Models," *Econometrica* 72, 1295-1319.
- HALL, P. AND C.C. HEYDE (1980): *Martingale Limit Theory and Its Applications*. Academic Press, New York.
- MERLEVÈDE, F., M. PEILGRAD, M., AND E. RIO (2009): "Bernstein Inequality and Moderate Deviations under Strong Mixing Conditions," *IMS collections: High Dimensional Probability V.*, 273-292.
- WHITE, H. (2001): *Asymptotic Theory for Econometricians*. Emerald, UK.