# Supplementary Appendix for "Discretizing Unobserved Heterogeneity"<sup>1</sup>

In this supplementary appendix we provide details on various aspects of the theory, and we report additional results.

## **B** Conditional first step

In this section of the appendix we study some theoretical properties of GFE based on a conditional first-step, under a series regression specification. We first study the case where  $\alpha_{i0}$  does not vary over time, and then turn to the time-varying case.

### **B.1** Convergence rate for the first step

We make the following assumption, where  $g(x, \alpha) = \mathbb{E}_{X_{it}=x, \alpha_{i0}=\alpha}[h(Y_{it}, X_{it})]$ , and  $f_{\mu}$  denotes the conditional density of  $X_{it}$  given  $\mu_{i0} = \mu$ .

#### Assumption B1.

(i) There exists an integrable function  $\omega_X$  such that  $f_{\mu}(x) \leq \omega_X(x)$  for all  $x, \mu$ .

The maximum eigenvalue of  $\int P_q(x)P_q(x)'\omega_X(x)dx$  is O(1), and  $\sup_{\alpha} \int ||g(x,\alpha)||^2 \omega_X(x)dx$  is bounded.

There exists a constant  $\underline{c} > 0$ , independent of q and T, such that the minimum eigenvalue of  $\frac{1}{T} \sum_{t=1}^{T} P_q(X_{it}) P_q(X_{it})'$  is larger than  $\underline{c} + o_p(1)$ , uniformly in *i*.

(ii) Let  $\mathcal{A}$  denote the parameter space for  $\alpha_i$ . For all  $\alpha \in \mathcal{A}$ ,  $g(\cdot, \alpha) \in \mathcal{G} \subset L^2(\omega_X)$ .  $\|g(x, \alpha') - g(x, \alpha)\| \leq C(x) \|\alpha' - \alpha\|$ , where  $\frac{1}{NT} \sum_{i=1}^{N} \sum_{t=1}^{T} C(X_{it})^2 = O_p(1)$ . There exist a constant a > 0 and a sequence of functions  $\beta_q(\cdot)$  such that  $\sup_{x,\alpha} \|g(x, \alpha) - P_q(x)'\beta_q(\alpha)\| = O(q^{-a})$ .

(*iii*) Let 
$$\varepsilon_{it} = h(Y_{it}, X_{it}) - g(X_{it}, \alpha_{i0})$$
. Then  $\frac{1}{N} \sum_{i=1}^{N} \left\| \frac{1}{T} \sum_{t=1}^{T} P_q(X_{it}) \varepsilon_{it} \right\|^2 = O_p(q/T)$ .

(iv) There exists a Lipschitz-continuous mapping  $\psi : \mathcal{G} \to \mathcal{A}$  such that  $\psi(g(\cdot, \alpha_{i0})) = \alpha_{i0}$ .

Parts (i) and (ii) in Assumption B1 will be satisfied when using a suitable family  $P_q$ . The condition on the maximum eigenvalue of  $\int P_q(x)P_q(x)'\omega_X(x)dx$  in part (i) is without loss of generality. In part (ii), a = s/r for splines and power series, where s is the number of continuous derivatives in x of

<sup>&</sup>lt;sup>1</sup>We thank the IFAU for access to, and help with, the Swedish administrative data.

 $g(x, \alpha)$ , and r is the dimension of  $X_{it}$ . Part (iii) holds under standard conditions on moments and serial dependence. Part (iv) is an injectivity condition.

To illustrate these conditions, suppose  $X_{it}$  is scalar, and consider a family of polynomials that are orthonormal with respect to the uniform distribution on [-1, 1]. If  $\omega_X$  is bounded by a constant B > 0, then  $\int P_q(x)P_q(x)'\omega_X(x)dx$  is bounded by  $BI_q$ , where  $I_q$  is the  $q \times q$  identity matrix, so its maximum eigenvalue is bounded. If in addition we assume that  $f_{\mu}(x) \geq b$  for all  $x, \mu$ , where b > 0 is a constant, then  $\mathbb{E}_{\mu_{i0}=\mu}\left(P_q(X_{it})P_q(X_{it})'\right) \geq bI_q$  for all  $\mu$  and q. Newey (1997, Theorem 1) provides conditions under which, for given i, the minimum eigenvalue of  $\frac{1}{T}\sum_{t=1}^T P_q(X_{it})P_q(X_{it})'$  is larger than  $\underline{c} + o_p(1)$ , for some constant  $\underline{c} > 0$ . Specifically, Newey's conditions require that  $X_{it}$  be i.i.d. over time (here, conditional on  $\mu_{i0}$ ), and  $\sup_x \|P_q(x)\|_F \leq \zeta_q$  with  $q\zeta_q^2/T = o(1)$ , where  $\|A\|_F = (\operatorname{Tr} A'A)^{1/2}$ denotes the Fröbenius norm of A. From this he obtains that  $\|\frac{1}{T}\sum_{t=1}^T Z_{it}\|_F = o_p(1)$ , where  $Z_{it} = P_q(X_{it})P_q(X_{it})' - \mathbb{E}_{\mu_{i0}}\left(P_q(X_{it})P_q(X_{it})'\right)$ .

In the present context, to guarantee that the minimum eigenvalue of  $\frac{1}{T} \sum_{t=1}^{T} P_q(X_{it}) P_q(X_{it})'$  be larger than  $\underline{c} + o_p(1)$  uniformly in *i*, we need the stronger uniform condition  $\max_i \| \frac{1}{T} \sum_{t=1}^{T} Z_{it} \|_F = o_p(1)$ . Results from random matrix theory are helpful to provide primitive conditions for this. For example, using the matrix Bernstein theorem (Theorem 1.4. in Tropp, 2012), under the assumption that  $X_{it}$  are independent over time conditional on  $\mu_{i0}$ , we have for all  $\epsilon > 0$ ,  $\mu$ , q, and T:

$$\Pr\left(\left\|\frac{1}{T}\sum_{t=1}^{T}Z_{it}\right\| \ge \epsilon \left\|\mu_{i0} = \mu\right) \le q \exp\left(-T\frac{\epsilon^2}{2(\zeta_q^4 + \frac{\zeta_q^2\epsilon}{3})}\right),\tag{B1}$$

where recall that ||A|| denotes the spectral norm of A. A related inequality holds under a form of serial dependence ( $\beta$ -mixing), see Banna, Merlevède and Youssef (2016). Using the union bound and the fact that  $||A||_F \leq \sqrt{q} ||A||$ , (B1) implies that  $\max_i ||\frac{1}{T} \sum_{t=1}^T Z_{it}||_F = o_p(1)$ , provided that  $q\zeta_q^4 \ln(Nq)/T \to 0$ .

Part (iv) is an injectivity condition that generalizes the one in the unconditional case. As an example, consider the probit model of Example 2:  $Y_{it} = \mathbf{1}\{X'_{it}\theta_0 + \alpha_{i0} + U_{it} \ge 0\}$ , where  $U_{it}$  are i.i.d. standard normal, independent of all  $X_{it}$ 's. Taking  $h(Y_{it}, X_{it}) = Y_{it}$ , we have  $g(X_{it}, \alpha_{i0}) = \Phi(X'_{it}\theta_0 + \alpha_{i0})$ . Provided that  $\alpha_{i0}$  and  $X_{it}$  are such that g is bounded between  $\epsilon > 0$  and  $1 - \epsilon$ , and denoting  $c_0 = \int \omega_X(x) dx > 0$  and  $c_1 = c_0^{-1} \int x \omega_X(x) dx$ , the following function satisfies part (iv):

$$\psi(g(\cdot,\alpha)) = c_0^{-1} \int \Phi^{-1}(g(x,\alpha))\omega_X(x)dx - c_1'\theta_0$$

We now state our main result on the convergence rate of the conditional kmeans estimator, where:

$$B_{\alpha}(K) = \min_{(\tilde{\alpha}, k_1, \dots, k_N)} \frac{1}{N} \sum_{i=1}^{N} \|\alpha_{i0} - \tilde{\alpha}(k_i)\|^2$$

denotes the approximation error associated with  $\alpha_{i0}$  alone.

**Theorem B1.** Suppose that  $\alpha_{i0}$  is time-invariant, of fixed dimension. Let parts (i) to (iii) in Assumption B1 hold. Then, as N, T, K, q tend to infinity:

$$\frac{1}{N}\sum_{i=1}^{N}\int \left\|P_{q}(x)'\hat{\beta}_{q}(\hat{k}_{i}) - g(x,\alpha_{i0})\right\|^{2}\omega_{X}(x)dx = O_{p}(q/T) + O_{p}(q^{-2a}) + O_{p}(B_{\alpha}(K))$$

If in addition part (iv) in Assumption B1 holds, then:

$$\frac{1}{N}\sum_{i=1}^{N} \left\| \psi \left( P_q(\cdot)'\hat{\beta}_q(\hat{k}_i) \right) - \alpha_{i0} \right\|^2 = O_p(q/T) + O_p(q^{-2a}) + O_p(B_\alpha(K)).$$

Note that the second part in Theorem B1 implies (19) in the main text. As a special case, Theorem B1 implies that, when  $g(x, \alpha) = P_q(x)'\beta_q(\alpha)$  (i.e., there is no remainder term), then the rate is  $O_p(1/T) + O_p(B_\alpha(K))$ . As an example, this rate is achieved when  $X_{it}$  has finite support and we use indicator functions as  $P_q$ . More generally, when g is nonparametric and covariates are continuous, a larger number of covariates increases  $q^{-2a}$ , thus requiring choosing a larger q and incurring a larger q/T. Nevertheless, even in this nonparametric case, the resulting rate may improve relative to the one of our baseline two-step estimation.

To give a concrete example, assume that there are two covariates and g has four continuous derivatives, so a = 2. Suppose also that  $\alpha_{i0}$  is scalar and  $\mu_{i0}$  has dimension two, and that  $(\alpha_{i0}, \mu'_{i0})$ has an absolutely continuous density, so the underlying dimension of heterogeneity is d = 3. Taking  $q \propto T^{1/5}$  we obtain a rate of  $O_p(T^{-4/5}) + O_p(K^{-2})$  in the conditional case. In contrast, the corresponding rate in the unconditional case is  $O_p(T^{-1}) + O_p(K^{-2/3})$ . When N and T grow at the same rate, the best rates achievable are  $O_p(T^{-4/5})$  and  $O_p(T^{-2/3})$ , respectively (since  $K \leq N$ ). In addition, the required numbers of groups to achieve these rates are  $K \propto T^{2/5}$  and T, respectively. This shows that the conditional method achieves a better rate, and that the number of groups needed to achieve that rate is smaller. As in other applications of series methods, smoothness assumptions are critical to ensure good performance. Under sufficient smoothness, the theoretical difference between the two methods increases as the dimension of  $X_{it}$  increases.

#### Proof of Theorem **B1**. Let:

$$(\underline{\beta}_q, \{\underline{k}_i\}) = \underset{(b,\{k_i\})}{\operatorname{argmin}} \sum_{i=1}^N \sum_{t=1}^T \left\| g(X_{it}, \alpha_{i0}) - P_q(X_{it})' b(k_i) \right\|^2.$$

Let also:

$$(\underline{\alpha}, \{\underline{k}_i^*\}) = \underset{(\widetilde{\alpha}, \{k_i\})}{\operatorname{argmin}} \sum_{i=1}^N \|\alpha_{i0} - \widetilde{\alpha}(k_i)\|^2.$$

By Assumption B1 (ii), we have, for a and  $\beta_q(\cdot)$  defined in that assumption:

$$\frac{1}{NT} \sum_{i=1}^{N} \sum_{t=1}^{T} \left\| g(X_{it}, \alpha_{i0}) - P_q(X_{it})' \underline{\beta}_q(\underline{k}_i) \right\|^2 \\
\leq \frac{1}{NT} \sum_{i=1}^{N} \sum_{t=1}^{T} \left\| g(X_{it}, \alpha_{i0}) - P_q(X_{it})' \beta_q(\underline{\alpha}(\underline{k}_i^*)) \right\|^2 \\
\leq \frac{1}{NT} \sum_{i=1}^{N} \sum_{t=1}^{T} \left\| g(X_{it}, \alpha_{i0}) - g(X_{it}, \underline{\alpha}(\underline{k}_i^*)) \right\|^2 + O_p(q^{-2a}) \\
= O_p(B_\alpha(K)) + O_p(q^{-2a}).$$
(B2)

Moreover, since:

$$\frac{1}{NT}\sum_{i=1}^{N}\sum_{t=1}^{T}\left\|h(Y_{it}, X_{it}) - P_q(X_{it})'\widehat{\beta}_q(\widehat{k}_i)\right\|^2 \le \frac{1}{NT}\sum_{i=1}^{N}\sum_{t=1}^{T}\left\|h(Y_{it}, X_{it}) - P_q(X_{it})'\underline{\beta}_q(\underline{k}_i)\right\|^2,$$

we have, using that  $h(Y_{it}, X_{it}) = g(X_{it}, \alpha_{i0}) + \varepsilon_{it}$ :

$$\frac{1}{NT} \sum_{i=1}^{N} \sum_{t=1}^{T} \left\| g(X_{it}, \alpha_{i0}) - P_q(X_{it})' \widehat{\beta}_q(\widehat{k}_i) \right\|^2 \\
\leq \frac{1}{NT} \sum_{i=1}^{N} \sum_{t=1}^{T} \left\| g(X_{it}, \alpha_{i0}) - P_q(X_{it})' \underline{\beta}_q(\underline{k}_i) \right\|^2 \\
+ \frac{2}{NT} \sum_{i=1}^{N} \sum_{t=1}^{T} \varepsilon'_{it} \left( P_q(X_{it})' \widehat{\beta}_q(\widehat{k}_i) - P_q(X_{it})' \underline{\beta}_q(\underline{k}_i) \right).$$
(B3)

Hence, using parts (ii) and (iii) in Assumption B1, equation (B2), and the Cauchy Schwarz and triangle inequalities, we obtain:

$$\begin{split} &\frac{1}{NT}\sum_{i=1}^{N}\sum_{t=1}^{T}\left\|P_{q}(X_{it})'\beta_{q}(\alpha_{i0})-P_{q}(X_{it})'\widehat{\beta}_{q}(\widehat{k}_{i})\right\|^{2} \\ &\leq O_{p}(1)\left(\frac{1}{N}\sum_{i=1}^{N}\left\|\frac{1}{T}\sum_{t=1}^{T}P_{q}(X_{it})'\varepsilon_{it}\right\|^{2}\right)^{\frac{1}{2}}\left(\frac{1}{N}\sum_{i=1}^{N}\left\|\widehat{\beta}_{q}(\widehat{k}_{i})-\underline{\beta}_{q}(\underline{k}_{i})\right\|^{2}\right)^{\frac{1}{2}} \\ &+O_{p}(q^{-2a})+O_{p}(B_{\alpha}(K)) \\ &\leq O_{p}\left(\sqrt{\frac{q}{T}}\right)\left(\frac{2}{N}\sum_{i=1}^{N}\left\|\widehat{\beta}_{q}(\widehat{k}_{i})-\beta_{q}(\alpha_{i0})\right\|^{2}+\frac{2}{N}\sum_{i=1}^{N}\left\|\underline{\beta}_{q}(\underline{k}_{i})-\beta_{q}(\alpha_{i0})\right\|^{2}\right)^{\frac{1}{2}} \\ &+O_{p}(q^{-2a})+O_{p}(B_{\alpha}(K)). \end{split}$$

It thus follows from part (i) in Assumption B1 that:

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^{N} \left\| \widehat{\beta}_{q}(\widehat{k}_{i}) - \beta_{q}(\alpha_{i0}) \right\|^{2} \\ &\leq O_{p}\left(\sqrt{\frac{q}{T}}\right) \left( \frac{2}{N} \sum_{i=1}^{N} \left\| \widehat{\beta}_{q}(\widehat{k}_{i}) - \beta_{q}(\alpha_{i0}) \right\|^{2} + \frac{2}{N} \sum_{i=1}^{N} \left\| \underline{\beta}_{q}(\underline{k}_{i}) - \beta_{q}(\alpha_{i0}) \right\|^{2} \right)^{\frac{1}{2}} \\ &+ O_{p}(q^{-2a}) + O_{p}(B_{\alpha}(K)), \end{aligned}$$

hence that:

$$\frac{1}{N} \sum_{i=1}^{N} \left\| \widehat{\beta}_{q}(\widehat{k}_{i}) - \beta_{q}(\alpha_{i0}) \right\|^{2} \\
\leq O_{p} \left( \frac{1}{N} \sum_{i=1}^{N} \left\| \underline{\beta}_{q}(\underline{k}_{i}) - \beta_{q}(\alpha_{i0}) \right\|^{2} \right) + O_{p} \left( \frac{q}{T} \right) + O_{p}(q^{-2a}) + O_{p}(B_{\alpha}(K)) \\
= O_{p} \left( \frac{q}{T} \right) + O_{p}(q^{-2a}) + O_{p}(B_{\alpha}(K)),$$

where the last identity comes from the fact that, by (B2) and parts (i) and (ii) in Assumption B1:

$$\frac{1}{N}\sum_{i=1}^{N} \left\| \underline{\beta}_{q}(\underline{k}_{i}) - \beta_{q}(\alpha_{i0}) \right\|^{2} = O_{p}(q^{-2a}) + O_{p}(B_{\alpha}(K)).$$

Hence, using parts (i) and (ii) in Assumption B1 we have:

$$\frac{1}{N}\sum_{i=1}^{N}\int \left\|P_q(x)'\widehat{\beta}_q(\widehat{k}_i) - g(x,\alpha_{i0})\right\|^2 \omega_X(x)dx$$

$$\leq \frac{1}{N}\sum_{i=1}^{N}\int \left\|P_q(x)'\widehat{\beta}_q(\widehat{k}_i) - P_q(x)'\beta_q(\alpha_{i0})\right\|^2 \omega_X(x)dx + O_p(q^{-2a})$$

$$\leq O_p(1)\frac{1}{N}\sum_{i=1}^{N}\left\|\widehat{\beta}_q(\widehat{k}_i) - \beta_q(\alpha_{i0})\right\|^2 + O_p(q^{-2a})$$

$$= O_p\left(\frac{q}{T}\right) + O_p(q^{-2a}) + O_p(B_\alpha(K)).$$

Lastly, we prove the second part in Theorem B1. We have, using the first part and part (iv) in Assumption B1:

$$\frac{1}{N} \sum_{i=1}^{N} \left\| \psi \left( P_q(\cdot)' \widehat{\beta}_q(\widehat{k}_i) \right) - \alpha_{i0} \right\|^2$$
$$= \frac{1}{N} \sum_{i=1}^{N} \left\| \psi \left( P_q(\cdot)' \widehat{\beta}_q(\widehat{k}_i) \right) - \psi(g(\cdot, \alpha_{i0})) \right\|^2$$
$$\leq O_p(1) \frac{1}{N} \sum_{i=1}^{N} \int \left\| P_q(x)' \widehat{\beta}_q(\widehat{k}_i) - g(x, \alpha_{i0}) \right\|^2 \omega_X(x) dx$$
$$= O_p\left(\frac{q}{T}\right) + O_p(q^{-2a}) + O_p(B_\alpha(K)).$$

## **B.2** Choice of *K*

In the conditional case, we define:

$$\widehat{Q}_X(K) = \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T \left\| P_q(X_{it})'\widehat{\beta}(\widehat{k}_i) - P_q(X_{it})'\widehat{\beta}_{qi} \right\|^2,$$

where  $\widehat{\beta}_{qi} = \operatorname{argmin}_b \sum_{t=1}^T \|h(Y_{it}, X_{it}) - P_q(X_{it})'b\|^2$ . We let:

$$\widehat{K} = \min_{K \ge 1} \left\{ K : \widehat{Q}_X(K) \le \gamma \widehat{V}_X \right\},\tag{B4}$$

where  $\widehat{V}_X = \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T \left\| P_q(X_{it})' \widehat{\beta}_{qi} - P_q(X_{it})' \beta_q(\alpha_{i0}) \right\|^2 + o_p(q/T).$ For example, when  $\varepsilon_{it}$  are independent over time and homoskedastic, one can take:

$$\widehat{V}_X = \frac{q}{NT^2} \sum_{i=1}^N \sum_{t=1}^T \|h(Y_{it}, X_{it}) - P_q(X_{it})'\widehat{\beta}_{qi}\|^2.$$

**Corollary B1.** Let Assumption **B1** hold. Suppose that  $(1 + \gamma)\hat{V}_X = O_p(q/T)$ . Take  $K \ge \hat{K}$ . Then:

$$\frac{1}{N}\sum_{i=1}^{N}\left\|\psi\left(P_{q}(\cdot)'\widehat{\beta}_{q}(\widehat{k}_{i})\right)-\alpha_{i0}\right\|^{2}=O_{p}(q/T)+O_{p}(q^{-2a}).$$

Note that, to implement the formula for  $\hat{K}$ , we need to compute  $\hat{\beta}_{qi}$ . An alternative approach, which we use in the simulations, is to replace  $\hat{\beta}_{qi}$  by a conditional kmeans estimate  $\hat{\beta}(\hat{k}_i)$  obtained using a large number of groups  $K_{\text{max}}$ . In all the models that we consider in the simulation section (Section G), we use this approach with  $K_{\text{max}} = 30$ .

**Proof of Corollary B1.** Using (B4) and the triangle inequality, we have:

$$\frac{1}{NT}\sum_{i=1}^{N}\sum_{t=1}^{T}\left\|P_q(X_{it})'\beta_q(\alpha_{i0}) - P_q(X_{it})'\widehat{\beta}_q(\widehat{k}_i)\right\|^2$$

$$\leq \frac{2}{NT}\sum_{i=1}^{N}\sum_{t=1}^{T}\left\|P_q(X_{it})'\beta_q(\alpha_{i0}) - P_q(X_{it})'\widehat{\beta}_{qi}\right\|^2$$

$$+ \frac{2}{NT}\sum_{i=1}^{N}\sum_{t=1}^{T}\left\|P_q(X_{it})'\widehat{\beta}_{qi} - P_q(X_{it})'\widehat{\beta}_q(\widehat{k}_i)\right\|^2$$

$$= O_p\left(\widehat{V}_X\right) + o_p(q/T) + O_p\left(\gamma\widehat{V}_X\right) = O_p\left(\frac{q}{T}\right).$$

The result then follows from the same arguments as in the proof of Theorem B1.

## B.3 Time-varying heterogeneity

The conditional kmeans first step can be adapted to models where  $\alpha_{it0}$  varies over time. Let us consider two cases. A first situation is when  $\alpha_{it0} = \alpha(\xi_{i0}, \lambda_t^{(o)})$ , where  $\lambda_t^{(o)}$  is a vector of observed time-varying factors. In our theory we require time-stationarity, so  $\lambda_t^{(o)}$  needs to be a stationary factor. In this case the baseline conditional kmeans method is unchanged, subject to including  $\lambda_t^{(o)}$  in the set of covariates  $X_{it}$ . Moreover, Theorem B1 holds exactly as in the time-invariant case.

A more general situation is when  $\alpha_{it0} = \alpha(\xi_{i0}, \lambda_{t0})$ , where  $\lambda_{t0}$  is a vector of unobserved factors. Here we focus on the case where the factors may depend unrestrictedly on time (i.e., we take p = T), although it would be easy to modify the setup to allow for heterogeneity varying only between subperiods but not within, as in Theorem 1. We modify the first step algorithm in the following way.

#### Algorithm 3. (conditional kmeans, time-varying heterogeneity)

- Given initial values for  $b_q(1,1), ..., b_q(K,T)$ , iterate between the following two steps until convergence:
- Given  $b_q(1,1), ..., b_q(K,T)$ , compute  $k_i = \operatorname{argmin}_{k=1,...,K} \sum_{t=1}^T \|h(Y_{it}, X_{it}) P_q(X_{it})' b_q(k,t)\|^2$ for all *i*.
- Given  $k_1, ..., k_N$ , compute  $b_q(k, t) = \operatorname{argmin}_b \sum_{i=1}^N \mathbf{1}\{k_i = k\} \|h(Y_{it}, X_{it}) P_q(X_{it})'b\|^2$  for all k, t.

Adapting the theory requires additional assumptions and different arguments. We start with a definition.

**Definition 1.** (sub-Gaussianity) A random vector Z is sub-Gaussian if there exists a scalar constant  $\lambda > 0$  such that  $\mathbb{E}[\exp(\tau'Z)] \leq \exp(\lambda \|\tau\|^2)$  for all  $\tau \in \mathbb{R}^{\dim Z}$ .

We make the following assumption, where  $g(x, \alpha) = \mathbb{E}_{X_{it}=x, \alpha_{it0}=\alpha}[h(Y_{it}, X_{it})].$ 

#### Assumption B2.

- (i) Let  $\varepsilon_{it} = h(Y_{it}, X_{it}) g(X_{it}, \alpha_{it0})$ . The NTq×1 random vector  $(P_q(X_{it})\varepsilon_{it})_{i=1,...,N,t=1,...,T}$  satisfies Definition 1.
- (ii) There exists a Lipschitz-continuous mapping  $\psi: \mathcal{G} \to \mathcal{A}$  such that  $\psi(g(\cdot, \alpha_{it0})) = \alpha_{it0}$ .

Part (i) in Assumption B2 requires  $P_q(X_{it})\varepsilon_{it}$  to be sub-Gaussian (e.g., Vershynin, 2010). This is stronger than the corresponding requirements in Theorem 1. For example, i.i.d. Gaussian random variables and i.i.d. bounded random variables are sub-Gaussian. More generally, this assumption allows for dependence across observations. As an example, a random vector  $W \sim \mathcal{N}(0, \Sigma)$  is sub-Gaussian when the maximal eigenvalue of  $\Sigma$  is bounded from above by  $2\lambda$ . This allows for weak dependence, across both individual units and time periods.

We have the following result.

**Corollary B2.** Let Assumption 1 hold, for time-varying  $\alpha_{it0}$  and  $\alpha$  Lipschitz-continuous in its first argument. Let parts (i) and (ii) in Assumption B1 hold, and let part (i) in Assumption B2 hold. Then, as N, T, K, q tend to infinity:

$$\frac{1}{NT}\sum_{i=1}^{N}\sum_{t=1}^{T}\int \left\|P_q(x)'\widehat{\beta}_q(\widehat{k}_i,t) - g(x,\alpha_{it0})\right\|^2 \omega_X(x)dx$$
$$= O_p((\ln K)/T) + O_p(q^{-2a}) + O_p(qK/N) + O_p(B_\alpha(K)).$$

If in addition part (ii) in Assumption B2 holds, then:

$$\frac{1}{NT}\sum_{i=1}^{N}\sum_{t=1}^{T}\left\|\psi\left(P_{q}(\cdot)'\widehat{\beta}_{q}(\widehat{k}_{i},t)\right)-\alpha_{it0}\right\|^{2}$$
$$=O_{p}((\ln K)/T)+O_{p}(q^{-2a})+O_{p}(qK/N)+O_{p}(B_{\alpha}(K)).$$

From Corollary B2 we obtain a convergence rate for conditional kmeans that depends on the underlying dimension  $d_{\alpha}$  of  $\alpha_{it0}$ , instead of the dimension d that appears in Theorem 1. Also, note that, in the time-invariant case, under the assumptions of Corollary B2 one can show that the conditional kmeans estimator satisfies the following convergence rate:  $O_p(\min(q, \ln K)/T) + O_p(qK/(NT)) + O_p(q^{-2a}) + O_p(B_{\alpha}(K))$ , which is as least as fast as the rate in Theorem B1.

Proof of Corollary B2. Let:

$$(\underline{\beta}_q, \{\underline{k}_i\}) = \underset{(b,\{k_i\})}{\operatorname{argmin}} \sum_{i=1}^N \sum_{t=1}^T \left\| g(X_{it}, \alpha_{it0}) - P_q(X_{it})' b(k_i, t) \right\|^2,$$

and let:

$$(\underline{\alpha}, \{\underline{k}_i^*\}) = \underset{(\widetilde{\alpha}, \{k_i\})}{\operatorname{argmin}} \sum_{i=1}^N \sum_{t=1}^T \|\alpha_{it0} - \widetilde{\alpha}(k_i, t)\|^2.$$

We have, as in the proof of Theorem B1:

$$\frac{1}{NT} \sum_{i=1}^{N} \sum_{t=1}^{T} \left\| g(X_{it}, \alpha_{it0}) - P_q(X_{it})' \underline{\beta}_q(\underline{k}_i, t) \right\|^2 = O_p(B_\alpha(K)) + O_p(q^{-2a}).$$
(B5)

Moreover, we have:

$$\frac{1}{NT} \sum_{i=1}^{N} \sum_{t=1}^{T} \left\| g(X_{it}, \alpha_{it0}) - P_q(X_{it})' \widehat{\beta}_q(\widehat{k}_i, t) \right\|^2 \\
\leq \frac{1}{NT} \sum_{i=1}^{N} \sum_{t=1}^{T} \left\| g(X_{it}, \alpha_{it0}) - P_q(X_{it})' \underline{\beta}_q(\underline{k}_i, t) \right\|^2 \\
+ \frac{2}{NT} \sum_{i=1}^{N} \sum_{t=1}^{T} \left( P_q(X_{it}) \varepsilon_{it} \right)' \left( \widehat{\beta}_q(\widehat{k}_i, t) - \underline{\beta}_q(\underline{k}_i, t) \right). \tag{B6}$$

The main difference with the proof of Theorem B1 is how to bound the cross-product term in (B6). To do so, we apply a version the Hanson-Wright tail inequality for quadratic forms, due to Hsu, Kakade and Zhang (2012, Theorem 2.1), which allows for dependent data.

**Lemma B1.** (Hsu, Kadade and Zhang, 2012) Let Z be a random vector such that, for some  $\lambda > 0$ ,  $\mathbb{E} [\exp(\tau' Z)] \leq \exp(\lambda ||\tau||^2)$  for all  $\tau \in \mathbb{R}^{\dim Z}$ . Let Q be a positive semi-definite matrix. Then, for all s > 0:

$$\Pr\left[Z'QZ > 2\lambda \operatorname{tr} Q + 4\lambda \sqrt{s \operatorname{tr} Q^2} + s4\lambda \|Q\|\right] \le \exp(-s).$$

Let  $V_{it} = P_q(X_{it})\varepsilon_{it}$ . We will bound:

$$\frac{1}{NT}\sum_{i=1}^{N}\sum_{t=1}^{T}\left(P_{q}(X_{it})\varepsilon_{it}\right)'\left(\widehat{\beta}_{q}(\widehat{k}_{i},t)-\underline{\beta}_{q}(\underline{k}_{i},t)\right) = \frac{1}{NT}\sum_{i=1}^{N}\sum_{t=1}^{T}\overline{V}_{t}(\widehat{k}_{i},\underline{k}_{i})'\left(\widehat{\beta}_{q}(\widehat{k}_{i},t)-\underline{\beta}_{q}(\underline{k}_{i},t)\right),$$

where  $\overline{V}_t(k, k')$  denotes the *t*-specific linear projection of  $V_{it}$  on group indicators  $\mathbf{1}\{\hat{k}_i = k\}$  and  $\mathbf{1}\{\underline{k}_i = k'\}$ . For this, we will use the Cauchy Schwarz inequality and bound:

$$\frac{1}{NT}\sum_{i=1}^{N}\sum_{t=1}^{T}\left\|\overline{V}_{t}(\widehat{k}_{i},\underline{k}_{i})\right\|^{2}.$$

Let, for given partitions  $\{k_{i1}\}, \{k_{i2}\}: \frac{1}{NT} \sum_{i=1}^{N} \sum_{t=1}^{T} \left\| \overline{V}_t(k_{1i}, k_{2i}) \right\|^2 = \frac{v'Qv}{NT}$ , where we have defined  $v = (V'_{11}, \dots, V'_{NT})', Q$  is an  $NTq \times NTq$  projection matrix with tr  $Q \leq 2KTq, Q^2 = Q$ , and  $\|Q\| = 1$ . By Lemma B1 and part (i) in Assumption B2 we have, for all s:

$$\Pr\left[v'Qv > 4\lambda KTq + 4\lambda\sqrt{2KTqs} + 4\lambda s\right] \le \exp(-s)$$

so, using that  $2\sqrt{ab} \le a + b$ :

$$\Pr\left[v'Qv > 8\lambda KTq + 6\lambda s\right] \le \exp(-s),$$

hence, for all b > 0:

$$\Pr\left[\frac{v'Qv}{NT} > b\right] \le \exp\left[-\left(\frac{bNT}{6\lambda} - \frac{4KTq}{3}\right)\right].$$

Lastly, by the union bound, given that the number of partitions  $\{k_{i1}\} \cap \{k_{i2}\}$  is bounded by  $K^{2N}$ :

$$\Pr\left[\frac{1}{NT}\sum_{i=1}^{N}\sum_{t=1}^{T}\|\overline{V}_{t}(\widehat{k}_{i},\underline{k}_{i})\|^{2} > b\right] \leq K^{2N}\max_{(\{k_{i1}\},\{k_{i2}\})} \Pr\left[\frac{1}{NT}\sum_{i=1}^{N}\sum_{t=1}^{T}\|\overline{V}_{t}(k_{i1},k_{i2})\|^{2} > b\right] \\ \leq \exp\left[2N\ln K + \frac{4KTq}{3} - \frac{bNT}{6\lambda}\right].$$

This implies that:

$$\frac{1}{NT} \sum_{i=1}^{N} \sum_{t=1}^{T} \|\overline{V}_t(\hat{k}_i, \underline{k}_i)\|^2 = O_p((\ln K)/T) + O_p(qK/N).$$

The rest of the proof is as in Theorem B1.

## C Model-based iteration and one-step GFE estimator

In this section we provide details on model-based iterated GFE, and we outline a continuously updated counterpart, the "one-step" GFE estimator.

## C.1 Model-based iteration

We are going to derive a bound on the within-group mean squared error of  $\alpha_{i0}$ , similarly to Corollary 1, for the case of iterated GFE.

Let  $(\hat{\theta}, \hat{\alpha})$  denote the two-step GFE estimator, with partition  $\{\hat{k}_i\}$ , and let  $\{\hat{k}_i^{(2)}\}$  denote the partition after one iteration. We use the same notation as in the proof of Theorem 1. In particular, we define  $\delta$  in the same way. We have:

$$\sum_{i=1}^{N}\sum_{j=1}^{p}\ell_{ij}\left(\widehat{\alpha}^{j}(\widehat{k}_{i},\widehat{\theta}),\widehat{\theta}\right) \leq \sum_{i=1}^{N}\sum_{j=1}^{p}\ell_{ij}\left(\widehat{\alpha}^{j}(\widehat{k}_{i}^{(2)},\widehat{\theta}),\widehat{\theta}\right).$$

Hence, expanding, we obtain:

$$-\frac{1}{2}\sum_{i=1}^{N}\sum_{j=1}^{p}\left(\widehat{\alpha}^{j}(\widehat{k}_{i}^{(2)},\widehat{\theta})-\widehat{\alpha}^{j}(\widehat{k}_{i},\widehat{\theta})\right)'v_{ij}^{\alpha}\left(a_{ij},\widehat{\theta}\right)\left(\widehat{\alpha}^{j}(\widehat{k}_{i}^{(2)},\widehat{\theta})-\widehat{\alpha}^{j}(\widehat{k}_{i},\widehat{\theta})\right)$$
$$\leq\sum_{i=1}^{N}\sum_{j=1}^{p}v_{ij}\left(\widehat{\alpha}^{j}(\widehat{k}_{i},\widehat{\theta}),\widehat{\theta}\right)\left(\widehat{\alpha}^{j}(\widehat{k}_{i}^{(2)},\widehat{\theta})-\widehat{\alpha}^{j}(\widehat{k}_{i},\widehat{\theta})\right),$$

where  $a_{ij}$  is between  $\widehat{\alpha}^{j}(\widehat{k}_{i},\widehat{\theta})$  and  $\widehat{\alpha}^{j}(\widehat{k}_{i}^{(2)},\widehat{\theta})$ .

Let us first assume that p does not grow with the sample size. Using similar arguments as in the proof of Theorem 1 we then have:

$$\frac{1}{Np}\sum_{i=1}^{N}\sum_{j=1}^{p}\left\|\widehat{\alpha}^{j}(\widehat{k}_{i}^{(2)},\widehat{\theta})-\widehat{\alpha}^{j}(\widehat{k}_{i},\widehat{\theta})\right\|^{2}=O_{p}\left(\frac{1}{Np}\sum_{i=1}^{N}\sum_{j=1}^{p}\left\|v_{ij}\left(\widehat{\alpha}^{j}(\widehat{k}_{i},\widehat{\theta}),\widehat{\theta}\right)\right\|^{2}\right).$$

Using Theorem 1, we then have:

$$\frac{1}{Np}\sum_{i=1}^{N}\sum_{j=1}^{p}\left\|\widehat{\alpha}^{j}(\widehat{k}_{i}^{(2)},\widehat{\theta})-\widehat{\alpha}^{j}(\widehat{k}_{i},\widehat{\theta})\right\|^{2}=O_{p}\left(\frac{1}{Np}\sum_{i=1}^{N}\sum_{j=1}^{p}\left\|v_{ij}\left(\alpha_{i0}^{j},\theta_{0}\right)\right\|^{2}\right)+O_{p}(\delta).$$

Moreover, since p does not grow with the sample size:

$$\frac{1}{Np}\sum_{i=1}^{N}\sum_{j=1}^{p}\left\|v_{ij}\left(\alpha_{i0}^{j},\theta_{0}\right)\right\|^{2}=O_{p}\left(\frac{1}{T}\right),$$

so:

$$\frac{1}{Np}\sum_{i=1}^{N}\sum_{j=1}^{p}\left\|\widehat{\alpha}^{j}(\widehat{k}_{i}^{(2)},\widehat{\theta})-\widehat{\alpha}^{j}(\widehat{k}_{i},\widehat{\theta})\right\|^{2}=O_{p}(\delta).$$

Using again Theorem 1, we obtain:

$$\frac{1}{Np}\sum_{i=1}^{N}\sum_{j=1}^{p}\left\|\widehat{\alpha}^{j}(\widehat{k}_{i}^{(2)},\widehat{\theta})-\alpha_{i0}^{j}\right\|^{2}=O_{p}(\delta).$$

It follows that:

$$\frac{1}{Np} \sum_{i=1}^{N} \sum_{j=1}^{p} \left\| \overline{\alpha}_{0}^{j}(\widehat{k}_{i}^{(2)}) - \alpha_{i0}^{j} \right\|^{2} = O_{p}(\delta).$$
(C1)

Let us now turn to the case where p grows with the sample size. In that case, one can bound:

$$\begin{aligned} \left\| \frac{1}{Np} \sum_{i=1}^{N} \sum_{j=1}^{p} v_{ij} \left( \widehat{\alpha}^{j}(\widehat{k}_{i},\widehat{\theta}),\widehat{\theta} \right) \left( \widehat{\alpha}^{j}(\widehat{k}_{i}^{(2)},\widehat{\theta}) - \widehat{\alpha}^{j}(\widehat{k}_{i},\widehat{\theta}) \right) \right\|^{2} \\ &\leq \left\| \frac{1}{Np} \sum_{i=1}^{N} \sum_{j=1}^{p} v_{ij} \left( \alpha_{i0}^{j},\theta_{0} \right) \left( \widehat{\alpha}^{j}(\widehat{k}_{i}^{(2)},\widehat{\theta}) - \widehat{\alpha}^{j}(\widehat{k}_{i},\widehat{\theta}) \right) \right\|^{2} + O_{p}(\delta) \frac{1}{Np} \sum_{i=1}^{N} \sum_{j=1}^{p} \left\| \widehat{\alpha}^{j}(\widehat{k}_{i}^{(2)},\widehat{\theta}) - \widehat{\alpha}^{j}(\widehat{k}_{i},\widehat{\theta}) \right\|^{2} \\ &= \left\| \frac{1}{Np} \sum_{i=1}^{N} \sum_{j=1}^{p} \overline{v}_{j} \left( \widehat{k}_{i},\widehat{k}_{i}^{(2)} \right) \left( \widehat{\alpha}^{j}(\widehat{k}_{i}^{(2)},\widehat{\theta}) - \widehat{\alpha}^{j}(\widehat{k}_{i},\widehat{\theta}) \right) \right\|^{2} + O_{p}(\delta) \frac{1}{Np} \sum_{i=1}^{N} \sum_{j=1}^{p} \left\| \widehat{\alpha}^{j}(\widehat{k}_{i}^{(2)},\widehat{\theta}) - \widehat{\alpha}^{j}(\widehat{k}_{i},\widehat{\theta}) \right\|^{2} \\ &\leq \left[ \frac{1}{Np} \sum_{i=1}^{N} \sum_{j=1}^{p} \left\| \overline{v}_{j} \left( \widehat{k}_{i},\widehat{k}_{i}^{(2)} \right) \right\|^{2} + O_{p}(\delta) \right] \frac{1}{Np} \sum_{i=1}^{N} \sum_{j=1}^{p} \left\| \widehat{\alpha}^{j}(\widehat{k}_{i}^{(2)},\widehat{\theta}) - \widehat{\alpha}^{j}(\widehat{k}_{i},\widehat{\theta}) \right\|^{2}, \end{aligned}$$

where  $\overline{v}_j(k,k')$  denotes the *j*-specific linear projection of  $v_{ij}\left(\alpha_{i0}^j,\theta_0\right)$  on group indicators  $\mathbf{1}\{\widehat{k}_i=k\}$ and  $\mathbf{1}\{\widehat{k}_i^{(2)}=k'\}$ .

Moreover, under a sub-Gaussianity assumption, and using a union bound argument together with Lemma B1, we will obtain the following bound:

$$\frac{1}{Np}\sum_{i=1}^{N}\sum_{j=1}^{p}\left\|\overline{v}_{j}\left(\widehat{k}_{i},\widehat{k}_{i}^{(2)}\right)\right\|^{2} = O_{p}\left(\frac{\ln K}{T}\right) + O_{p}\left(\frac{Kp}{NT}\right).$$

Combining results, and using Theorem 1, we obtain a bound that, similarly to (C1), depends on  $K^{-2/d}$ , where d is the underlying dimension of  $(\alpha'_{i0}, \mu'_{i0})$ . Hence we obtain a convergence rate, but no improvement, for iterated GFE relative to the baseline two-step approach.

### C.2 One-step GFE estimator

A continuously updated counterpart to the iterated estimator is the *one-step* GFE estimator, which is defined as follows:

$$\left(\widehat{\theta}^{\text{1step}}, \widehat{\alpha}^{\text{1step}}, \{\widehat{k}_i^{\text{1step}}\}\right) = \underset{(\theta, \alpha, \{k_i\})}{\operatorname{argmax}} \sum_{i=1}^N \ln f_i\left(\alpha\left(k_i\right), \theta\right),$$
(C2)

where the maximum is taken with respect to all possible parameter values  $(\theta, \alpha)$  and all possible partitions  $\{k_i\}$  of  $\{1, ..., N\}$  into at most K groups. This corresponds to the classification maximum likelihood estimator of Bryant and Williamson (1978); see also Hahn and Moon (2010) and Bonhomme and Manresa (2015). Unlike in two-step GFE, (C2) requires optimizing the likelihood function with respect to every partition and parameter value. This poses two difficulties. First, the estimator may be substantially more computationally intensive than two-step methods. Second, this complicates the statistical analysis since the discrete classification depends on parameter values and the objective function of the one-step estimator is therefore not smooth. In the case of the kmeans estimator, Pollard (1981, 1982) derived asymptotic properties for fixed K and T, as N tends to infinity. Deriving the properties of one-step estimators in (C2) as N, T, K tend jointly to infinity is an interesting avenue for future work.

## D Two-way grouped fixed-effects

In this section of the appendix we consider two-way GFE estimators, and we derive an expansion in the spirit of Theorem 1.

We have the following lemma, the proof of which is analogous to that of Lemma 1, and omitted for brevity. Here we consider the general case where the precision of  $h_i$  is S, and the precision of  $w_t$ is J, while in the main text we focus on the case where S = T and J = N.

**Lemma D1.** Suppose that there exist random vectors  $h_i$  and  $w_t$ , with fixed dimensions, and Lipschitzcontinuous functions  $\varphi$  and  $\phi$ , such that  $h_i = \varphi(\xi_{i0}) + o_p(1)$ ,  $\frac{1}{N} \sum_{i=1}^N ||h_i - \varphi(\xi_{i0})||^2 = O_p(1/S)$ ,  $w_t = \phi(\lambda_{t0}) + o_p(1)$ , and  $\frac{1}{T} \sum_{t=1}^T ||w_t - \phi(\lambda_{t0})||^2 = O_p(1/J)$  as N, T, S, J tend to infinity. Then we have, as N, S, K tend to infinity:

$$\frac{1}{N}\sum_{i=1}^{N}\left\|\widehat{h}(\widehat{k}_{i})-\varphi(\xi_{i0})\right\|^{2}=O_{p}\left(\frac{1}{S}\right)+O_{p}\left(B_{\xi}(K)\right)$$

and, as T, J, p tend to infinity:

$$\frac{1}{T}\sum_{t=1}^{T}\left\|\widehat{w}(\widehat{l}_{t})-\phi(\lambda_{t0})\right\|^{2}=O_{p}\left(\frac{1}{J}\right)+O_{p}\left(B_{\lambda}(p)\right),$$

where  $B_{\lambda}(p)$  is defined analogously as  $B_{\xi}(K)$ .

Assumption D1. (regularity, two-way)

(i)  $(Y'_{it}, X'_{it}, \xi'_{i0}, \lambda'_{t0})', i = 1, ..., N, t = 1, ..., T, are i.i.d.$ 

 $\ell_{it}(\alpha_{it},\theta)$  is three times differentiable in both its arguments, for all i, t.

The parameter space  $\Theta$  for  $\theta_0$  is compact, the spaces for  $\xi_{i0}$  and  $\lambda_{t0}$  are compact, and  $\theta_0$  belongs to the interior of  $\Theta$ .

(ii) N, T, S, J, K, p tend jointly to infinity.

 $\sup_{\xi,\lambda,\alpha,\theta} |\mathbb{E}_{\xi_{i0}=\xi,\lambda_{t0}=\lambda}(\ell_{it}(\alpha,\theta))| = O(1), \text{ and similarly for the first three derivatives of } \ell_{it} \text{ in both its arguments.}$ 

The minimum (respectively, maximum) eigenvalue of  $\left(-\frac{\partial^2 \ell_{it}(\alpha,\theta)}{\partial \alpha \partial \alpha'}\right)$  is bounded away from zero (resp., infinity) with probability one, uniformly in  $i, t, \alpha, \theta$ .

The third derivatives of  $\ell_{it}(\alpha, \theta)$  are  $O_p(1)$ , uniformly in  $i, t, \alpha, \theta$ .

 $\frac{1}{NT}\sum_{i=1}^{N}\sum_{t=1}^{T} [\ell_{it}(\alpha_{it0},\theta_0) - \mathbb{E}_{\xi_{i0},\lambda_{t0}}(\ell_{it}(\alpha_{it0},\theta_0))]^2 = O_p(1), \text{ and similarly for the first three derivatives of } \ell_{it} \text{ in both its arguments.}$ 

- (iii) For all  $\theta$ ,  $\xi$ , and  $\lambda$ , let  $\overline{\alpha}(\theta, \xi, \lambda) = \operatorname{argmax}_{\alpha} \mathbb{E}_{\xi_{i0} = \xi, \lambda_{t0} = \lambda}(\ell_{it}(\alpha, \theta))$ .  $\inf_{\xi,\lambda,\theta} \mathbb{E}_{\xi_{i0} = \xi, \lambda_{t0} = \lambda}(-\frac{\partial^2 \ell_{it}(\overline{\alpha}(\theta, \xi, \lambda), \theta)}{\partial \alpha \partial \alpha'})$  is positive definite.  $\mathbb{E}\left[\ell_{it}(\overline{\alpha}(\theta, \xi_{i0}, \lambda_{t0}), \theta)\right]$  has a unique maximum at  $\theta_0$  on  $\Theta$ , and its second derivative -H is negative definite.
- $$\begin{split} (iv) & \sup_{\tilde{\xi},\lambda,\alpha} \left\| \frac{\partial}{\partial \xi'} \right|_{\xi = \tilde{\xi}} \mathbb{E}_{\xi_{i0} = \xi,\lambda_{t0} = \lambda} (\operatorname{vec} \frac{\partial^2 \ell_{it}(\alpha,\theta_0)}{\partial \theta \partial \alpha'}) \right\| = O(1). \\ & \sup_{\xi,\tilde{\lambda},\alpha} \left\| \frac{\partial}{\partial \lambda'} \right|_{\lambda = \tilde{\lambda}} \mathbb{E}_{\xi_{i0} = \xi,\lambda_{t0} = \lambda} (\operatorname{vec} \frac{\partial^2 \ell_{it}(\alpha,\theta_0)}{\partial \theta \partial \alpha'}) \right\| = O(1). \\ & \sup_{\tilde{\xi},\lambda,\alpha} \left\| \frac{\partial}{\partial \xi'} \right|_{\xi = \tilde{\xi}} \mathbb{E}_{\xi_{i0} = \xi,\lambda_{t0} = \lambda} (\operatorname{vec} \frac{\partial^2 \ell_{it}(\alpha,\theta_0)}{\partial \alpha \partial \alpha'}) \right\| = O(1). \\ & \sup_{\xi,\tilde{\lambda},\alpha} \left\| \frac{\partial}{\partial \lambda'} \right|_{\lambda = \tilde{\lambda}} \mathbb{E}_{\xi_{i0} = \xi,\lambda_{t0} = \lambda} (\operatorname{vec} \frac{\partial^2 \ell_{it}(\alpha,\theta_0)}{\partial \alpha \partial \alpha'}) \right\| = O(1). \\ & \sup_{\tilde{\xi},\lambda,\theta} \left\| \frac{\partial}{\partial \xi'} \right|_{\xi = \tilde{\xi}} \mathbb{E}_{\xi_{i0} = \xi,\lambda_{t0} = \lambda} (\frac{\partial \ell_{it}(\overline{\alpha}(\theta,\xi,\lambda),\theta)}{\partial \alpha}) \right\| = O(1). \\ & \sup_{\xi,\tilde{\lambda},\theta} \left\| \frac{\partial}{\partial \lambda'} \right|_{\lambda = \tilde{\lambda}} \mathbb{E}_{\xi_{i0} = \xi,\lambda_{t0} = \lambda} (\frac{\partial \ell_{it}(\overline{\alpha}(\theta,\xi,\lambda),\theta)}{\partial \alpha}) \right\| = O(1). \end{split}$$
- (v)  $\mathbb{E}_{h_i=h,\xi_{i0}=\xi,w_t=w,\lambda_{t0}=\lambda}\left(\frac{\partial \ell_{it}(\overline{\alpha}(\theta,\xi,\lambda),\theta)}{\partial \alpha}\right)$  and  $\mathbb{E}_{h_i=h,\xi_{i0}=\xi,w_t=w,\lambda_{t0}=\lambda}\left(\operatorname{vec}\frac{\partial}{\partial \theta'}\Big|_{\theta_0}\frac{\partial \ell_{it}(\overline{\alpha}(\theta,\xi,\lambda),\theta)}{\partial \alpha}\right)$  are twice differentiable with respect to h and w, with first and second derivatives that are uniformly bounded in  $h \in \mathcal{H}, w \in \mathcal{W}, \xi, \lambda, and \theta \in \Theta$ , where  $\mathcal{H}$  and  $\mathcal{W}$  denote the supports of  $h_i$  and  $w_t$ , respectively.  $\|\operatorname{Var}_{h_i=h,\xi_{i0}=\xi,w_t=w,\lambda_{t0}=\lambda}\left(\frac{\partial \ell_{it}(\overline{\alpha}(\theta,\xi,\lambda),\theta)}{\partial \alpha}\right)\| = O(1), uniformly in h, w, \xi, \lambda, \theta.$  $\|\operatorname{Var}_{h_i=h,\xi_{i0}=\xi,w_t=w,\lambda_{t0}=\lambda}\left(\operatorname{vec}\frac{\partial}{\partial \theta'}\Big|_{\theta_0}\frac{\partial \ell_{it}(\overline{\alpha}(\theta,\xi,\lambda),\theta)}{\partial \alpha}\right)\| = O(1), uniformly in h, w, \xi, \lambda.$

Let us denote:

$$\widetilde{s}_{it} = \frac{\partial \ell_{it}(\alpha_{it0}, \theta_0)}{\partial \theta} + \mathbb{E}_{\xi_{i0}, \lambda_{t0}} \left( \frac{\partial^2 \ell_{it}(\alpha_{it0}, \theta_0)}{\partial \theta \partial \alpha'} \right) \left[ \mathbb{E}_{\xi_{i0}, \lambda_{t0}} \left( -\frac{\partial^2 \ell_{it}(\alpha_{it0}, \theta_0)}{\partial \alpha \partial \alpha'} \right) \right]^{-1} \frac{\partial \ell_{it}(\alpha_{it0}, \theta_0)}{\partial \alpha}, \quad (D1)$$

and:

$$\widetilde{H} = \mathbb{E}\left[\mathbb{E}_{\xi_{i0},\lambda_{t0}}\left(-\frac{\partial^{2}\ell_{it}(\alpha_{it0},\theta_{0})}{\partial\theta\partial\theta'}\right) - \mathbb{E}_{\xi_{i0},\lambda_{t0}}\left(\frac{\partial^{2}\ell_{it}(\alpha_{it0},\theta_{0})}{\partial\theta\partial\alpha'}\right)\left[\mathbb{E}_{\xi_{i0},\lambda_{t0}}\left(-\frac{\partial^{2}\ell_{it}(\alpha_{it0},\theta_{0})}{\partial\alpha\partial\alpha'}\right)\right]^{-1}\mathbb{E}_{\xi_{i0},\lambda_{t0}}\left(\frac{\partial^{2}\ell_{it}(\alpha_{it0},\theta_{0})}{\partial\alpha\partial\theta'}\right)\right]. \quad (D2)$$

**Theorem D1.** Let the conditions in Lemma D1 hold. Suppose that  $B_{\xi}(K) = O_p(K^{-\frac{2}{d}})$  and  $B_{\lambda}(p) = O_p(p^{-\frac{2}{d_{\lambda}}})$ . Suppose that  $\alpha$  and  $\mu$  are Lipschitz-continuous in both arguments. Suppose that there exist two Lipschitz-continuous functions  $\psi$  and  $\Psi$  such that  $\xi_{i0} = \psi(\varphi(\xi_{i0}))$  and  $\lambda_{i0} = \Psi(\phi(\lambda_{i0}))$ . Lastly, let Assumption D1 hold. Then, as N, T, S, J, K, p tend to infinity such that Kp/(NT) tends to zero, we have:

$$\widehat{\theta} = \theta_0 + \widetilde{H}^{-1} \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T \widetilde{s}_{it} + O_p\left(\frac{1}{S}\right) + O_p\left(\frac{1}{J}\right) + O_p\left(\frac{Kp}{NT}\right) + O_p\left(K^{-\frac{2}{d}}\right) + O_p\left(p^{-\frac{2}{d_\lambda}}\right) + o_p\left(\frac{1}{\sqrt{NT}}\right).$$
(D3)

#### Proof of Theorem D1

In the proof we closely follow the steps of the proof of Theorem 1. Let, for all  $\theta \in \Theta$ ,  $k \in \{1, ..., K\}$ , and  $l \in \{1, ..., p\}$ :

$$\widehat{\alpha}(k,l,\theta) = \operatorname{argmax}_{\alpha} \sum_{i=1}^{N} \sum_{t=1}^{T} \mathbf{1}\{\widehat{k}_{i} = k\} \mathbf{1}\{\widehat{l}_{t} = l\} \ell_{it}(\alpha,\theta), \qquad (D4)$$

and define  $\overline{\alpha}(\theta,\xi,\lambda)$  according to Assumption D1 (*iii*). Let also  $\delta = \frac{1}{S} + \frac{1}{J} + \frac{Kp}{NT} + K^{-\frac{2}{d}} + p^{-\frac{2}{d_{\lambda}}}$ .

A key step in the proof is to establish the following main intermediate results:

$$\frac{1}{NT}\sum_{i=1}^{N}\sum_{t=1}^{T}\frac{\partial\ell_{it}(\widehat{\alpha}(\widehat{k}_{i},\widehat{l}_{t},\theta_{0}),\theta_{0})}{\partial\theta} = \frac{1}{NT}\sum_{i=1}^{N}\sum_{t=1}^{T}\frac{\partial}{\partial\theta}\Big|_{\theta_{0}}\ell_{it}\left(\overline{\alpha}(\theta,\xi_{i0},\lambda_{t0}),\theta\right) + O_{p}\left(\delta\right), \quad (D5)$$

$$\frac{1}{NT}\sum_{i=1}^{N}\sum_{t=1}^{T}\frac{\partial^{2}}{\partial\theta\partial\theta'}\Big|_{\theta_{0}}\left(\ell_{it}\left(\widehat{\alpha}(\widehat{k}_{i},\widehat{l}_{t},\theta),\theta\right)-\ell_{it}\left(\overline{\alpha}(\theta,\xi_{i0},\lambda_{t0}),\theta\right)\right)=o_{p}(1).$$
(D6)

**Consistency of**  $\hat{\theta}$ . The consistency proof follows Theorem 1 closely. First, we show that  $\overline{\alpha}(\theta, \xi, \lambda)$  is Lipschitz-continuous, with coefficients that are uniformly bounded with respect to  $\theta$ ,  $\xi$  and  $\lambda$ .

Next, we define  $a(k, l, \theta) = \overline{\alpha}(\theta, \psi(\widehat{h}(k)), \Psi(\widehat{w}(l)))$ , and we use Lemma D1 to show that:

$$\sup_{\theta \in \Theta} \frac{1}{NT} \sum_{i=1}^{N} \sum_{t=1}^{T} \left\| a(\widehat{k}_i, \widehat{l}_t, \theta) - \overline{\alpha} \left(\theta, \xi_{i0}, \lambda_{t0}\right) \right\|^2 = O_p(\delta).$$
(D7)

Next, we use that, for all  $\theta$ :

$$\sum_{i=1}^{N} \sum_{t=1}^{T} \ell_{it} \left( a(\widehat{k}_i, \widehat{l}_t, \theta), \theta \right) \le \sum_{i=1}^{N} \sum_{t=1}^{T} \ell_{it} \left( \widehat{\alpha}(\widehat{k}_i, \widehat{l}_t, \theta), \theta \right),$$
(D8)

and we expand both sides of this inequality as in the proof of Theorem 1 to obtain, for all  $\theta \in \Theta$ :

$$\frac{1}{NT}\sum_{i=1}^{N}\sum_{t=1}^{T}\left\|\widehat{\alpha}(\widehat{k}_{i},\widehat{l}_{t},\theta)-\overline{\alpha}(\theta,\xi_{i0},\lambda_{t0})\right\|^{2} \leq O_{p}\left[\left(\frac{1}{NT}\sum_{i=1}^{N}\sum_{t=1}^{T}\|\overline{v}(\widehat{k}_{i},\widehat{l}_{t},\theta)\|^{2}\right)^{\frac{1}{2}}\left(\frac{1}{NT}\sum_{i=1}^{N}\sum_{t=1}^{T}\left\|\widehat{\alpha}(\widehat{k}_{i},\widehat{l}_{t},\theta)-a(\widehat{k}_{i},\widehat{l}_{t},\theta)\right\|^{2}\right)^{\frac{1}{2}}\right]+O_{p}(\delta),$$

where  $\overline{v}(k, l, \theta)$  denotes the mean of  $v_{it}(\overline{\alpha}(\theta, \xi_{i0}, \lambda_{t0}), \theta)$  in the intersection of groups  $\hat{k}_i = k$  and  $\hat{l}_t = l$ . As in the proof of Theorem 1, the key step is to show:

$$\frac{1}{NT}\sum_{i=1}^{N}\sum_{t=1}^{T}\left\|\overline{v}(\hat{k}_{i},\hat{l}_{t},\theta)\right\|^{2} = O_{p}\left(\delta\right).$$
(D9)

Let, for all  $\theta, h, w, \xi, \lambda$ :

$$\rho(h,\xi,w,\lambda,\theta) = \mathbb{E}_{h_i=h,\xi_{i0}=\xi,w_t=w,\lambda_{t0}=\lambda}(v_{it}(\overline{\alpha}(\theta,\xi,\lambda),\theta)),$$

and let, for all  $i, t, \theta$ :

$$\zeta_{it}(\theta) = v_{it}(\overline{\alpha}(\theta, \xi_{i0}, \lambda_{t0}), \theta) - \rho(h_i, \xi_{i0}, w_t, \lambda_{t0}, \theta).$$

Expanding  $\rho(h_i, \xi_{i0}, w_t, \lambda_{t0}, \theta)$  around  $h_i = \varphi(\xi_{i0})$  and  $w_t = \phi(\lambda_{t0})$ , we have:

$$\rho(h_i, \xi_{i0}, w_t, \lambda_{t0}, \theta) = \rho(\varphi(\xi_{i0}), \xi_{i0}, \phi(\lambda_{t0}), \lambda_{t0}, \theta) + \frac{\partial \rho(\varphi(\xi_{i0}), \xi_{i0}, \phi(\lambda_{t0}), \lambda_{t0}, \theta)}{\partial h'} (h_i - \varphi(\xi_{i0})) + \frac{\partial \rho(\varphi(\xi_{i0}), \xi_{i0}, \phi(\lambda_{t0}), \lambda_{t0}, \theta)}{\partial w'} (w_t - \phi(\lambda_{t0})) + O_p\left(\frac{1}{S}\right) + O_p\left(\frac{1}{J}\right).$$

Hence, taking expectations:

$$\begin{aligned} 0 &= \mathbb{E}_{\xi_{i0},\lambda_{t0}} \left( v_{it}(\overline{\alpha}(\theta,\xi_{i0},\lambda_{t0}),\theta)) \right) \\ &= \mathbb{E}_{\xi_{i0},\lambda_{t0}} \left[ \rho(h_i,\xi_{i0},w_t,\lambda_{t0},\theta) \right] \\ &= \rho(\varphi(\xi_{i0}),\xi_{i0},\phi(\lambda_{t0}),\lambda_{t0},\theta) + \frac{\partial \rho(\varphi(\xi_{i0}),\xi_{i0},\phi(\lambda_{t0}),\lambda_{t0},\theta)}{\partial h'} \mathbb{E}_{\xi_{i0},\lambda_{t0}}(h_i - \varphi(\xi_{i0})) \\ &+ \frac{\partial \rho(\varphi(\xi_{i0}),\xi_{i0},\phi(\lambda_{t0}),\lambda_{t0},\theta)}{\partial w'} \mathbb{E}_{\xi_{i0},\lambda_{t0}}(w_t - \phi(\lambda_{t0})) + O_p\left(\frac{1}{S}\right) + O_p\left(\frac{1}{J}\right). \end{aligned}$$

It follows that:

$$\frac{1}{NT} \sum_{i=1}^{N} \sum_{t=1}^{T} \|\rho(h_i, \xi_{i0}, w_t, \lambda_{t0}, \theta)\|^2 = O_p\left(\frac{1}{S}\right) + O_p\left(\frac{1}{J}\right).$$

We thus only need to bound:

$$\mathbb{E}\left[\frac{1}{NT}\sum_{i=1}^{N}\sum_{t=1}^{p}\|\overline{\zeta}(\hat{k}_{i},\hat{l}_{t},\theta)\|^{2}\right] = \frac{1}{NT}\sum_{k=1}^{K}\sum_{l=1}^{p}\mathbb{E}\left[\frac{\sum_{i=1}^{N}\sum_{t=1}^{T}\mathbf{1}\{\hat{k}_{i}=k\}\mathbf{1}\{\hat{l}_{t}=l\}\mathbb{E}_{h_{i},\xi_{i0},w_{t},\lambda_{t0}}\left(\zeta_{it}(\theta)'\zeta_{it}(\theta)\right)}{\sum_{i=1}^{N}\sum_{t=1}^{T}\mathbf{1}\{\hat{k}_{i}=k\}\mathbf{1}\{\hat{l}_{t}=l\}}\right],$$

where we have used that observations are independent across i and t. To bound this quantity, we use part (v) in Assumption D1. We thus obtain (D9).

The rest of the consistency part is as in the proof of Theorem 1.

**Proof of (D5).** To show (D5), we are now going to show that:

$$\frac{1}{NT}\sum_{i=1}^{N}\sum_{t=1}^{T}\left\{v_{it}^{\theta}\left(\widehat{\alpha}(\widehat{k}_{i},\widehat{l}_{t})-\alpha_{it0}\right)+\mathbb{E}_{\xi_{i0},\lambda_{t0}}\left(v_{it}^{\theta}\right)\left[\mathbb{E}_{\xi_{i0},\lambda_{t0}}\left(v_{it}^{\alpha}\right)\right]^{-1}v_{it}\right\}=O_{p}\left(\delta\right),\tag{D10}$$

where we omit references to  $\theta_0$  and  $\alpha_{it0}$ .

We will bound, in turn:

$$A \equiv \frac{1}{NT} \sum_{i=1}^{N} \sum_{t=1}^{T} \mathbb{E}_{\xi_{i0},\lambda_{t0}} \left( v_{it}^{\theta} \right) \left[ \mathbb{E}_{\xi_{i0},\lambda_{t0}} \left( v_{it}^{\alpha} \right) \right]^{-1} v_{it}^{\alpha} \left( \widehat{\alpha}(\widehat{k}_{i},\widehat{l}_{t}) - \alpha_{it0} + (v_{it}^{\alpha})^{-1} v_{it} \right),$$
  
$$B \equiv \frac{1}{NT} \sum_{i=1}^{N} \sum_{t=1}^{T} \left( v_{it}^{\theta} \left( v_{it}^{\alpha} \right)^{-1} - \mathbb{E}_{\xi_{i0},\lambda_{t0}} \left( v_{it}^{\theta} \right) \left[ \mathbb{E}_{\xi_{i0},\lambda_{t0}} \left( v_{it}^{\alpha} \right) \right]^{-1} \right) v_{it}^{\alpha} \left( \widehat{\alpha}(\widehat{k}_{i},\widehat{l}_{t}) - \alpha_{it0} \right).$$

To bound A, the key term is:

$$A_3 \equiv \frac{1}{NT} \sum_{i=1}^{N} \sum_{t=1}^{T} \mathbb{E}_{\xi_{i0},\lambda_{t0}} \left( v_{it}^{\theta} \right) \left[ \mathbb{E}_{\xi_{i0},\lambda_{t0}} \left( v_{it}^{\alpha} \right) \right]^{-1} \left( -v_{it}^{\alpha} \right) \left( \left( -v_{it}^{\alpha} \right)^{-1} v_{it} - \widetilde{v}(\widehat{k}_i, \widehat{l}_t) \right),$$

where  $\tilde{v}$  is defined analogously as in the proof of Theorem 1.

Let 
$$z(\xi, \lambda)' = \mathbb{E}_{\xi_{i0}=\xi, \lambda_{t0}=\lambda} \left( v_{it}^{\theta} \right) \left[ \mathbb{E}_{\xi_{i0}=\xi, \lambda_{t0}=\lambda} \left( v_{it}^{\alpha} \right) \right]^{-1}$$
. We have:  

$$A_{3} = \frac{1}{NT} \sum_{i=1}^{N} \sum_{t=1}^{T} \left( z(\xi_{i0}, \lambda_{t0})' - z^{*} \left( \widehat{k}_{i}, \widehat{l}_{t} \right)' \right) v_{it} + \frac{1}{NT} \sum_{i=1}^{N} \sum_{t=1}^{T} \left( z^{*} \left( \widehat{k}_{i}, \widehat{l}_{t} \right)' - \widetilde{z} \left( \widehat{k}_{i}, \widehat{l}_{t} \right)' \right) v_{it},$$

where  $\widetilde{z}(k,l)$  and  $z^{*}(k,l)$  are defined analogously as in the proof of Theorem 1.

To see that the first term in  $A_3$  is  $O_p(\delta)$ , we use an argument similar to the one we used in the proof of Theorem 1. Let:  $\zeta_{it} = v_{it} - \mathbb{E}_{h_i,\xi_{i0},w_t,\lambda_{t0}}(v_{it})$ . We have, as in the proof of Theorem 1:

$$\frac{1}{NT} \sum_{i=1}^{N} \sum_{t=1}^{T} \left\| \mathbb{E}_{h_i, \xi_{i0}, w_t, \lambda_{t0}}(v_{it}) \right\|^2 = O_p\left(\frac{1}{S}\right) + O_p\left(\frac{1}{J}\right).$$
(D11)

Moreover, we have:

$$\frac{1}{NT} \sum_{i=1}^{N} \sum_{t=1}^{T} \left\| z(\xi_{i0}, \lambda_{t0}) - z^*(\widehat{k}_i, \widehat{l}_t) \right\|^2 = O_p(\delta).$$
(D12)

Since  $\zeta_{it}$  are independent across *i* and *t*, with zero mean conditional on  $h_1, ..., h_N, g_1, ..., g_T, \xi_{10}, ..., \xi_{N0}$ , and  $\lambda_{10}, ..., \lambda_{T0}$ , we thus have:

$$\mathbb{E}\left[\left\|\frac{1}{NT}\sum_{i=1}^{N}\sum_{t=1}^{T}\left(z(\xi_{i0},\lambda_{t0})'-z^{*}\left(\hat{k}_{i},\hat{l}_{t}\right)'\right)v_{it}\right\|^{2}\right] \\
\leq 2\left(O\left(\frac{1}{S}\right)+O\left(\frac{1}{J}\right)\right)\mathbb{E}\left[\frac{1}{NT}\sum_{i=1}^{N}\sum_{t=1}^{T}\left\|z(\xi_{i0},\lambda_{t0})'-z^{*}\left(\hat{k}_{i},\hat{l}_{t}\right)'\right\|^{2}\right] \\
+2\mathbb{E}\left[\left\|\frac{1}{NT}\sum_{i=1}^{N}\sum_{t=1}^{T}\left(z(\xi_{i0},\lambda_{t0})'-z^{*}\left(\hat{k}_{i},\hat{l}_{t}\right)'\right)\zeta_{it}\right\|^{2}\right] \\
=O\left(\frac{\delta}{S}\right)+O\left(\frac{\delta}{J}\right) \\
+\frac{1}{N^{2}T^{2}}\sum_{i=1}^{N}\sum_{t=1}^{T}\mathbb{E}\left[\left(z(\xi_{i0},\lambda_{t0})'-z^{*}\left(\hat{k}_{i},\hat{l}_{t}\right)'\right)\mathbb{E}_{h_{i},\xi_{i0},w_{t},\lambda_{t0}}\left[\zeta_{it}\zeta_{it}'\right]\left(z(\xi_{i0},\lambda_{t0})'-z^{*}\left(\hat{k}_{i},\hat{l}_{t}\right)'\right)\right] \\
=O\left(\frac{\delta}{S}\right)+O\left(\frac{\delta}{J}\right)+O\left(\frac{\delta}{NT}\right)=O(\delta^{2}).$$

The second term in  $A_3$  is also  $O_p(\delta)$ , using similar arguments as in the proof of Theorem 1 and equation (D9). This shows that  $A = O_p(\delta)$ .

Let us now turn to *B*. Let  $\pi'_{it} = v^{\theta}_{it} (v^{\alpha}_{it})^{-1} - \mathbb{E}_{\xi_{i0},\lambda_{t0}} (v^{\theta}_{it}) [\mathbb{E}_{\xi_{i0},\lambda_{t0}} (v^{\alpha}_{it})]^{-1}$ . As in the proof of Theorem 1, the key step is to bound:

$$B_{3} \equiv \frac{1}{NT} \sum_{i=1}^{N} \sum_{t=1}^{T} \pi_{it}' v_{it}^{\alpha} \left( \widetilde{\alpha}(\widehat{k}_{i}, \widehat{l}_{t}) - \alpha_{it0} \right)$$
$$= \frac{1}{NT} \sum_{i=1}^{N} \sum_{t=1}^{T} \pi_{it}' v_{it}^{\alpha} \left( \alpha^{*}(\widehat{k}_{i}, \widehat{l}_{t}) - \alpha_{it0} \right) + \frac{1}{NT} \sum_{i=1}^{N} \sum_{t=1}^{T} \pi_{it}' v_{it}^{\alpha} \left( \widetilde{\alpha}(\widehat{k}_{i}, \widehat{l}_{t}) - \alpha^{*}(\widehat{k}_{i}, \widehat{l}_{t}) \right),$$

where  $\tilde{\alpha}(k,l)$  and  $\alpha^*(k,l)$  are defined analogously as in the proof of Theorem 1.

The first term is  $O_p(\delta)$ , since the  $\tau_{it} = \pi'_{it} v_{it}^{\alpha}$  are independent across *i* and *t* with zero mean given  $\xi_{i0}, \lambda_{t0}$ . The second term is:

$$\frac{1}{NT}\sum_{i=1}^{N}\sum_{t=1}^{T}\pi_{it}'v_{it}^{\alpha}\left(\widetilde{\alpha}(\widehat{k}_{i},\widehat{l}_{t})-\alpha^{*}(\widehat{k}_{i},\widehat{l}_{t})\right) = \frac{1}{NT}\sum_{i=1}^{N}\sum_{t=1}^{T}\widetilde{\pi}(\widehat{k}_{i},\widehat{l}_{t})'v_{it}^{\alpha}\left(\widetilde{\alpha}(\widehat{k}_{i},\widehat{l}_{t})-\alpha^{*}(\widehat{k}_{i},\widehat{l}_{t})\right).$$

As in the proof of Theorem 1,  $\frac{1}{NT}\sum_{i=1}^{N}\sum_{t=1}^{T} \|\widetilde{\pi}(\widehat{k}_i,\widehat{l}_t)\|^2 = O_p(\delta)$ . Moreover, we have:

$$\frac{1}{NT}\sum_{i=1}^{N}\sum_{t=1}^{T}\|\widetilde{\alpha}(\widehat{k}_{i},\widehat{l}_{t})-\alpha^{*}(\widehat{k}_{i},\widehat{l}_{t})\|^{2}=O_{p}(\delta)$$

This shows that  $B_3 = O_p(\delta)$ , hence that  $B = O_p(\delta)$ . This implies that (D5) holds.

**Proof of (D6).** The proof of (D6) follows the same steps as the corresponding part in the proof of Theorem 1. The key step is to show that:

$$\frac{1}{NT}\sum_{i=1}^{N}\sum_{t=1}^{T}\left\|\frac{\partial\widehat{\alpha}(\widehat{k}_{i},\widehat{l}_{t},\theta_{0})}{\partial\theta'} - \frac{\partial\overline{\alpha}^{j}(\theta_{0},\xi_{i0},\lambda_{t0})}{\partial\theta'}\right\|^{2} = o_{p}\left(1\right).$$
(D13)

The proof of (D13) follows similar arguments as in the proof of Theorem 1.

## **E** Further results for fixed p

In this section we derive two additional properties of GFE estimators, in the case where p is fixed; that is, when heterogeneity is time-invariant and its dimension does not grow with the sample size. Without loss of generality we set p = 1. For simplicity, in this section we assume that S = T, and we remove p and j from the notation. Finally, we denote  $\eta_{i0} = (\alpha'_{i0}, \mu'_{i0})'$  the full vector of heterogeneity.

## E.1 Expansion of two-step GFE

We start with a result that shows that the expansion in Theorem 1 continues to hold under suitable conditions, when the concavity in  $\alpha$  of Assumption 3 part (iia) fails to hold.

**Corollary E1.** Let the conditions in Theorem 1 hold, with  $\xi_{i0} = \eta_{i0}$ , except that Assumption 3 part *(iia)* is replaced by the following:

- (i) For all  $\epsilon > 0$ ,  $\inf_{\theta,\eta=(\alpha',\mu')'} \inf_{\|(\widetilde{\alpha},\widetilde{\theta})-(\alpha,\theta)\|>\epsilon} \mathbb{E}_{\eta_{i0}=\eta}[\ell_i(\alpha,\theta)-\ell_i(\widetilde{\alpha},\widetilde{\theta})] > 0$ .
- (ii) The function  $\hat{\ell}_i(\theta) = \ell_i(\hat{\alpha}(\hat{k}_i, \theta), \theta)$  is three times differentiable on a neighborhood of  $\theta_0$  (that is, the function is three times differentiable for almost all sample realizations), where  $\hat{\alpha}(k, \theta)$  for all k is the solution of (3) given any  $\theta \in \Theta$ . Moreover,  $\frac{1}{N} \sum_{i=1}^{N} ||\frac{\partial^2 \hat{\ell}_i(\theta)}{\partial \theta \partial \theta'}||^2 = O_p(1)$  uniformly in a neighborhood of  $\theta_0$ , and similarly for the third derivative of  $\hat{\ell}_i$ .

Then, for p fixed as N, T, K tend to infinity, we have:

$$\widehat{\theta} = \theta_0 + H^{-1} \frac{1}{N} \sum_{i=1}^N s_i + O_p\left(\frac{1}{T}\right) + O_p\left(K^{-\frac{2}{d}}\right) + o_p\left(\frac{1}{\sqrt{NT}}\right), \quad (E1)$$

and:

$$\frac{1}{N}\sum_{i=1}^{N} \left\|\widehat{\alpha}(\widehat{k}_{i}) - \alpha_{i0}\right\|^{2} = O_{p}\left(\frac{1}{T}\right) + O_{p}\left(K^{-\frac{2}{d}}\right).$$
(E2)

Condition (i) in Corollary E1 is an identification condition. Condition (ii) is a differentiability condition on the sample GFE objective function. Such an assumption is not needed in order to characterize the first-order properties of fixed-effects estimators, since, under suitable conditions, fixed-effects estimators of individual effects are uniformly consistent in the sense that  $\max_{i=1,...,N} \|\hat{\alpha}_i - \hat{\alpha}_i\|$ 

 $\alpha_{i0} \parallel = o_p(1)$ ; see, e.g., Hahn and Kuersteiner (2011). In contrast, our characterization of GFE estimators is based on establishing a rate of convergence for the average  $\frac{1}{N} \sum_{i=1}^{N} \|\widehat{\alpha}(\widehat{k}_i) - \alpha_{i0}\|^2$ . Lastly, the score  $s_i$  and Hessian H in Corollary E1 are the same as in Theorem 1.

**Proof of Corollary E1.** The proof follows closely the proof of Theorem 1. The beginning of the consistency proof is identical. Here,  $\delta = 1/T + K^{-2/d}$ . One difference is that we introduce the fixed-effects estimator of  $\alpha_i$ , for given  $\theta$ , which is defined as  $\hat{\alpha}_i(\theta) = \operatorname{argmax}_{\alpha} \ell_i(\alpha, \theta)$ . We have, for all  $\theta$ :

$$\sum_{i=1}^{N} \ell_i \left( a(\widehat{k}_i, \theta), \theta \right) \le \sum_{i=1}^{N} \ell_i \left( \widehat{\alpha}(\widehat{k}_i, \theta), \theta \right) \le \sum_{i=1}^{N} \ell_i \left( \widehat{\alpha}_i(\theta), \theta \right).$$
(E3)

Moreover, expanding the fixed-effects log-likelihood around  $\overline{\alpha}_i(\theta) \equiv \overline{\alpha}(\theta, \eta_{i0})$ , we have (as in Arellano and Hahn, 2007, for example):

$$\frac{1}{N}\sum_{i=1}^{N}\ell_{i}\left(\widehat{\alpha}_{i}(\theta),\theta\right) = \frac{1}{N}\sum_{i=1}^{N}\ell_{i}\left(\overline{\alpha}_{i}\left(\theta\right),\theta\right) + O_{p}\left(\frac{1}{T}\right),$$

uniformly in  $\theta$ .

Now, for some  $a_i(\theta)$  between  $\widehat{\alpha}_i(\theta)$  and  $a(\widehat{k}_i, \theta)$ , we have:

$$\frac{1}{N}\sum_{i=1}^{N}\ell_{i}\left(a(\widehat{k}_{i},\theta),\theta\right) - \frac{1}{N}\sum_{i=1}^{N}\ell_{i}\left(\widehat{\alpha}_{i}\left(\theta\right),\theta\right) = \frac{1}{2N}\sum_{i=1}^{N}\left(a(\widehat{k}_{i},\theta) - \widehat{\alpha}_{i}\left(\theta\right)\right)'v_{i}^{\alpha}\left(a_{i}\left(\theta\right),\theta\right)\left(a(\widehat{k}_{i},\theta) - \widehat{\alpha}_{i}\left(\theta\right)\right).$$

We then obtain, using similar arguments as in the proof of Theorem 1, that the following equality holds pointwise in  $\theta$ :

$$\left|\frac{1}{N}\sum_{i=1}^{N}\ell_{i}\left(\widehat{\alpha}(\widehat{k}_{i},\theta),\theta\right)-\frac{1}{N}\sum_{i=1}^{N}\ell_{i}\left(\widehat{\alpha}_{i}\left(\theta\right),\theta\right)\right|=O_{p}(\delta).$$
(E4)

In turn, from the uniform convergence of the fixed-effects log-likelihood:

$$\sup_{\theta \in \Theta} \left| \frac{1}{N} \sum_{i=1}^{N} \ell_i \left( \widehat{\alpha}_i(\theta), \theta \right) - \frac{1}{N} \sum_{i=1}^{N} \ell_i \left( \overline{\alpha}_i \left( \theta \right), \theta \right) \right| = o_p(1),$$

we obtain, using similar arguments as above:

$$\sup_{\theta \in \Theta} \left| \frac{1}{N} \sum_{i=1}^{N} \ell_i \left( \widehat{\alpha}(\widehat{k}_i, \theta), \theta \right) - \frac{1}{N} \sum_{i=1}^{N} \ell_i \left( \widehat{\alpha}_i \left( \theta \right), \theta \right) \right| \le \sup_{\theta \in \Theta} \left| \frac{1}{N} \sum_{i=1}^{N} \ell_i \left( a(\widehat{k}_i, \theta), \theta \right) - \frac{1}{N} \sum_{i=1}^{N} \ell_i \left( \widehat{\alpha}_i \left( \theta \right), \theta \right) \right| \le \sup_{\theta \in \Theta} \left| \frac{1}{N} \sum_{i=1}^{N} \ell_i \left( a(\widehat{k}_i, \theta), \theta \right) - \frac{1}{N} \sum_{i=1}^{N} \ell_i \left( \overline{\alpha}_i \left( \theta \right), \theta \right) \right| + o_p(1) = o_p(1).$$

Consistency of  $\hat{\theta}$  then follows, as in the proof of Theorem 1.

Let us now show equation (A3). From (E4) evaluated at  $\theta_0$  we have, for some  $a_i$  between  $\widehat{\alpha}(\widehat{k}_i, \theta_0)$ and  $\widehat{\alpha}_i(\theta_0)$ , and omitting from now on the reference to  $\theta_0$  for conciseness:

$$O_p(\delta) = \frac{1}{N} \sum_{i=1}^N \ell_i(\widehat{\alpha}_i) - \frac{1}{N} \sum_{i=1}^N \ell_i(\widehat{\alpha}(\widehat{k}_i)) = \frac{1}{2N} \sum_{i=1}^N \left(\widehat{\alpha}(\widehat{k}_i) - \widehat{\alpha}_i\right)' \left(-v_i^\alpha\left(a_i\right)\right) \left(\widehat{\alpha}(\widehat{k}_i) - \widehat{\alpha}_i\right) \ge 0.$$
(E5)

Now, by the assumptions of Theorem 1 part (a), there exists a constant  $\epsilon > 0$  and a positive definite matrix  $\underline{\Sigma}$  such that:

$$\inf_{\eta=(\alpha',\mu')'} \inf_{\|\widetilde{\alpha}-\alpha\| \le \epsilon} \mathbb{E}_{\eta_{i0}=\eta} \left( -v_i^{\alpha}(\widetilde{\alpha}) \right) \ge \underline{\Sigma}.$$

For this  $\epsilon$  we will first show that:

$$\frac{1}{N}\sum_{i=1}^{N} \mathbf{1}\left\{\|\widehat{\alpha}(\widehat{k}_{i}) - \alpha_{i0}\| > \epsilon\right\} = O_{p}(\delta).$$
(E6)

Showing (E6) will allow us to control the difference  $\hat{\alpha}(\hat{k}_i) - \alpha_{i0}$  in an average sense.

To see that (E6) holds, let  $\iota_i = \mathbf{1} \left\{ \|\widehat{\alpha}(\widehat{k}_i) - \alpha_{i0}\| \le \epsilon \right\}$ , and note that by (E5) we have, since  $\ell_i(\widehat{\alpha}_i) \ge \ell_i(\widehat{\alpha}(\widehat{k}_i))$  for all *i*:

$$0 \leq \frac{1}{N} \sum_{i=1}^{N} (1 - \iota_i) \left( \ell_i \left( \widehat{\alpha}_i \right) - \ell_i \left( \widehat{\alpha}(\widehat{k}_i) \right) \right) = O_p(\delta).$$

Now, by part (i) in the assumptions of the corollary, and using that  $\max_{i=1,\dots,N} \|\hat{\alpha}_i - \alpha_{i0}\| = o_p(1)$ and Assumption 3 part (iia), we have:

$$\min_{i=1,\dots,N} \inf_{\|\alpha_i - \alpha_{i0}\| > \epsilon} \ell_i(\widehat{\alpha}_i) - \ell_i(\alpha_i) \ge \inf_{\eta = (\alpha', \mu')'} \inf_{\|\widetilde{\alpha} - \alpha\| > \epsilon} \mathbb{E}_{\eta_{i0} = \eta} \left[\ell_i(\alpha) - \ell_i(\widetilde{\alpha})\right] + o_p(1) \ge \zeta + o_p(1),$$

where  $\zeta > 0$  is a constant. Hence  $\frac{1}{N} \sum_{i=1}^{N} (1 - \iota_i)(\zeta + o_p(1)) = O_p(\delta)$ , from which (E6) follows.

Next, by Assumption 3 part (iia),  $\sup_{\eta = (\alpha', \mu')'} \sup_{\widetilde{\alpha}} \left\| v_i^{\alpha}(\widetilde{\alpha}) - \mathbb{E}_{\eta_{i0} = \eta} \left( v_i^{\alpha}(\widetilde{\alpha}) \right) \right\| = o_p(1)$ . We thus have:

$$\min_{i=1,\dots,N} \inf_{\|\alpha_i - \alpha_{i0}\| \le \epsilon} (-v_i^{\alpha}(\alpha_i)) \ge \underline{\Sigma} + o_p(1).$$
(E7)

Using (E5) this implies that:  $\frac{1}{N} \sum_{i=1}^{N} \iota_i \left\| \widehat{\alpha}(\widehat{k}_i) - \widehat{\alpha}_i \right\|^2 = O_p(\delta)$ . Hence, using in addition (E6) and the fact that the parameter space for  $\alpha_i$  is compact, we have:

$$\frac{1}{N}\sum_{i=1}^{N}\left\|\widehat{\alpha}(\widehat{k}_{i})-\widehat{\alpha}_{i}\right\|^{2} = \frac{1}{N}\sum_{i=1}^{N}\iota_{i}\left\|\widehat{\alpha}(\widehat{k}_{i})-\widehat{\alpha}_{i}\right\|^{2} + \frac{1}{N}\sum_{i=1}^{N}(1-\iota_{i})\left\|\widehat{\alpha}(\widehat{k}_{i})-\widehat{\alpha}_{i}\right\|^{2} = O_{p}(\delta).$$
(E8)

Given (E8), the rest of the proof of (A3) is as in the proof of Theorem 1, with some simplifications due to the fact that here p does not grow with the sample size.

However, the proof of (A4) is different from the one in Theorem 1. To proceed, let:

$$\widetilde{\iota}(k) = \mathbf{1} \left\{ \sum_{i'=1}^{N} \mathbf{1} \{ \widehat{k}_{i'} = k \} \left( -v_{i'}^{\alpha} \left( \widehat{\alpha}(\widehat{k}_{i'}) \right) \right) \ge \frac{1}{2} \underline{\Sigma} \right\}.$$

We are first going to show that:

$$\frac{1}{N}\sum_{i=1}^{N}\left(1-\widetilde{\iota}\left(\widehat{k}_{i}\right)\right)=O_{p}(\delta).$$
(E9)

Let  $\epsilon > 0$  as in (E6), and define  $\iota_i$  accordingly. From (E6) it suffices to show that:

$$\frac{1}{N}\sum_{i=1}^{N}\iota_{i}\left(1-\widetilde{\iota}\left(\widehat{k}_{i}\right)\right)=O_{p}(\delta).$$

With probability approaching one we have:  $\min_{i:\iota_i=1} \left(-v_i^{\alpha}\left(\widehat{\alpha}(\widehat{k}_i)\right)\right) \geq \frac{2}{3}\underline{\Sigma}$ . When this condition is satisfied we have:

$$\begin{split} \iota_{i}\left(1-\widetilde{\iota}\left(\widehat{k}_{i}\right)\right) &= \iota_{i}\left(1-1\left\{\sum_{i'=1}^{N}\mathbf{1}\left\{\widehat{k}_{i'}=\widehat{k}_{i}\right\}\left(-v_{i'}^{\alpha}\left(\widehat{\alpha}(\widehat{k}_{i'})\right)\right)-\frac{1}{2}\underline{\Sigma}\geq0\right\}\right)\\ &\leq \iota_{i}\left(1-1\left\{\left(-v_{i}^{\alpha}\left(\widehat{\alpha}(\widehat{k}_{i})\right)\right)-\frac{1}{2}\underline{\Sigma}+\sum_{i'\neq i}\mathbf{1}\left\{\widehat{k}_{i'}=\widehat{k}_{i}\right\}\left(-v_{i'}^{\alpha}\left(\widehat{\alpha}(\widehat{k}_{i'})\right)\right)\geq0\right\}\right)\\ &\leq \iota_{i}\left(1-1\left\{\sum_{i'=1}^{N}\iota_{i'}\mathbf{1}\left\{\widehat{k}_{i'}=\widehat{k}_{i}\right\}\frac{1}{6}\underline{\Sigma}\geq-\sum_{i'=1}^{N}(1-\iota_{i'})\mathbf{1}\left\{\widehat{k}_{i'}=\widehat{k}_{i}\right\}\left(-v_{i'}^{\alpha}\left(\widehat{\alpha}(\widehat{k}_{i'})\right)\right)\right\}\right)\\ &\leq \mathbf{1}\left\{\sum_{i'=1}^{N}\iota_{i'}\mathbf{1}\left\{\widehat{k}_{i'}=\widehat{k}_{i}\right\}\leq\frac{6}{\underline{\sigma}}\left(\max_{i'=1,\dots,N}\left\|-v_{i'}^{\alpha}\left(\widehat{\alpha}(\widehat{k}_{i'})\right)\right\|\right)\sum_{i'=1}^{N}(1-\iota_{i'})\mathbf{1}\left\{\widehat{k}_{i'}=\widehat{k}_{i}\right\}\right\}\\ &\leq \mathbf{1}\left\{\sum_{i'=1}^{N}\mathbf{1}\left\{\widehat{k}_{i'}=\widehat{k}_{i}\right\}\leq\left(1+\frac{6}{\underline{\sigma}}\max_{i'=1,\dots,N}\left\|-v_{i'}^{\alpha}\left(\widehat{\alpha}(\widehat{k}_{i'})\right)\right\|\right)\sum_{i'=1}^{N}(1-\iota_{i'})\mathbf{1}\left\{\widehat{k}_{i'}=\widehat{k}_{i}\right\}\right\},\end{split}$$

where  $\underline{\sigma}$  denotes the minimum eigenvalue of  $\underline{\Sigma}$ . Hence we have, with probably approaching one:

$$0 \leq \frac{1}{N} \sum_{i=1}^{N} \iota_{i} \left( 1 - \widetilde{\iota} \left( \widehat{k}_{i} \right) \right)$$

$$\leq \frac{1}{N} \sum_{i=1}^{N} \mathbf{1} \left\{ \sum_{i'=1}^{N} \mathbf{1} \{ \widehat{k}_{i'} = \widehat{k}_{i} \} \leq \left( 1 + \frac{6}{\underline{\sigma}} \max_{i'=1,\dots,N} \left\| -v_{i'}^{\alpha} \left( \widehat{\alpha}(\widehat{k}_{i'}) \right) \right\| \right) \sum_{i'=1}^{N} (1 - \iota_{i'}) \mathbf{1} \{ \widehat{k}_{i'} = \widehat{k}_{i} \} \right\}$$

$$= \frac{1}{N} \sum_{k=1}^{K} \sum_{i=1}^{N} \mathbf{1} \{ \widehat{k}_{i} = k \} \mathbf{1} \left\{ \sum_{i'=1}^{N} \mathbf{1} \{ \widehat{k}_{i'} = k \} \leq \left( 1 + \frac{6}{\underline{\sigma}} \max_{i'=1,\dots,N} \left\| -v_{i'}^{\alpha} \left( \widehat{\alpha}(\widehat{k}_{i'}) \right) \right\| \right) \sum_{i'=1}^{N} (1 - \iota_{i'}) \mathbf{1} \{ \widehat{k}_{i'} = k \} \right\}$$

$$\leq \frac{1}{N} \sum_{k=1}^{K} \left( 1 + \frac{6}{\underline{\sigma}} \max_{i'=1,\dots,N} \left\| -v_{i'}^{\alpha} \left( \widehat{\alpha}(\widehat{k}_{i'}) \right) \right\| \right) \sum_{i'=1}^{N} (1 - \iota_{i'}) \mathbf{1} \{ \widehat{k}_{i'} = k \} = O_{p} \left( \frac{1}{N} \sum_{i'=1}^{N} (1 - \iota_{i'}) \right) = O_{p}(\delta).$$

This shows (E9).

We are now going to show (A4). By part (ii) in the assumptions of the corollary, the Cauchy

Schwarz inequality, and (E9), we have:

$$\left\|\frac{1}{N}\sum_{i=1}^{N}\left(1-\widetilde{\iota}(\widehat{k}_{i})\right)\frac{\partial^{2}}{\partial\theta\partial\theta'}\Big|_{\theta_{0}}\ell_{i}\left(\widehat{\alpha}(\widehat{k}_{i},\theta),\theta\right)\right\|^{2}$$

$$\leq \frac{1}{N}\sum_{i=1}^{N}(1-\widetilde{\iota}(\widehat{k}_{i}))\times\frac{1}{N}\sum_{i=1}^{N}\left\|\frac{\partial^{2}}{\partial\theta\partial\theta'}\Big|_{\theta_{0}}\ell_{i}\left(\widehat{\alpha}(\widehat{k}_{i},\theta),\theta\right)\right\|^{2} = o_{p}(1).$$

Let k such that  $\tilde{\iota}(k) = 1$ . Differentiating with respect to  $\theta$ :  $\sum_{i=1}^{N} \mathbf{1}\{\hat{k}_i = k\}v_i(\hat{\alpha}(k,\theta),\theta) = 0$  we obtain, at  $\theta = \theta_0$ :

$$\frac{\partial \widehat{\alpha}(k)}{\partial \theta'} = \left(\sum_{i'=1}^{N} \mathbf{1}\{\widehat{k}_{i'} = k\} \left(-v_{i'}^{\alpha}\left(\widehat{\alpha}(\widehat{k}_{i'})\right)\right)\right)^{-1} \sum_{i'=1}^{N} \mathbf{1}\{\widehat{k}_{i'} = k\} \left(v_{i'}^{\theta}\left(\widehat{\alpha}(\widehat{k}_{i'})\right)\right)', \quad (E10)$$

where we note that, since  $\tilde{\iota}(k) = 1$ ,  $\sum_{i'=1}^{N} \mathbf{1}\{\hat{k}_{i'} = k\}\left(-v_{i'}^{\alpha}\left(\hat{\alpha}(\hat{k}_{i'})\right)\right)$  is bounded from below by  $\underline{\Sigma}/2$ . Let now:

$$D \equiv \frac{1}{N} \sum_{i=1}^{N} \widetilde{\iota}(\widehat{k}_{i}) \frac{\partial^{2}}{\partial \theta \partial \theta'} \Big|_{\theta_{0}} \ell_{i} \left( \widehat{\alpha}(\widehat{k}_{i},\theta), \theta \right) - \frac{1}{N} \sum_{i=1}^{N} \frac{\partial^{2}}{\partial \theta \partial \theta'} \Big|_{\theta_{0}} \ell_{i} \left( \overline{\alpha}_{i}(\theta), \theta \right)$$

We have, at  $\theta_0$  (omitting again the reference to  $\theta_0$  from the notation):

$$D = \frac{1}{N} \sum_{i=1}^{N} \left\{ \widetilde{\iota}(\widehat{k}_{i}) \frac{\partial^{2} \ell_{i}\left(\widehat{\alpha}(\widehat{k}_{i})\right)}{\partial \theta \partial \theta'} + \widetilde{\iota}(\widehat{k}_{i}) v_{i}^{\theta}\left(\widehat{\alpha}(\widehat{k}_{i})\right) \frac{\partial \widehat{\alpha}(\widehat{k}_{i})}{\partial \theta'} - \frac{\partial^{2} \ell_{i}\left(\alpha_{i0}, \theta_{0}\right)}{\partial \theta \partial \theta'} - v_{i}^{\theta} \frac{\partial \overline{\alpha}_{i}}{\partial \theta'} - \left(\frac{\partial \overline{\alpha}_{i}}{\partial \theta'}\right)' (v_{i}^{\theta})' - \left(\frac{\partial \overline{\alpha}_{i}}{\partial \theta'}\right)' v_{i}^{\alpha} \frac{\partial \overline{\alpha}_{i}}{\partial \theta'} - \frac{\partial^{2}}{\partial \theta \partial \theta'} \Big|_{\theta_{0}} \left(\overline{\alpha}_{i}(\theta)' v_{i}\right) \right\}$$
$$= \frac{1}{N} \sum_{i=1}^{N} \widetilde{\iota}(\widehat{k}_{i}) \frac{\partial^{2} \ell_{i}\left(\widehat{\alpha}(\widehat{k}_{i})\right)}{\partial \theta \partial \theta'} + \widetilde{\iota}(\widehat{k}_{i}) v_{i}^{\theta}\left(\widehat{\alpha}(\widehat{k}_{i})\right) \frac{\partial \widehat{\alpha}(\widehat{k}_{i})}{\partial \theta'} - \frac{\partial^{2} \ell_{i}\left(\alpha_{i0}, \theta_{0}\right)}{\partial \theta \partial \theta'} - v_{i}^{\theta} \frac{\partial \overline{\alpha}_{i}}{\partial \theta'} + o_{p}(1),$$

where we have used that  $\mathbb{E}_{\eta_{i0}}(v_i) = 0$ , and that  $\frac{\partial \overline{\alpha}_i}{\partial \theta'} = \left[\mathbb{E}_{\eta_{i0}}(-v_i^{\alpha})\right]^{-1} \mathbb{E}_{\eta_{i0}}(v_i^{\theta})'$ . Hence, using (E8) and (E9):

$$D = \frac{1}{N} \sum_{i=1}^{N} \widetilde{\iota}(\widehat{k}_{i}) v_{i}^{\theta} \left( \frac{\partial \widehat{\alpha}(\widehat{k}_{i}, \theta_{0})}{\partial \theta'} - \frac{\partial \overline{\alpha}_{i}(\theta_{0})}{\partial \theta'} \right) + o_{p}(1),$$

so:

$$D = \frac{1}{N} \sum_{i=1}^{N} \widetilde{\iota}(\widehat{k}_{i}) \mathbb{E}_{\eta_{i0}}\left(v_{i}^{\theta}\right) \left(\frac{\partial \widehat{\alpha}(\widehat{k}_{i}, \theta_{0})}{\partial \theta'} - \frac{\partial \overline{\alpha}_{i}(\theta_{0})}{\partial \theta'}\right) + o_{p}\left(1\right).$$

Next, defining  $z(\eta)' = \mathbb{E}_{\eta_{i0}=\eta} \left( v_i^{\theta} \right) \left[ \mathbb{E}_{\eta_{i0}=\eta} \left( -v_i^{\alpha} \right) \right]^{-1}$  and  $\widetilde{z}(k)$  as in the proof of Theorem 1 we have:

$$D = \frac{1}{N} \sum_{i=1}^{N} \widetilde{\iota}(\widehat{k}_{i}) z(\eta_{i0})' \mathbb{E}_{\eta_{i0}} (-v_{i}^{\alpha}) \left( \frac{\partial \widehat{\alpha}(\widehat{k}_{i})}{\partial \theta'} - \frac{\partial \overline{\alpha}_{i}}{\partial \theta'} \right) + o_{p} (1)$$

$$= \frac{1}{N} \sum_{i=1}^{N} \widetilde{\iota}(\widehat{k}_{i}) z(\eta_{i0})' (-v_{i}^{\alpha}) \left( \frac{\partial \widehat{\alpha}(\widehat{k}_{i})}{\partial \theta'} - \frac{\partial \overline{\alpha}_{i}}{\partial \theta'} \right) + o_{p} (1)$$

$$= \frac{1}{N} \sum_{i=1}^{N} \widetilde{\iota}(\widehat{k}_{i}) z(\eta_{i0})' \left( -v_{i}^{\alpha} \left( \widehat{\alpha}(\widehat{k}_{i}) \right) \right) \left( \frac{\partial \widehat{\alpha}(\widehat{k}_{i})}{\partial \theta'} - \frac{\partial \overline{\alpha}_{i}}{\partial \theta'} \right) + o_{p} (1)$$

$$= \frac{1}{N} \sum_{i=1}^{N} \widetilde{\iota}(\widehat{k}_{i}) \widetilde{z}(\widehat{k}_{i})' \left( -v_{i}^{\alpha} \left( \widehat{\alpha}(\widehat{k}_{i}) \right) \right) \left( \frac{\partial \widehat{\alpha}(\widehat{k}_{i})}{\partial \theta'} - \frac{\partial \overline{\alpha}_{i}}{\partial \theta'} \right) + o_{p} (1)$$

$$= \frac{1}{N} \sum_{i=1}^{N} \widetilde{\iota}(\widehat{k}_{i}) \widetilde{z}(\widehat{k}_{i})' \left( \left( v_{i}^{\theta} \left( \widehat{\alpha}(\widehat{k}_{i}) \right) \right)' - \left( -v_{i}^{\alpha} \left( \widehat{\alpha}(\widehat{k}_{i}) \right) \right) \frac{\partial \overline{\alpha}_{i}}{\partial \theta'} \right) + o_{p} (1)$$

$$= \frac{1}{N} \sum_{i=1}^{N} \widetilde{\iota}(\widehat{k}_{i}) \widetilde{z}(\widehat{k}_{i})' \left( \left( \underbrace{\mathbb{E}_{\eta_{i0}} \left( v_{i}^{\theta} \right) \right)' - \left( \mathbb{E}_{\eta_{i0}} \left( -v_{i}^{\alpha} \right) \right) \frac{\partial \overline{\alpha}_{i}}{\partial \theta'} \right) + o_{p} (1) = o_{p} (1)$$

where we have used (E8) in the third equality, (A33) and the expression of  $\partial \hat{\alpha}(k) / \partial \theta'$  in the fifth one, and we have expanded around  $\alpha_{i0}$  and used (E8) in the last equality.

Finally, to show (E2), let us define, analogously to the beginning of the proof of (A3):

$$\widehat{\iota}_i = \mathbf{1} \left\{ \|\widehat{\alpha}(\widehat{k}_i, \widehat{\theta}) - \alpha_{i0}\| \le \epsilon \right\},\$$

where  $\epsilon$  is such that:  $\inf_{\theta,\eta=(\alpha',\mu')'} \inf_{\|(\widetilde{\alpha},\widetilde{\theta})-(\alpha,\theta)\|\leq 2\epsilon} \mathbb{E}_{\eta_{i0}=\eta}\left(-v_i^{\alpha}(\widetilde{\alpha},\widetilde{\theta})\right) \geq \underline{\Sigma}$ . Using that  $\widehat{\theta}$  is consistent it is easy to verify that:

$$\left|\frac{1}{N}\sum_{i=1}^{N}\ell_{i}\left(\widehat{\alpha}(\widehat{k}_{i},\widehat{\theta}),\widehat{\theta}\right) - \frac{1}{N}\sum_{i=1}^{N}\ell_{i}\left(\widehat{\alpha}_{i}(\widehat{\theta}),\widehat{\theta}\right)\right| \leq \left|\frac{1}{N}\sum_{i=1}^{N}\ell_{i}\left(a(\widehat{k}_{i},\widehat{\theta}),\widehat{\theta}\right) - \frac{1}{N}\sum_{i=1}^{N}\ell_{i}\left(\widehat{\alpha}_{i}(\widehat{\theta}),\widehat{\theta}\right)\right| = O_{p}(\delta).$$
(E11)

Using similar arguments as at the beginning of the proof of (A3), but now at  $\hat{\theta}$ , we obtain that:

$$\frac{1}{N}\sum_{i=1}^{N}(1-\hat{\iota}_i) = O_p(\delta), \quad \frac{1}{N}\sum_{i=1}^{N}\hat{\iota}_i \left\|\widehat{\alpha}(\widehat{k}_i,\widehat{\theta}) - \widehat{\alpha}_i(\widehat{\theta})\right\|^2 = O_p(\delta).$$

hence that:

$$\frac{1}{N}\sum_{i=1}^{N}\left\|\widehat{\alpha}(\widehat{k}_{i},\widehat{\theta})-\widehat{\alpha}_{i}(\widehat{\theta})\right\|^{2}=O_{p}(\delta).$$

(E2) then comes from the fact that  $\frac{1}{N} \sum_{i=1}^{N} \|\widehat{\alpha}_i(\widehat{\theta}) - \alpha_{i0}\|^2 = O_p(\delta).$ This ends the proof of Corollary E1.

### E.2 Analytical expression of the asymptotic bias

The next result provides an explicit characterization of the  $O_p(1/T)$  bias in (9), in cases where K grows relatively fast relative to T. We recall that here p is fixed.

**Corollary E2.** Suppose that the conditions of Theorem 1 are satisfied. Let  $\hat{\alpha}_i(\theta) = \operatorname{argmax}_{\alpha_i} \ell_i(\alpha_i, \theta)$ , and let  $\hat{g}_i(\theta) = \frac{\partial^2 \ell_i(\hat{\alpha}_i(\theta), \theta)}{\partial \theta \partial \alpha'} (\frac{\partial^2 \ell_i(\hat{\alpha}_i(\theta), \theta)}{\partial \alpha \partial \alpha'})^{-1}$ . Suppose in addition:

- (i) There is an v > 0 such that  $T B_{\xi}(K^{1-v}) \xrightarrow{p} 0$  as N, T, K tend to infinity. Moreover, for any diverging  $K_{N,T}$  sequence,  $T B_{\varepsilon}(K_{N,T}) \xrightarrow{p} 0$  as N, T tend to infinity.
- (ii)  $\ell_i$  is four times differentiable, and its fourth derivatives satisfy similar uniform boundedness properties as the first three.
- (iii)  $\gamma(h) = \mathbb{E}_{h_i=h}[\widehat{\alpha}_i(\theta_0)]$  and  $\lambda(h) = \mathbb{E}_{h_i=h}[\widehat{g}_i(\theta_0)]$  are differentiable with respect to h, uniformly bounded with uniformly bounded first derivatives. Moreover, uniformly in h,  $\mathbb{E}_{h_i=h}[\|\widehat{\alpha}_i(\theta_0) - \gamma(h_i)\|^2] = O(T^{-1})$ ,  $\mathbb{E}_{h_i=h}[\|\widehat{g}_i(\theta_0) - \lambda(h_i)\|^2] = O(T^{-1})$ ,  $\mathbb{E}_{h_i=h}[\|\widehat{\alpha}_i(\theta_0) - \gamma(h_i)\|^3] = o(T^{-1})$ , and  $\mathbb{E}_{h_i=h}[\|\widehat{g}_i(\theta_0) - \lambda(h_i)\|^3] = o(T^{-1})$ .

Then, for p fixed as N, T, K tend to infinity such that K/N tends to zero we have:

$$\widehat{\theta} = \theta_0 + H^{-1} \frac{1}{N} \sum_{i=1}^N s_i + \frac{C}{T} + o_p \left(\frac{1}{T}\right) + o_p \left(\frac{1}{\sqrt{NT}}\right), \quad (E12)$$

where the expression of the constant C is given by (E24)-(E25) in the proof.

Under the conditions of Corollary E2, the kmeans objective is:

$$\frac{1}{N}\sum_{i=1}^{N} \|h_i - \hat{h}(\hat{k}_i)\|^2 = o_p\left(\frac{1}{T}\right).$$
(E13)

This happens when K grows sufficiently fast relative to T. As an example, when  $\xi_{i0}$  scalar and  $B_{\xi}(K) = O_p(K^{-2})$ , the condition requires  $TK^{-2}$  to tend to zero. The condition on  $B_{\varepsilon}$  is satisfied when  $\varepsilon_i$  is normal with zero mean and variance  $\Sigma/T$  for some  $\Sigma > 0$ , for example.

Corollary E2 shows that, when K is sufficiently large so that the approximation error  $B_{\xi}(K)$  is small relative to 1/T, and when in addition K/N tends to zero, the GFE estimator of  $\theta_0$  satisfies an expansion similar to that of the fixed-effects estimator, with a different first-order bias term; see, e.g., Hahn and Newey (2004, p.1302) for an expression of the bias of fixed-effects.

**Proof of Corollary E2.** We first show (E13). Let  $\{k_i\} = \{k_{i1}\} \cap \{k_{i2}\}$  be the intersection of two partitions of  $\{1, ..., N\}$ : a first partition with (the integer part of)  $K_1 = K^{1-v}$  groups, and a

second partition with (the integer part of)  $K_2 = K^{\nu}$  groups. Since  $(\hat{h}, \{\hat{k}_i\})$  solves (2), we have, using Condition (*i*) in the corollary:

$$\frac{1}{N}\sum_{i=1}^{N} \left\|h_{i} - \widehat{h}(\widehat{k}_{i})\right\|^{2} = \frac{1}{N}\sum_{i=1}^{N} \left\|\varphi(\xi_{i0}) + \varepsilon_{i} - \widehat{h}(\widehat{k}_{i})\right\|^{2}$$

$$\leq \min_{(\widetilde{h}_{1},\widetilde{h}_{2},\{k_{i1}\},\{k_{i2}\})} \frac{1}{N}\sum_{i=1}^{N} \left\|\varphi(\xi_{i0}) - \widetilde{h}_{1}(k_{i1}) + \varepsilon_{i} - \widetilde{h}_{2}(k_{i2})\right\|^{2}$$

$$= O_{p}\left(B_{\xi}(K_{1})\right) + O_{p}\left(B_{\varepsilon}\left(K_{2}\right)\right) = O_{p}\left(\frac{1}{T}\right).$$

To show that (E12) holds, we will follow a likelihood approach as in Arellano and Hahn (2007, 2016). Consider the difference between the GFE and fixed-effects concentrated likelihoods:

$$\Delta L(\theta) = \frac{1}{N} \sum_{i=1}^{N} \ell_i(\widehat{\alpha}(\widehat{k}_i, \theta), \theta) - \frac{1}{N} \sum_{i=1}^{N} \ell_i(\widehat{\alpha}_i(\theta), \theta).$$

We are going to derive an expansion for the derivative of  $\Delta L(\theta)$  at  $\theta_0$ . From there, we will characterize the first-order bias of the GFE estimator  $\hat{\theta}$ .

Specifically, we are going to show that:

$$\frac{\partial}{\partial \theta}\Big|_{\theta_0} \Delta L(\theta) = -\frac{\partial}{\partial \theta}\Big|_{\theta_0} \frac{1}{2N} \sum_{i=1}^N \nu_i(\theta)' \mathbb{E}_{\eta_{i0}} \left[-v_i^{\alpha}\left(\overline{\alpha}_i(\theta), \theta\right)\right] \nu_i(\theta) + o_p\left(\frac{1}{T}\right), \quad (E14)$$

where  $\nu_i(\theta) = \widehat{\alpha}_i(\theta) - \mathbb{E}_{h_i}(\widehat{\alpha}_i(\theta))$ , and we denote  $\overline{\alpha}_i(\theta) \equiv \overline{\alpha}(\theta, \eta_{i0})$ . In this case also, we can define  $\xi_{i0} = \eta_{i0}$ .

To show (E14) we are first going to establish several preliminary results. Together with fourthorder differentiability in Condition (*ii*) of the corollary, those will allow us to derive the required expansions. In the following we will evaluate all functions at  $\theta_0$ , and omit  $\theta_0$  for the notation. In particular,  $\hat{\alpha}_i$  will be a shorthand for  $\hat{\alpha}_i(\theta_0)$ .

First, note that, from the proof of Theorem 1, and using (E13), we have:

$$\frac{1}{N}\sum_{i=1}^{N} \|\widehat{\alpha}(\widehat{k}_i) - \widehat{\alpha}_i\|^2 = O_p\left(\frac{1}{T}\right).$$
(E15)

Next, let  $\widehat{\alpha}_i = \gamma(h_i) + \nu_i$ , where  $\gamma(h) = \mathbb{E}_{h_i = h}(\widehat{\alpha}_i)$ . We have:

$$\widehat{\alpha}(k) = \left(\sum_{i=1}^{N} \mathbf{1}\{\widehat{k}_i = k\}(-v_i^{\alpha}(a_i(k)))\right)^{-1} \left(\sum_{i=1}^{N} \mathbf{1}\{\widehat{k}_i = k\}(-v_i^{\alpha}(a_i(k)))\widehat{\alpha}_i\right),$$
(E16)

for some  $a_i(k)$  between  $\hat{\alpha}_i$  and  $\hat{\alpha}(k)$ . Note that, by the conditions of Theorem 1,  $(-v_i^{\alpha}(\alpha))$  is uniformly bounded away from zero with probability approaching one. Let  $\hat{\gamma}(k)$  and  $\hat{\nu}(k)$  denote the weighted means of  $\gamma(h_i)$  and  $\nu_i$  in group  $\hat{k}_i = k$ , respectively, where the weight is  $(-v_i^{\alpha}(a_i(k)))$ . Note that  $\hat{\alpha}(k) = \hat{\gamma}(k) + \hat{\nu}(k)$ . By (E13), and since  $\gamma$  is Lipschitz-continuous, we have:

$$\frac{1}{N}\sum_{i=1}^{N} \|\gamma(h_i) - \widehat{\gamma}(\widehat{k}_i)\|^2 = o_p\left(\frac{1}{T}\right).$$
(E17)

Moreover, since by Condition (iii) in the corollary the  $\sqrt{T}\nu_i$ , which are mean independent of the  $\hat{k}_{i'}$ 's and have zero mean, have bounded conditional variance, and denoting as  $\overline{\nu}(k)$  the unweighted mean of  $\nu_i$  in group  $\hat{k}_i = k$ , we have:  $\frac{1}{N}\sum_{i=1}^N \|\overline{\nu}(\hat{k}_i)\|^2 = O_p\left(\frac{K}{NT}\right) = o_p\left(\frac{1}{T}\right)$ , where we have used that p is fixed and K/N tends to zero. Hence:

$$\frac{1}{N}\sum_{i=1}^{N}\|\widehat{\nu}(\widehat{k}_i)\|^2 = o_p\left(\frac{1}{T}\right).$$
(E18)

Let  $\widehat{g}_i = v_i^{\theta}(\widehat{\alpha}_i)(-v_i^{\alpha}(\widehat{\alpha}_i))^{-1} = \lambda(h_i) + \tau_i$ , where  $\lambda(h) = \mathbb{E}_{h_i=h}(\widehat{g}_i)$ . Similarly to (E18), we have, using analogous notations for weighted group means:

$$\frac{1}{N}\sum_{i=1}^{N} \|\lambda(h_i) - \widehat{\lambda}(\widehat{k}_i)\|^2 = o_p\left(\frac{1}{T}\right), \quad \frac{1}{N}\sum_{i=1}^{N} \|\widehat{\tau}(\widehat{k}_i)\|^2 = o_p\left(\frac{1}{T}\right).$$
(E19)

Further, denote as  $\tilde{\gamma}(k)$ ,  $\tilde{\nu}(k)$ ,  $\tilde{\lambda}(k)$ , and  $\tilde{\tau}(k)$  the weighted means of  $\gamma(h_i)$ ,  $\nu_i$ ,  $\lambda(h_i)$ , and  $\tau_i$  in group  $\hat{k}_i = k$ , respectively, where the weight is  $(-v_i^{\alpha}(\hat{\alpha}_i))$ . Using similar arguments we obtain:

$$\frac{1}{N}\sum_{i=1}^{N} \|\gamma(h_i) - \widetilde{\gamma}(\widehat{k}_i)\|^2 = o_p\left(\frac{1}{T}\right), \quad \frac{1}{N}\sum_{i=1}^{N} \|\widetilde{\nu}(\widehat{k}_i)\|^2 = o_p\left(\frac{1}{T}\right), \quad (E20)$$

$$\frac{1}{N}\sum_{i=1}^{N} \|\lambda(h_i) - \widetilde{\lambda}(\widehat{k}_i)\|^2 = o_p\left(\frac{1}{T}\right), \quad \frac{1}{N}\sum_{i=1}^{N} \|\widetilde{\tau}(\widehat{k}_i)\|^2 = o_p\left(\frac{1}{T}\right).$$
(E21)

Next, using (E16), (E17), and (E18), in addition to the parameter space for  $\alpha_{i0}$  being compact and  $\gamma$  being bounded, we have that:  $\frac{1}{N} \sum_{i=1}^{N} \|\widehat{\alpha}(\widehat{k}_i) - \widehat{\alpha}_i\|^3 = -\frac{1}{N} \sum_{i=1}^{N} \|\nu_i\|^3 + o_p(1/T)$ . Hence, by Condition (iii) in the corollary:

$$\frac{1}{N}\sum_{i=1}^{N}\|\widehat{\alpha}(\widehat{k}_{i}) - \widehat{\alpha}_{i}\|^{3} = o_{p}\left(\frac{1}{T}\right).$$
(E22)

To see that (E14) holds, first note that, denoting  $a^{\otimes 2} = a \otimes a$ :

$$\begin{split} \frac{\partial}{\partial \theta} \Big|_{\theta_0} \Delta L(\theta) &= \frac{1}{N} \sum_{i=1}^N \frac{\partial \ell_i(\widehat{\alpha}(\widehat{k}_i))}{\partial \theta} - \frac{1}{N} \sum_{i=1}^N \frac{\partial \ell_i(\widehat{\alpha}_i)}{\partial \theta} \\ &= \frac{1}{N} \sum_{i=1}^N v_i^{\theta}(\widehat{\alpha}_i) \left(\widehat{\alpha}(\widehat{k}_i) - \widehat{\alpha}_i\right) + \frac{1}{2N} \sum_{i=1}^N v_i^{\theta\alpha}(a_i) \left(\widehat{\alpha}(\widehat{k}_i) - \widehat{\alpha}_i\right)^{\otimes 2} \\ &= \underbrace{\frac{1}{N} \sum_{i=1}^N v_i^{\theta}(\widehat{\alpha}_i) \left(\widetilde{\alpha}(\widehat{k}_i) - \widehat{\alpha}_i\right)}_{\equiv A_1} + \underbrace{\frac{1}{2N} \sum_{i=1}^N v_i^{\theta\alpha}(a_i) \left(\widehat{\alpha}(\widehat{k}_i) - \widehat{\alpha}_i\right)^{\otimes 2}}_{\equiv A_2} \\ &+ \underbrace{\frac{1}{2N} \sum_{i=1}^N v_i^{\theta}(\widehat{\alpha}_i) \left(\sum_{i'=1}^N \mathbf{1}\{\widehat{k}_{i'} = \widehat{k}_i\}(-v_{i'}^{\alpha}(\widehat{\alpha}_{i'}))\right)^{-1} \sum_{i'=1}^N \mathbf{1}\{\widehat{k}_{i'} = \widehat{k}_i\} v_{i'}^{\alpha\alpha} \left(a_{i'}(\widehat{k}_{i'})\right) \left(\widehat{\alpha}(\widehat{k}_{i'}) - \widehat{\alpha}_{i'}\right)^{\otimes 2}, \\ &= \underbrace{A_3} \end{split}$$

where  $a_i$  lies between  $\hat{\alpha}_i$  and  $\hat{\alpha}(\hat{k}_i)$  and so does  $a_i(\hat{k}_i)$ ,  $v_i^{\theta\alpha}(a_i)$  is a matrix of third derivatives with  $(\dim \alpha_{i0})^2$  columns, and we have defined, for all k:

$$\widetilde{\alpha}(k) = \left(\sum_{i=1}^{N} \mathbf{1}\{\widehat{k}_i = k\}(-v_i^{\alpha}(\widehat{\alpha}_i))\right)^{-1} \sum_{i=1}^{N} \mathbf{1}\{\widehat{k}_i = k\}(-v_i^{\alpha}(\widehat{\alpha}_i))\widehat{\alpha}_i,$$
(E23)

where we note that  $(-v_i^{\alpha}(\hat{\alpha}_i))$  is uniformly bounded away from zero with probability approaching one.

Let us consider the three terms in turn. First, we have:

$$\begin{split} A_1 &= \frac{1}{N} \sum_{i=1}^N \widehat{g}_i \left( -v_i^{\alpha}(\widehat{\alpha}_i) \right) \left( \widetilde{\alpha}(\widehat{k}_i) - \widehat{\alpha}_i \right) \\ &= \frac{1}{N} \sum_{i=1}^N \left( \widehat{g}_i - \widetilde{g}(\widehat{k}_i) \right) \left( -v_i^{\alpha}(\widehat{\alpha}_i) \right) \left( \widetilde{\alpha}(\widehat{k}_i) - \widehat{\alpha}_i \right) \\ &= -\frac{1}{N} \sum_{i=1}^N \left( \lambda(h_i) - \widetilde{\lambda}(\widehat{k}_i) + \tau_i - \widetilde{\tau}(\widehat{k}_i) \right) \left( -v_i^{\alpha}(\widehat{\alpha}_i) \right) \left( \gamma(h_i) - \widetilde{\gamma}(\widehat{k}_i) + \nu_i - \widetilde{\nu}(\widehat{k}_i) \right) \\ &= -\frac{1}{N} \sum_{i=1}^N \tau_i (-v_i^{\alpha}(\widehat{\alpha}_i)) \nu_i + o_p \left( \frac{1}{T} \right) \\ &= -\frac{1}{N} \sum_{i=1}^N \tau_i \mathbb{E}_{\eta_{i0}} (-v_i^{\alpha}(\alpha_{i0})) \nu_i + o_p \left( \frac{1}{T} \right), \end{split}$$

where we have used (E16), (E20), and (E21).

Next, we have, using in addition (E22):

$$A_{2} = \frac{1}{2N} \sum_{i=1}^{N} \mathbb{E}_{\eta_{i0}} \left( v_{i}^{\theta\alpha}(\alpha_{i0}) \right) \left( \widehat{\alpha}(\widehat{k}_{i}) - \widehat{\alpha}_{i} \right)^{\otimes 2} + o_{p} \left( \frac{1}{T} \right)$$
$$= \frac{1}{2N} \sum_{i=1}^{N} \mathbb{E}_{\eta_{i0}} \left( v_{i}^{\theta\alpha}(\alpha_{i0}) \right) \left( \widehat{\gamma}(\widehat{k}_{i}) - \gamma(h_{i}) + \widehat{\nu}(\widehat{k}_{i}) - \nu_{i} \right)^{\otimes 2} + o_{p} \left( \frac{1}{T} \right)$$
$$= \frac{1}{2N} \sum_{i=1}^{N} \mathbb{E}_{\eta_{i0}} \left( v_{i}^{\theta\alpha}(\alpha_{i0}) \right) \nu_{i}^{\otimes 2} + o_{p} \left( \frac{1}{T} \right).$$

Lastly, defining  $\tilde{g}(k)$  the weighted mean of  $\hat{g}_i$  in group  $\hat{k}_i = k$  with weight  $(-v_i^{\alpha}(\hat{\alpha}_i))$ , we have:

$$\begin{split} A_{3} &= \frac{1}{2N} \sum_{i=1}^{N} \widehat{g}_{i}(-v_{i}^{\alpha}(\widehat{\alpha}_{i})) \left( \sum_{i'=1}^{N} \mathbf{1}\{\widehat{k}_{i'} = \widehat{k}_{i}\}(-v_{i'}^{\alpha}(\widehat{\alpha}_{i'})) \right)^{-1} \\ &\times \sum_{i'=1}^{N} \mathbf{1}\{\widehat{k}_{i'} = \widehat{k}_{i}\}v_{i'}^{\alpha\alpha} \left(a_{i'}(\widehat{k}_{i'})\right) \left(\widehat{\alpha}(\widehat{k}_{i'}) - \widehat{\alpha}_{i'}\right)^{\otimes 2} \\ &= \frac{1}{2N} \sum_{i=1}^{N} \widetilde{g}(\widehat{k}_{i})(-v_{i}^{\alpha}(\widehat{\alpha}_{i})) \left( \sum_{i'=1}^{N} \mathbf{1}\{\widehat{k}_{i'} = \widehat{k}_{i}\}(-v_{i'}^{\alpha}(\widehat{\alpha}_{i'})) \right)^{-1} \\ &\times \sum_{i'=1}^{N} \mathbf{1}\{\widehat{k}_{i'} = \widehat{k}_{i}\}v_{i'}^{\alpha\alpha} \left(a_{i'}(\widehat{k}_{i'})\right) \left(\widehat{\alpha}(\widehat{k}_{i'}) - \widehat{\alpha}_{i'}\right)^{\otimes 2} + o_{p}\left(\frac{1}{T}\right) \\ &= \frac{1}{2N} \sum_{i=1}^{N} \widetilde{g}(\widehat{k}_{i})v_{i}^{\alpha\alpha} \left(a_{i}(\widehat{k}_{i})\right) \left(\widehat{\alpha}(\widehat{k}_{i}) - \widehat{\alpha}_{i}\right)^{\otimes 2} + o_{p}\left(\frac{1}{T}\right) \\ &= \frac{1}{2N} \sum_{i=1}^{N} \mathbb{E}_{\eta_{i0}} \left(v_{i}^{\theta}(\alpha_{i0})\right) \left[\mathbb{E}_{\eta_{i0}}(-v_{i}^{\alpha}(\alpha_{i0}))\right]^{-1} \mathbb{E}_{\eta_{i0}} \left[v_{i}^{\alpha\alpha} \left(\alpha_{i0}\right)\right] \nu_{i}^{\otimes 2} + o_{p}\left(\frac{1}{T}\right). \end{split}$$

Combining results, we get:

$$\begin{aligned} \frac{\partial}{\partial \theta} \Big|_{\theta_0} \Delta L(\theta) &= -\frac{1}{N} \sum_{i=1}^N \tau_i \mathbb{E}_{\eta_{i0}} \left( -v_i^{\alpha}(\alpha_{i0}) \right) \nu_i \\ &+ \frac{1}{2N} \sum_{i=1}^N \left[ \mathbb{E}_{\eta_{i0}} \left( v_i^{\theta\alpha}(\alpha_{i0}) \right) + \mathbb{E}_{\eta_{i0}} \left( v_i^{\theta}(\alpha_{i0}) \right) \left[ \mathbb{E}_{\eta_{i0}} (-v_i^{\alpha}(\alpha_{i0})) \right]^{-1} \mathbb{E}_{\eta_{i0}} \left[ v_i^{\alpha\alpha}(\alpha_{i0}) \right] \nu_i^{\otimes 2} + o_p \left( \frac{1}{T} \right). \end{aligned}$$

This shows (E14), since  $\frac{\partial \widehat{\alpha}_i(\theta_0)}{\partial \theta'} = \widehat{g}'_i$ , and:

$$\frac{\partial}{\partial \theta'}\Big|_{\theta_0} \operatorname{vec} \mathbb{E}_{\eta_{i0}} \left[ -v_i^{\alpha} \left( \overline{\alpha}_i(\theta), \theta \right) \right] = -\left( \mathbb{E}_{\eta_{i0}} \left( v_i^{\theta\alpha}(\alpha_{i0}) \right) + \mathbb{E}_{\eta_{i0}} \left( v_i^{\theta}(\alpha_{i0}) \right) \left[ \mathbb{E}_{\eta_{i0}}(-v_i^{\alpha}(\alpha_{i0})) \right]^{-1} \mathbb{E}_{\eta_{i0}} \left[ v_i^{\alpha\alpha}(\alpha_{i0}) \right] \right)'.$$

Equation (E14) readily delivers an expression for the first-order bias term of the GFE estimator.

Indeed, using that (e.g., Arellano and Hahn, 2007):

$$\begin{split} &\frac{\partial}{\partial \theta}\Big|_{\theta_0} \frac{1}{N} \sum_{i=1}^N \ell_i(\overline{\alpha}_i(\theta), \theta) - \frac{1}{N} \sum_{i=1}^N \ell_i(\widehat{\alpha}_i(\theta), \theta) \\ &= -\frac{\partial}{\partial \theta}\Big|_{\theta_0} \frac{1}{2N} \sum_{i=1}^N v_i\left(\overline{\alpha}_i(\theta), \theta\right)' \mathbb{E}_{\eta_{i0}} \left[-v_i^{\alpha}\left(\overline{\alpha}_i(\theta), \theta\right)\right]^{-1} v_i\left(\overline{\alpha}_i(\theta), \theta\right) + o_p\left(\frac{1}{T}\right), \end{split}$$

it follows that Corollary E2 holds, with:

$$C = H^{-1} \lim_{N, T \to \infty} \frac{1}{N} \sum_{i=1}^{N} \frac{\partial}{\partial \theta} \Big|_{\theta_0} T b_i(\theta),$$
(E24)

for:

$$b_{i}(\theta) = -\frac{1}{2} \left( \widehat{\alpha}_{i}(\theta) - \mathbb{E}_{h_{i}}\left( \widehat{\alpha}_{i}(\theta) \right) \right)' \mathbb{E}_{\eta_{i0}} \left[ -v_{i}^{\alpha} \left( \overline{\alpha}_{i}(\theta), \theta \right) \right] \left( \widehat{\alpha}_{i}(\theta) - \mathbb{E}_{h_{i}}\left( \widehat{\alpha}_{i}(\theta) \right) \right) + \frac{1}{2} v_{i} \left( \overline{\alpha}_{i}(\theta), \theta \right)' \mathbb{E}_{\eta_{i0}} \left[ -v_{i}^{\alpha} \left( \overline{\alpha}_{i}(\theta), \theta \right) \right]^{-1} v_{i} \left( \overline{\alpha}_{i}(\theta), \theta \right).$$
(E25)

Bias in a regression example. Consider the following model for a scalar outcome:

$$Y_{it} = \rho_0 Y_{i,t-1} + X'_{it} \beta_0 + \alpha_{i0} + U_{it},$$
(E26)

where  $|\rho_0| < 1$ . In this case,  $\widehat{\alpha}_i(\theta) = (1-\rho)\overline{Y}_i - \overline{X}'_i\beta + o_p\left(T^{-\frac{1}{2}}\right)$ . Hence, when classifying individuals based on  $h_i = \left(\overline{Y}_i, \overline{X}'_i\right)', \widehat{\alpha}_i(\theta)$  belongs to the span of  $h_i$ , up to small-order terms. It follows that C/T is identical to the first-order bias  $C^{\text{FE}}/T$  of the fixed-effects estimator, and fixed-effects and two-step GFE are first-order equivalent.

However, this equivalence does not hold generally. As an example, let us suppose the unobservables  $(\alpha_{i0}, \mu'_{i0})'$  follow a one-factor structure with  $\mu_{i0} = \lambda \alpha_{i0}$  for a vector  $\lambda$ , and let us classify individual units based on  $h_i = \overline{Y}_i$  only. Injectivity is satisfied in this example, due to the low underlying dimensionality of  $\eta_{i0} = (\alpha_{i0}, \mu'_{i0})'$ . In this case it can be verified that:

$$\mathbb{E}_{h_i}\left(\widehat{\alpha}_i(\theta)\right) = \left(\frac{\frac{1-\rho}{1-\rho_0} + \left(\frac{1-\rho}{1-\rho_0}\beta_0 - \beta\right)'\lambda}{\frac{1}{1-\rho_0} + \frac{\beta_0'\lambda}{1-\rho_0}}\right)\overline{Y}_i + o_p\left(\frac{1}{\sqrt{T}}\right),$$

and, letting  $V_{it} = X_{it} - \lambda \alpha_{i0}$ :

$$\nu_i(\theta) = \beta' \frac{\lambda U_i - V_i}{1 + \beta'_0 \lambda} + o_p\left(\frac{1}{\sqrt{T}}\right)$$

As a result, the first-order bias term on  $\rho_0$  is the same for GFE and fixed-effects, while for  $\beta_0$  we have, letting  $\frac{1}{T} \sum_{t=1}^{T} \mathbb{E}(V_{it}V'_{it}) = \Sigma > 0$ :

$$C = C^{\text{FE}} - \Sigma^{-1} \mathbb{E} \left[ T \left( \lambda \overline{U}_i - \overline{V}_i \right) \left( \lambda \overline{U}_i - \overline{V}_i \right)' \right] \frac{\beta_0}{\left( 1 + \beta'_0 \lambda \right)^2},$$

so  $C \neq C^{\text{FE}}$  in general.

## **F** Verification of assumptions in linear models

In this section we verify the assumptions of Theorem 1 in the linear model:

$$\begin{cases} Y_{it} = \rho_0 Y_{i,t-1} + X'_{it} \beta_0 + \alpha_{it0} + U_{it}, \\ X_{it} = \mu_{it0} + V_{it}, \end{cases}$$
(F1)

where  $Y_{it}$  is a scalar outcome and  $X_{it}$  is a vector of strictly exogenous covariates. We consider two cases: time-invariant  $\alpha_{i0}, \mu_{i0}$ , and time-varying  $\alpha_{it0}, \mu_{it0}$ . Throughout, we use the notation  $Z_i = (Z'_{i1}, ..., Z'_{iT})'$ and  $\overline{Z}_i = \frac{1}{T} \sum_{t=1}^T Z_{it}$ .

### F.1 Time-invariant heterogeneity

Assumption F1. (regularity in linear model, time-invariant heterogeneity)

- (i) The parameter space  $\Theta$  for  $\theta_0 = (\rho_0, \beta'_0)'$ , and the space for  $(\alpha_{i0}, \mu'_{i0})$ , are compact.  $\theta_0$  belongs to the interior of  $\Theta$ , and  $|\rho_0| < 1$ .
- (ii) U<sub>it</sub> ~ N(0, σ<sub>0</sub><sup>2</sup>), i.i.d. across individuals and over time, independent of X<sub>i</sub>, α<sub>i0</sub>, μ<sub>i0</sub> and Y<sub>i0</sub>.
  V<sub>it</sub> ~ N(0, Σ<sub>0</sub><sup>2</sup>), i.i.d. across individuals and over time, independent of X<sub>i</sub>, α<sub>i0</sub>, μ<sub>i0</sub>, U<sub>i</sub>, and Y<sub>i0</sub>.
  Y<sub>i0</sub> is drawn from its stationary distribution conditional on α<sub>i0</sub> and μ<sub>i0</sub>.
- (iii) Let  $W_{it} = (Y_{i,t-1}, X'_{it})'$ .  $\mathbb{E}((W_{it} \mathbb{E}_{\alpha_{i0},\mu_{i0}}(W_{it})) (W_{it} \mathbb{E}_{\alpha_{i0},\mu_{i0}}(W_{it}))')$  is positive definite.

(iv) T tends to infinity, and GFE is based on the moment  $h_i = (\overline{Y}_i, \overline{X}'_i)'$ .

In Assumption F1 we suppose that  $U_{it}$  and  $V_{it}$  are normal homoskedastic. This can be relaxed, and we could work instead under stationary mixing conditions as in Hahn and Kuersteiner (2011). In that case the likelihood function would be interpreted as a pseudo-likelihood, and the formula for the asymptotic distribution should be adapted.

We are now going to verify the assumptions of Theorem 1, part (b). Assumption 1 holds letting  $\xi_{i0} = (\alpha_{i0}, \mu'_{i0})'$ , which has compact support. Note that in this model with time-invariant heterogeneity we can abstract from  $\lambda_{i0}$  without loss of generality. Next,  $\varphi(\xi_{i0}) = (\frac{\alpha_{i0} + \mu'_{i0}\beta_0}{1-\rho_0}, \mu'_{i0})'$ , and it is easy to check that  $\varphi$  is Lipschitz-continuous. Moreover, by part (ii) in Assumption F1 we have  $\frac{1}{N}\sum_{i=1}^{N} \|h_i - \varphi(\xi_{i0})\|^2 = O_p(1/T)$ . Lastly, letting  $\psi(h_1, h_2) = ((1-\rho_0)h_1 - h'_2\beta_0, h'_2)'$ , we have  $(\alpha_{i0}, \mu'_{i0})' = \psi(\varphi(\xi_{i0}))$ .  $\psi$  is Lipschitz-continuous by part (i) in Assumption F1. We have thus verified Assumptions 1 and 2, and the conditions of Lemmas 1 and 2.

Let us now verify Assumption 3, with part (iib). In the present case, p = 1 and R = T. We consider the log-likelihood function:

$$\ell_i(\alpha, \theta) = -\frac{1}{2T} \sum_{t=1}^T \left( Y_{it} - W'_{it}\theta - \alpha \right)^2,$$

where we have set  $\sigma_0^2 = 1$  without loss of generality. In addition, we have, for all  $\xi = (\alpha, \mu')'$ :

$$\overline{\alpha}(\theta,\xi) = \mathbb{E}_{\xi_{i0}=\xi}(\overline{Y}_i - \overline{W}'_i\theta) = \frac{1-\rho}{1-\rho_0}\alpha + \mu'\left(\frac{1-\rho}{1-\rho_0}\beta_0 - \beta\right).$$

Note that  $\overline{\alpha}(\theta_0, \xi_{i0}) = \alpha_{i0}$ .

It is easy to see that Assumption 3 part (i) is satisfied. Next, by part (i) in Assumption F1 the expected log-likelihood, and the expected log-likelihood derivatives, are bounded. Moreover,  $\frac{\partial^2 \ell_i(\alpha, \theta)}{\partial \alpha^2} = -1$ , and the third derivatives of  $\ell_i$  are zero. From the assumptions on  $U_{it}$ :

$$\frac{1}{N}\sum_{i=1}^{N} [\ell_i(\alpha_{i0},\theta_0) - \mathbb{E}_{\xi_{i0}}(\ell_i(\alpha_{i0},\theta_0))]^2 = \frac{1}{N}\sum_{i=1}^{N} \left(\frac{1}{T}\sum_{t=1}^{T}U_{it}^2 - \mathbb{E}_{\xi_{i0}}(U_{it}^2)\right)^2 = O_p(1/T).$$

For all first three derivatives of  $\ell_i$ , we obtain similar  $O_p(1/T)$  rates by using the conditions on  $U_{it}$ ,  $V_{it}$  and heterogeneity. This verifies Assumption 3 part (iib).

Next, the expected "target" log-likelihood  $\mathbb{E}\left[\ell_i(\overline{\alpha}(\theta,\xi_{i0}),\theta)\right]$  is quadratic in  $\theta$ , and its partial derivatives with respect to  $\rho$  and  $\beta$  are, respectively:

$$\mathbb{E}\left(\left(Y_{i,t-1} - \frac{\alpha_{i0} + \mu_{i0}'\beta_0}{1 - \rho_0}\right)\left(Y_{it} - \frac{\alpha_{i0} + \mu_{i0}'\beta_0}{1 - \rho_0} - \rho\left(Y_{i,t-1} - \frac{\alpha_{i0} + \mu_{i0}'\beta_0}{1 - \rho_0}\right) - (X_{it} - \mu_{i0})'\beta\right)\right),\$$

and:

$$\mathbb{E}\left((X_{it}-\mu_{i0})\left(Y_{it}-\frac{\alpha_{i0}+\mu_{i0}'\beta_{0}}{1-\rho_{0}}-\rho\left(Y_{i,t-1}-\frac{\alpha_{i0}+\mu_{i0}'\beta_{0}}{1-\rho_{0}}\right)-(X_{it}-\mu_{i0})'\beta\right)\right).$$

It is easy to verify that those are zero at  $\theta_0 = (\rho_0, \beta'_0)'$ . Moreover, the second derivative -H is negative definite by Assumption F1 (iii). Lastly, we have:

$$\sup_{\theta} \frac{1}{N} \sum_{i=1}^{N} \left\| \frac{\partial^2 \ell_i(\overline{\alpha}(\theta, \xi_{i0}), \theta)}{\partial \theta \partial \alpha} \right\|^2 = \frac{1}{N} \sum_{i=1}^{N} \|\overline{W}_i\|^2 = O_p(1),$$

where we have used the assumptions on  $U_{it}$ ,  $V_{it}$ , and the heterogeneity. This verifies Assumption 3 part (iii).

Turning to Assumption 3 part (iv), we have:

$$\frac{\partial}{\partial \xi'}\Big|_{\xi=\widetilde{\xi}} \mathbb{E}_{\xi_{i0}=\xi} \left(\frac{\partial^2 \ell_i(\alpha, \theta_0)}{\partial \theta \partial \alpha}\right) = -\frac{\partial}{\partial \xi'}\Big|_{\xi=\widetilde{\xi}} \mathbb{E}_{\xi_{i0}=\xi} \left(\overline{W}_i\right) = -\begin{pmatrix} \frac{1}{1-\rho_0} & \frac{\beta_0}{1-\rho_0} \\ 0 & I \end{pmatrix}$$

where I is the identity matrix of size dim  $\mu_{i0}$ . Similarly, we have:

$$\frac{\partial}{\partial \xi'}\Big|_{\xi=\widetilde{\xi}} \mathbb{E}_{\xi_{i0}=\xi} \left(\frac{\partial^2 \ell_i(\alpha, \theta_0)}{\partial \alpha^2}\right) = 0,$$

and:

$$\frac{\partial}{\partial \xi'}\Big|_{\xi=\widetilde{\xi}} \mathbb{E}_{\xi_{i0}=\xi} \left(\frac{\partial \ell_i(\overline{\alpha}(\theta,\widetilde{\xi}),\theta)}{\partial \alpha}\right) = \frac{\partial}{\partial \xi'}\Big|_{\xi=\widetilde{\xi}} \left(\overline{\alpha}(\theta,\xi) - \overline{\alpha}(\theta,\widetilde{\xi})\right) = \left(\frac{1-\rho}{1-\rho_0}, \frac{1-\rho}{1-\rho_0}\beta'_0 - \beta'\right),$$

which is bounded by Assumption F1 (i). This verifies Assumption 3 part (iv).

Lastly, turning to Assumption 3 part (v), we have:

$$\mathbb{E}_{h_i=h,\xi_{i0}=\xi}\left(\frac{\partial \ell_i(\overline{\alpha}(\theta,\xi),\theta)}{\partial \alpha}\right) = \mathbb{E}_{\varepsilon_i=h-\varphi(\xi),\xi_{i0}=\xi}\left(\overline{Y}_i - \overline{W}'_i\theta - \overline{\alpha}(\theta,\xi)\right).$$

It is easy to verify that, by Assumption F1 (ii), this is a linear function of  $h - \varphi(\xi)$  whose coefficients are uniformly bounded.

Likewise:

$$\mathbb{E}_{h_i=h,\xi_{i0}=\xi}\left(\operatorname{vec}\frac{\partial}{\partial\theta}\bigg|_{\theta_0}\frac{\partial\ell_i(\overline{\alpha}(\theta,\xi),\theta)}{\partial\alpha}\right)$$

is also linear in  $h - \varphi(\xi)$ , with coefficients that are uniformly bounded.

Moreover:

$$\operatorname{Var}_{h_{i}=h,\xi_{i0}=\xi}\left(\frac{\partial \ell_{i}(\overline{\alpha}(\theta,\xi),\theta)}{\partial \alpha}\right) = \operatorname{Var}_{h_{i}=h,\xi_{i0}=\xi}\left(\overline{Y}_{i}-\overline{W}_{i}^{\prime}\theta\right)$$

is non-stochastic by Assumption F1 (ii), and it can be shown that is satisfies:

$$\operatorname{Var}_{h_i=h,\xi_{i0}=\xi}\left(\overline{Y}_i-\overline{W}'_i\theta\right)=O_p\left(\frac{1}{T}\right),$$

uniformly in h,  $\xi$ , and  $\theta$ .

Finally, we have:

$$\operatorname{Var}_{h_{i}=h,\xi_{i0}=\xi}\left(\operatorname{vec}\frac{\partial}{\partial\theta}\Big|_{\theta_{0}}\frac{\partial\ell_{i}(\overline{\alpha}(\theta,\xi),\theta)}{\partial\alpha}\right)=\operatorname{Var}_{h_{i}=h,\xi_{i0}=\xi}\left(\overline{W}_{i}\right).$$

Now,  $\operatorname{Var}_{h_i=h,\xi_{i0}=\xi}\left(\overline{X}_i\right)=0$ . Moreover:

$$\operatorname{Var}_{h_i=h,\xi_{i0}=\xi}\left(\overline{Y}_{i,-1}\right) = O_p\left(\frac{1}{T}\right),$$

uniformly in h and  $\xi$ , where  $\overline{Y}_{i,-1} = \frac{1}{T} \sum_{t=1}^{T} Y_{i,t-1}$ .

This verifies Assumption  $3 \text{ part } (\mathbf{v})$ .

Applying Theorem 1, we obtain:

$$\widehat{\theta} = \theta_0 + H^{-1} \frac{1}{N} \sum_{i=1}^N s_i + O_p\left(\frac{1}{T}\right) + O_p\left(K^{-\frac{2}{\dim X_{it}+1}}\right) + o_p\left(\frac{1}{\sqrt{NT}}\right),$$

where  $H = \mathbb{E}\left(\left(W_{it} - \mathbb{E}_{\xi_{i0}}(W_{it})\right) \left(W_{it} - \mathbb{E}_{\xi_{i0}}(W_{it})\right)'\right)$ , and  $s_i = \frac{1}{T} \sum_{t=1}^{T} \left(W_{it} - \mathbb{E}_{\xi_{i0}}(W_{it})\right) U_{it}$ .

## F.2 Time-varying heterogeneity

We will verify the conditions in the time-varying case under the following assumptions, where for simplicity we focus on a scalar covariate  $X_{it}$ . The two unobservables have a factor structure, with a single factor in each equation.

Assumption F2. (regularity in linear model, time-varying heterogeneity)

- (i)  $\alpha_{it0} = \alpha_{i0}\lambda_{t0}^{\alpha}$  and  $\mu_{it0} = \mu_{i0}\lambda_{t0}^{\mu}$ , where  $\alpha_{i0}$ ,  $\mu_{i0}$ ,  $\lambda_{t0}^{\alpha}$ , and  $\lambda_{t0}^{\mu}$  are scalar.
- (ii)  $\lambda_{t0}^{\alpha}$  and  $\lambda_{t0}^{\mu}$  are stationary with bounded supports and non-zero means, and are independent of  $U_i$ and  $V_i$ . Moreover,  $\|\frac{1}{T}\sum_{t=1}^T \lambda_{t0}^{\alpha} - \mathbb{E}(\lambda_{t0}^{\alpha})\|^2 = O_p(1/T)$ , and  $\|\frac{1}{T}\sum_{t=1}^T \lambda_{t0}^{\mu} - \mathbb{E}(\lambda_{t0}^{\mu})\|^2 = O_p(1/T)$ .
- (iii) Assumption F1 holds, with part (iii) replaced by  $\mathbb{E}((W_{it} \mathbb{E}_{\alpha_{i0},\mu_{i0},\lambda_0}(W_{it})) (W_{it} \mathbb{E}_{\alpha_{i0},\mu_{i0},\lambda_0}(W_{it}))')$ being positive definite.

Assumption 1 holds by part (i) in Assumption F2. In addition, we have:

$$\varphi(\xi_{i0}) = \lim_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} \left( \begin{array}{c} \alpha_{i0} \lambda_{t0}^{\alpha} + \mu_{i0} \lambda_{t0}^{\mu} \beta_{0} + \sum_{s=1}^{+\infty} \rho_{0}^{s} \left[ \alpha_{i0} \lambda_{t-s,0}^{\alpha} + \mu_{i0} \lambda_{t-s,0}^{\mu} \beta_{0} \right] \\ \mu_{i0} \lambda_{t0}^{\mu} \end{array} \right),$$

where, as in the time-invariant case, we define  $\xi_{i0} = (\alpha_{i0}, \mu_{i0})'$ . Hence, by stationarity of the factors:

$$\varphi(\xi_{i0}) = \begin{pmatrix} \frac{\alpha_{i0} \mathbb{E}[\lambda_{t0}^{\alpha}] + \mu_{i0} \mathbb{E}[\lambda_{t0}^{\mu}]\beta_0}{1 - \rho_0} \\ \mu_{i0} \mathbb{E}[\lambda_{t0}^{\mu}] \end{pmatrix}.$$

It follows that  $\varphi$  is Lipschitz-continuous. Injectivity follows from  $\mathbb{E}[\lambda_{t0}^{\alpha}] \neq 0$  and  $\mathbb{E}[\lambda_{t0}^{\mu}] \neq 0$ .

Let  $\nu_{t0}^{\alpha} = \lambda_{t0}^{\alpha} - \mathbb{E}[\lambda_{t0}^{\alpha}]$ , and  $\nu_{t0}^{\mu} = \lambda_{t0}^{\mu} - \mathbb{E}[\lambda_{t0}^{\mu}]$ . We have:

$$\begin{split} h_{i} - \varphi(\xi_{i0}) = & \frac{1}{T} \sum_{t=1}^{T} \left( \begin{array}{c} U_{it} + V_{it}\beta_{0} + \sum_{s=1}^{+\infty} \rho_{0}^{s} \left[ U_{i,t-s} + V_{i,t-s}\beta_{0} \right] \\ V_{it} \end{array} \right) \\ &+ \frac{1}{T} \sum_{t=1}^{T} \left( \begin{array}{c} \alpha_{i0}\nu_{t0}^{\alpha} + \mu_{i0}\nu_{t0}^{\mu}\beta_{0} + \sum_{s=1}^{+\infty} \rho_{0}^{s} \left[ \alpha_{i0}\nu_{t-s,0}^{\alpha} + \mu_{i0}\nu_{t-s,0}^{\mu}\beta_{0} \right] \\ \mu_{i0}\nu_{t0}^{\mu} \end{array} \right). \end{split}$$

Using Assumption F1 and part (ii) in Assumption F2, we verify that  $\frac{1}{N}\sum_{i=1}^{N} \|h_i - \varphi(\xi_{i0})\|^2 = O_p(1/T)$ .

Let us now verify Assumption 3, with part (iib). In the present case, p = T and R = 1. We consider the log-likelihood function:

$$\ell_{it}(\alpha,\theta) = -\frac{1}{2} \left( Y_{it} - W'_{it}\theta - \alpha \right)^2,$$

where we have again set  $\sigma_0^2 = 1$ . We have, for all  $\xi = (\alpha, \mu)'$  and t:

$$\overline{\alpha}^{t}(\theta,\xi) = \mathbb{E}_{\xi_{i0}=\xi,\lambda_{0}^{\alpha},\lambda_{0}^{\mu}}(Y_{it} - W_{it}^{\prime}\theta)$$
$$= \alpha\lambda_{t0}^{\alpha} + (\beta_{0} - \beta)\mu\lambda_{t0}^{\mu} + (\rho_{0} - \rho)\alpha\sum_{s=1}^{\infty}\rho_{0}^{s-1}\lambda_{t-s,0}^{\alpha} + (\rho_{0} - \rho)\beta_{0}\mu\sum_{s=1}^{\infty}\rho_{0}^{s-1}\lambda_{t-s,0}^{\mu}$$

Note that  $\overline{\alpha}^t(\theta_0, \xi_{i0}) = \alpha_{i0}\lambda_{t0}^{\alpha} = \alpha_{it0}$ .

It is easy to see that Assumption 3 part (i) is satisfied. To check Assumption 3 part (iib), note that the expected log-likelihood is uniformly bounded by Assumption F1 and part (ii) in Assumption

F2, and similarly for its first three derivatives. As in the time-invariant case,  $\frac{\partial^2 \ell_{it}(\alpha,\theta)}{\partial \alpha^2} = -1$ , and the third derivatives of  $\ell_{it}$  are zero. Moreover:

$$\frac{1}{NT}\sum_{i=1}^{N}\sum_{t=1}^{T} [\ell_{it}(\alpha_{it0},\theta_0) - \mathbb{E}_{\xi_{i0},\lambda_0}(\ell_{it}(\alpha_{it0},\theta_0))]^2 = \frac{1}{NT}\sum_{i=1}^{N}\sum_{t=1}^{T} \left(U_{it}^2 - \mathbb{E}_{\alpha_{i0},\mu_{i0},\lambda_0^{\alpha},\lambda_0^{\mu}}(U_{it}^2)\right)^2 = O_p(1),$$

and all the first three derivatives of  $\ell_{it}$  are similarly  $O_p(1)$ . This verifies Assumption 3 part (iib).

Turning to part (iii), the function:

$$\mathbb{E}\left[\frac{1}{T}\sum_{t=1}^{T}\ell_{it}(\overline{\alpha}^{t}(\theta,\xi_{i0}),\theta)\right]$$

is quadratic in  $\theta$ , and its partial derivatives with respect to  $\rho$  and  $\beta$  are, respectively:

$$\mathbb{E}\left(\left(Y_{i,t-1} - \alpha_{i0}\sum_{s=1}^{+\infty}\rho_0^{s-1}\lambda_{t-s,0}^{\alpha} - \beta_0\mu_{i0}\sum_{s=1}^{+\infty}\rho_0^{s-1}\lambda_{t-s,0}^{\mu}\right)\left(Y_{it} - \rho Y_{i,t-1} - \beta X_{it} - \overline{\alpha}^t(\theta,\xi_{i0})\right)\right),$$

and:

$$\mathbb{E}\left(\left(X_{it}-\mu_{i0}\lambda_{t0}^{\mu}\right)\left(Y_{it}-\rho Y_{i,t-1}-\beta X_{it}-\overline{\alpha}^{t}(\theta,\xi_{i0})\right)\right)$$

These are zero at  $\theta_0$ . Moreover, the second derivative -H is negative definite by Assumption F2 (iii). Lastly, we have:

$$\sup_{\theta} \frac{1}{NT} \sum_{i=1}^{N} \sum_{t=1}^{T} \left\| \frac{\partial^2 \ell_{it}(\overline{\alpha}^t(\theta, \xi_{i0}), \theta)}{\partial \theta \partial \alpha} \right\|^2 = \frac{1}{NT} \sum_{i=1}^{N} \sum_{t=1}^{T} \|W_{it}\|^2 = O_p(1).$$

This verifies Assumption 3 part (iii).

Let us now check Assumption 3 part (iv). We have:

$$\begin{aligned} \frac{\partial}{\partial \xi'} \Big|_{\xi = \tilde{\xi}} \mathbb{E}_{\xi_{i0} = \xi, \lambda_0 = \lambda} \left( \frac{\partial^2 \ell_{it}(\alpha, \theta_0)}{\partial \theta \partial \alpha} \right) &= -\frac{\partial}{\partial \xi'} \Big|_{\xi = \tilde{\xi}} \mathbb{E}_{\xi_{i0} = \xi, \lambda_0 = \lambda} \left( W_{it} \right) \\ &= - \left( \begin{array}{cc} \sum_{s=1}^{\infty} \rho_0^{s-1} \lambda_{t-s}^{\alpha} & \beta_0 \sum_{s=1}^{\infty} \rho_0^{s-1} \lambda_{t-s}^{\mu} \\ 0 & \lambda_t^{\mu} \end{array} \right), \end{aligned}$$

which is uniformly bounded using Assumption F1 and part (ii) in Assumption F2. The second condition in Assumption 3 part (iv) is immediate to verify. Lastly, we have:

$$\begin{split} \frac{\partial}{\partial \xi'}\Big|_{\xi=\widetilde{\xi}} \mathbb{E}_{\xi_{i0}=\xi,\lambda_0=\lambda} \left(\frac{\partial \ell_{it}(\overline{\alpha}^t(\theta,\widetilde{\xi}),\theta)}{\partial \alpha}\right) &= \frac{\partial}{\partial \xi'}\Big|_{\xi=\widetilde{\xi}} \left(\overline{\alpha}^t(\theta,\xi) - \overline{\alpha}^t(\theta,\widetilde{\xi})\right) \\ &= \left(\lambda_t^{\alpha} + (\rho_0 - \rho)\sum_{s=1}^{+\infty} \rho_0^{s-1}\lambda_{t-s}^{\alpha} , \ (\beta_0 - \beta)\lambda_t^{\mu} + (\rho_0 - \rho)\beta_0\sum_{s=1}^{+\infty} \rho_0^{s-1}\lambda_{t-s}^{\mu}\right) \end{split}$$

which is also uniformly bounded. This verifies Assumption 3 part (iv).

Finally, let us verify Assumption 3 part (v). We have:

$$\mathbb{E}_{h_i=h,\xi_{i0}=\xi,\lambda_0=\lambda}\left(\frac{\partial\ell_{it}(\overline{\alpha}^t(\theta,\xi_{i0}),\theta)}{\partial\alpha}\right) = \mathbb{E}_{\varepsilon_i=h-\varphi(\xi),\xi_{i0}=\xi,\lambda_0=\lambda}\left(Y_{it}-W_{it}\theta-\overline{\alpha}^t(\theta,\xi_{i0})\right).$$

This is a linear function of  $h - \varphi(\xi)$  whose coefficients are uniformly bounded.

Likewise:

$$\mathbb{E}_{h_i=h,\xi_{i0}=\xi,\lambda_0=\lambda}\left(\frac{\partial}{\partial\theta}\bigg|_{\theta_0}\frac{\partial\ell_{it}(\overline{\alpha}^t(\theta,\xi_{i0}),\theta)}{\partial\alpha}\right)$$

is also linear in  $h - \varphi(\xi)$ , with coefficients that are uniformly bounded.

Moreover:

$$\operatorname{Var}_{h_{i}=h,\xi_{i0}=\xi,\lambda_{0}=\lambda}\left(\frac{\partial\ell_{it}(\overline{\alpha}^{t}(\theta,\xi_{i0}),\theta)}{\partial\alpha}\right) = \operatorname{Var}_{h_{i}=h,\xi_{i0}=\xi,\lambda_{0}=\lambda}\left(Y_{it}-W_{it}^{\prime}\theta\right)$$

is non-stochastic by Assumption F1 (ii), and it is easy to see that it is bounded uniformly in h,  $\xi$ ,  $\lambda$ , t, and  $\theta$ .

Similarly:

$$\operatorname{Var}_{h_{i}=h,\xi_{i0}=\xi,\lambda_{0}=\lambda}\left(\frac{\partial}{\partial\theta}\Big|_{\theta_{0}}\frac{\partial\ell_{it}(\overline{\alpha}^{t}(\theta,\xi),\theta)}{\partial\alpha}\right) = \operatorname{Var}_{h_{i}=h,\xi_{i0}=\xi,\lambda_{0}=\lambda}\left(W_{it}\right)$$

is uniformly bounded.

This verifies Assumption 3 part (v).

Applying Theorem 1, we obtain:

$$\widehat{\theta} = \theta_0 + \widetilde{H}^{-1} \frac{1}{N} \sum_{i=1}^N \widetilde{s}_i + O_p\left(\frac{1}{T}\right) + O_p\left(\frac{K}{N}\right) + O_p\left(K^{-\frac{2}{\dim X_{it}+1}}\right) + o_p\left(\frac{1}{\sqrt{NT}}\right),$$

where  $\widetilde{H} = \mathbb{E}\left(\left(W_{it} - \mathbb{E}_{\xi_{i0},\lambda_0}(W_{it})\right)\left(W_{it} - \mathbb{E}_{\xi_{i0},\lambda_0}(W_{it})\right)'\right)$ , and  $\widetilde{s}_i = \frac{1}{T}\sum_{t=1}^T \left(W_{it} - \mathbb{E}_{\xi_{i0},\lambda_0}(W_{it})\right)U_{it}$ .

## G Additional simulation results for the probit models

In this section we present additional results on the simulations of probit models. We start by describing the DGP in detail. We then show results about inference and estimation of K. Finally, we present two exercises related to computational aspects of GFE.

## G.1 Monte Carlo designs and estimators

We consider the following static probit model:

$$\begin{cases} Y_{it} = \mathbf{1} \{ X'_{it} \theta_{it0} + \alpha_{it0} + U_{it} > 0 \}, \\ X_{it} = \mu_{it0} + V_{it}, \end{cases}$$
(G1)

where  $U_{it}$  are i.i.d. standard normal, independent of  $X_i$ ,  $\alpha_{i0}$ ,  $\mu_{i0}$ , and  $V_{it}$  are i.i.d. normal with zero mean and identity covariance matrix, independent of all  $U_{it}$ 's,  $X_i$ ,  $\alpha_{i0}$ , and  $\mu_{i0}$ .

We consider three different specifications for  $\theta_{it0}$ ,  $\alpha_{it0}$ , and  $\mu_{it0}$ . In the *fixed-effects probit* specification, all components of  $\theta_{it0}$  are constant, equal to 1, and  $\mu_{it0} = \mu_{i0}$  are normal with zero mean and identity covariance matrix. In DGP 2 to 4, we set  $\alpha_{i0} = \frac{1}{\sqrt{r}} \sum_{r'=1}^{r} \mu_{i0}^{(r')} + \zeta_i$ , where  $\zeta_i$  are standard normal independent of  $\mu_{i0}$ , and  $\mu_{i0}^{(r')}$ , r' = 1, ..., r, denote the components of  $\mu_{i0}$ . In DGP 1 we set  $\alpha_{i0} = \mu_{i0}$ .

In the time-varying specification, all components of  $\theta_{it0}$  are constant, equal to 1,  $\alpha_{it0} = \alpha_{i0}f_t$  and  $\mu_{it0} = \mu_{i0}f_t$ , where  $\alpha_{i0}$  and  $\mu_{i0}$  are as in the fixed-effects specification, and  $f_t = -t$  is a time trend. Note that, given the one-dimensional factor structure, the underlying dimension of  $\alpha_{it0}$  is still one in this case.

Finally, in the random coefficients specification,  $\alpha_{it0}$  and  $\mu_{it0}$  are modeled as for fixed-effects probit. In addition, the first component of  $\theta_{it0}$  is constant, equal to 1, and all other components are identical, equal to  $\theta_{it0}^{(r')} = \theta_{i0}^{(r')} = 1 + \frac{\alpha_{i0}}{2}$ . While the number of heterogeneous parameters is increasing relative to the fixed-effects specification, the underlying dimension of  $(\alpha_{it0}, \theta'_{it0})$  is equal to one.

We estimate the fixed-effects probit specification using GFE methods based on the moment vectors  $h_i = (\overline{Y}_i, \overline{X}'_i)'$ , running a probit on group indicators in the second step. When estimating the timevarying specification we take  $h_i = (\overline{Y}_i, \overline{X}'_i)'$  as moments, and interact the group indicators with time dummies in the second step. Lastly, we estimate the random coefficients specification using GFE based on  $h_i = (\overline{Y}_i, \overline{X}'_i, \overline{YX}'_i)'$ , where in the second step we interact the group indicators with all covariates except the first one. We will compare GFE with fixed-effects estimators in two of the three specifications: fixed-effects probit and random coefficients. However, fixed-effects is not feasible in the time-varying specification. For every specification, we will vary the dimension of heterogeneity by varying the number of covariates.

In all simulations we take N = 1000 and T = 20, and report estimates of the first component of  $\theta_{it0}$ , which is always constant and equal to 1 irrespective of the design. We perform 500 simulations, and estimate kmeans using the matlab implementation, which augments Lloyd's updates with a local search rountine, using 100 randomly generated starting values. For iterated GFE, we start the iteration at the two-step estimate and iterate 10 times. To implement conditional kmeans we initialize Algorithm 2 at two-step estimates. In the cubic specification we include all linear, quadratic, and cubic powers of the elements of  $X_{it}$ , as well as a set of interactions, for  $(\dim X_{it} + 1)^2$  coefficients in total. In the neural network specification we use a single hidden layer feedforward model with three nodes, with a sigmoid input link and a linear output link. We train the neural network using matlab's "trainscg". Note that, in this case the algorithm is not guaranteed to decrease the objective function in each step. We stop the algorithm after 10 iterations.
# G.2 Inference

In Table S1 we report the analytical coverage for two-step GFE, iterated GFE, and fixed-effects, both uncorrected and using half-panel jackknife bias correction. To compute confidence intervals we use two methods: an analytical method based on  $\hat{H}$  in (17), and a method based on the bootstrap clustered at the individual level.

Consider first the fixed-effects probit model in Panel A. The results show that coverage accuracy for two-step GFE strongly depends on the dimension of heterogeneity. In DGP 1, where the dimension is one, both analytical and bootstrap methods give close to nominal coverage when using bias correction. This is true both when K = 10 and when we estimate K using our rule (DGP 1<sup>\*</sup>). This finding is in line with Corollary 4. However, in the other DGP the dimension of heterogeneity is larger and two-step GFE is substantially biased, as shown by Table 1. As a result, coverage probabilities are heavily distorted. The intermediate case of DGP 2 is interesting, since under-coverage is somewhat less severe when using our rule to estimate K, especially when using the bootstrap to construct confidence intervals.

In contrast, the iterated GFE estimator, while showing some coverage distortions as the dimension of X grows (even after bias correction), outperforms two-step GFE, and in fact it tends to outperform bias-corrected fixed-effects in terms of coverage when using the bootstrap. According to our simulations, bootstrap inference is more reliable than inference based on the analytical method. It is interesting to see that bootstrap inference based on iterated GFE with bias correction has relatively small finite sample distortions, even in DGP 3 and 4 where the dimension of heterogeneity is such that Corollary 4 does not apply.

Consider next the model with time-varying unobservables in Panel B of Table S1. In DGP 1, bias-corrected two-step GFE tends to undercover. Note that in this case the conditions of Corollary 4 are not met, since there are p = T time-varying parameters, and the DGP features a non-stationary time trend. In DGP 2 to 4 the dimension of heterogeneity is larger, and biases are higher, so the twostep method undercovers heavily. An exception is DGP 2 when we estimate K (DGP 2<sup>\*</sup>), in which case bias-reduced two-step GFE is not substantially biased, and the bootstrap coverage probability is good. In the various DGP, iterated GFE gives less distorted coverage probabilities. Indeed, while the analytical inference method undercovers, bias-reduced iterated GFE tends to be conservative when based on the bootstrap. Lastly, in this time-varying specification, fixed-effects is no longer consistent, hence we do not report inference in this case.

Finally, consider the random coefficients probit model in Panel C. In both DGP 3 and 4, where the dimension of heterogeneity is relatively large, coverage based on two-step GFE is heavily distorted. In such DGP, Corollary 4 is not applicable, since (16) is not met. In contrast, although iterated GFE undercovers substantially when using analytical confidence intervals, the bootstrap provides more accurate coverage. Note that for DGP 4 the bias-corrected fixed-effects estimator undercovers

less than in DGP 3, but in both cases the estimator is substantially biased and the higher coverage probability in DGP 4 only comes from the very large standard deviation of the estimator (see Table 2 in the main text).

### G.3 Additional results with estimated K

We next present results on the performance of two-step estimators when we estimate K. When the dimension of heterogeneity is small, such as in DGP 1 and 2, we use (4) with  $\gamma = 1$  to set K. We have used this method in DGP 1<sup>\*</sup> and 2<sup>\*</sup> in Tables 1, 2, and S1. In DGP with a larger dimension of heterogeneity, such as DGP 3 and 4, conditional methods perform substantially better than two-step GFE, and we use a different rule based on a conditional first step to set K, as described in Appendix B.2. Both rules are based on the first step alone.

In Table S2 we show the results for DGP 1 and 2, using the unconditional rule (4). The results illustrate how the choice of K adapts to the dimension of heterogeneity. A larger dimension d increases approximation error, hence leads to a larger value of K. We also report the estimated K in half-panels, since we compute them as part of the bias correction strategy. In half-panels the estimates of K decrease, since the noise level increases. Notice that, in the time-varying specification, K estimates differ for the two half-panels. This is probably due to the lack of stationarity of the DGP in this case. We commented on mean estimates when discussing the results of Tables 1 and 2 in the main text. In addition, standard deviations do not seem to increase when estimating K.

We turn now to the performance of GFE methods when the choice of K is based on a conditional first step, as described in Appendix B.2. We first compute the kmeans estimator with  $K_{\text{max}} = 30$ . We use this unconditional classification to initialize a linear conditional kmeans algorithm that uses first-order polynomials as a basis, also with  $K_{\text{max}} = 30$ . In the case of the time-varying probit model, we add a linear trend to the set of covariates. This step provides the value of the noise level. We then estimate the conditional kmeans estimator, and we compare the noise level to the value of the objective, for different values of K. The selected K is then the smallest value such that the noise level is larger than the conditional kmeans objective. This choice of K balances the time-series noise versus the approximation error solely due to  $\alpha_{i0}$ .

We show the results in Tables S3 and S4. Notice that the estimated K is now lower than when using our unconditional rule (see Table S2). As an example, in the fixed-effects probit model in DGP 2, K is around 6 in the conditional case, whereas K is around 20 in the unconditional case. In addition, the estimated K values remain stable when the number of covariates increases. In fact, we see a mild decrease in K, which might be due to the nonparametric approximation to the conditional expectation of Y given X: since we keep the number of terms in the series fixed, the nonparametric error increases with the number of covariates, hence driving the noise level up and giving a lower K value.

Tables S3 and S4 show that the performance of GFE methods is quite similar to the case with

K = 10. The case of DGP 2 in Table S3 is interesting, since when using the conditional rule to set K two-step GFE is severely biased in this case. This is due to the fact that the estimated number of groups ( $K \approx 5-6$ ) is small relative to the one given by our unconditional rule in Table S2 ( $K \approx 20$ ). Indeed, the conditional rule targets a uni-dimensional approximation error in this case, but the bias of two-step GFE is affected by a bi-dimensional approximation error. In contrast, we see that GFE methods with a conditional first step – for which the conditional rule is designed – tend to perform very well.

#### G.4 Computational aspects

Finding the global optimum in kmeans is challenging, and most available algorithms only provide approximate solutions. In Table S5 we study the impact of using different numbers of random starting values in the kmeans heuristic on the statistical performance of GFE estimators. We focus on a fixedeffects probit specification with N = 1000 and T = 20, and use K = 10 groups. Let n be a number of starting values. In each of the 500 simulated samples, we compute 100 kmeans partitions starting at n randomly generated values, and the resulting two-step GFE estimators. We then report the total variance across samples and kmeans runs, as well as the between-sample and within-sample variances. The share of within to total variance measures the contribution of numerical error to overall variability of the estimator.

Unsurprisingly, as the number of starting values increases, the share of variance due to numerical error decreases. Interestingly, it is quite small as soon as we use more than 100 starting values. With 100 values, the numerical variability of kmeans increases the variance of two-step GFE by 1.9% in DGP 2, and by 12.3% in DGP 4. With 1000 starting values, the increase is 0.3% and 7.5%, respectively. Although increasing the number of starting values adds to the overall computational burden, off-the-shelf parallelization is available in various kmeans routines.

Finally, in Table S6 we compare the computational cost of fixed-effects, two-step GFE, and iterated GFE estimators, for different cross-sectional sizes, from N = 100 to N = 10000, in a random coefficients probit model with three covariates (DGP 3). We separately record the time used by the kmeans algorithm in the total computation of GFE. We use 100 starting values to compute the kmeans partition in the first step. When the sample size is small, fixed-effects is faster than GFE, due to the overhead cost of starting the computation in kmeans. However, as N increases to moderate sizes, such as N = 1000, GFE becomes faster than fixed-effects. The reason is that kmeans computation only increases moderately with the sample size, and second-step computation involves much fewer parameters to estimate than fixed-effects. Iterating adds some computational cost, and fixed-effects is faster than iterated GFE (where we iterate 10 times) up to N = 1000. However, for N = 5000 and N = 10000, iterated GFE is faster than fixed-effects in this DGP.

# H Complements on the dynamic model of location choice

In this section of the appendix we provide details on our illustration to the dynamic structural model of migration, and we report additional results.

### H.1 Details on computation and estimation

Value functions. Let us denote the integrated value function as:

$$\overline{V}_t(S_{i,t-1}) = \mathbb{E}\left[\max_{j \in \{1,\dots,J\}} V_t(j, S_{i,t-1}) + \xi_{it}(j) \middle| S_{i,t-1}\right].$$

The alternative-specific value functions are:

$$V_t(j, S_{i,t-1}) = \mathbb{E}\left[\rho W_{it}(j) - c_i(j_{i,t-1})\mathbf{1}\{j \neq j_{i,t-1}\} + \beta \overline{V}_t(S_{it}) \middle| j_{it} = j, S_{i,t-1}\right],$$

where  $S_{it} = \left(j, \mathcal{J}_{i,t-1}^{j}, \alpha_i\left(\mathcal{J}_{i,t-1}^{j}\right), c_i\left(\mathcal{J}_{i,t-1}^{j}\right)\right)$  when  $j_{it} = j$ , for  $\mathcal{J}_{i,t-1}^{j} = \mathcal{J}_{i,t-1} \cup \{j\}$ . From the functional forms we obtain (as in Rust, 1994):

$$\overline{V}_t(S_{i,t-1}) = \ln\left(\sum_{j=1}^J \exp V_t(j, S_{i,t-1})\right) + \gamma, \tag{H1}$$

where  $\gamma \approx .57$  is Euler's constant. Moreover:

$$V_{t}(j, S_{i,t-1}) = \mathbb{E}\left[\rho \exp\left(\alpha_{i}(j) + \frac{\sigma^{2}}{2}\right) - c_{i}(j_{i,t-1})\mathbf{1}\{j \neq j_{i,t-1}\} + \beta \overline{V}_{t}\left(j, \mathcal{J}_{i,t-1}^{j}, \alpha_{i}\left(\mathcal{J}_{i,t-1}^{j}\right), c_{i}\left(\mathcal{J}_{i,t-1}^{j}\right)\right) \middle| j_{it} = j, S_{i,t-1}\right]$$
(H2)

where the expectation is taken with respect to the belief distribution, which is the conditional distribution of  $\alpha_i(j)$  given  $\alpha_i(\mathcal{J}_{i,t-1})$  and  $c_i(\mathcal{J}_{i,t-1})$ , conditional on  $j_{i,t-1}$  and  $j_{it} = j$ .

**Computation.** To compute the solution, we first calculate the full-information value functions and then proceed using backward induction. In the case where all locations have been visited,  $\mathcal{J}_{it} = \{1, ..., J\}$  so  $S_{it} = (j_{it}, \{1, ..., J\}, \{\alpha_i(1), ..., \alpha_i(J)\}, \{c_i(1), ..., c_i(J)\})$ . Denote the corresponding integrated value function given the most recent location j as  $\overline{V}^J(i, j)$ . From (H1) and (H2) we have:

$$\overline{V}^{J}(i,j) = \ln\left(\sum_{j'=1}^{J} \exp\left[\rho \exp\left(\alpha_{i}(j') + \frac{\sigma^{2}}{2}\right) - c_{i}(j)\mathbf{1}\{j' \neq j\} + \beta \overline{V}^{J}(i,j')\right]\right) + \gamma, \quad j = 1, \dots, J.$$

We solve this fixed-point system by successive iterations.

Consider now a case where the agent has visited s states in set  $\mathcal{J} \subsetneqq \{1, ..., J\}$ , and is currently at location j. Let  $\overline{V}^{s}(i, j, \mathcal{J})$  denote her integrated value function. The latter solves:

$$\overline{V}^{s}(i,j,\mathcal{J}) = \ln\left(\sum_{j'\notin\mathcal{J}} \exp\left[\mathbb{E}_{\mathcal{J},j,j'}\left(\rho \exp\left(\alpha_{i}(j') + \frac{\sigma^{2}}{2}\right) - c_{i}(j) + \beta \overline{V}^{s+1}(i,j',\mathcal{J}^{j'})\right)\right] + \sum_{j'\in\mathcal{J}} \exp\left[\rho \exp\left(\alpha_{i}(j') + \frac{\sigma^{2}}{2}\right) - c_{i}(j)\mathbf{1}\{j'\neq j\} + \beta \overline{V}^{s}(i,j',\mathcal{J})\right]\right) + \gamma_{s}$$

where  $\mathbb{E}_{\mathcal{J},j,j'}$  is taken with respect to the distribution of  $\alpha_i(j')$  given  $\alpha_i(\mathcal{J})$  and  $c_i(\mathcal{J})$ , conditional on moving from j to j'. In practice we discretize the values of each  $\alpha_i(j)$  on a 50-point grid. In the model with heterogeneous costs, this also implies a discretization of mobility costs given our modeling of costs and returns (see below). In the computation of the fixed points we set a  $10^{-11}$  numerical tolerance.

Estimation. The choice probabilities entering the likelihood are given by an estimated counterpart to (26), where the estimated value functions  $\hat{V}_t(j, j_{i,t-1}, \mathcal{J}_{i,t-1}, \hat{\alpha}(\hat{k}_i, \mathcal{J}_{i,t-1}), c(\hat{k}_i, \mathcal{J}_{i,t-1}), \theta)$  solve the system (H1)-(H2). Notice that, in the model with homogeneous costs, c does not depend on the group or the history of locations. We estimate the conditional expectation in (H2) as a conditional mean given  $\hat{\alpha}(\hat{k}_i, \mathcal{J}_{i,t-1})$  and  $c(\hat{k}_i, \mathcal{J}_{i,t-1})$ , based on all job movers from  $\mathcal{J}_{i,t-1}$  to  $j_{it} = j$ . Nonparametric or semi-parametric methods could be used for this purpose. We use an exponential regression estimator in the illustration, with degree one to compute expected returns, and degree two for expected value functions. We checked these specifications provided a good fit to the conditional expectations.

Iteration. To perform the iteration, we first estimate the idiosyncratic variance of log-wages  $\sigma^2$  as:

$$\widehat{\sigma}^2 = \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T \left( \ln W_{it} - \widehat{\alpha}(\widehat{k}_i, j_{it}) \right)^2.$$
(H3)

Then, individual groups are assigned as:

$$\widehat{k}_{i}^{(2)} = \operatorname{argmax}_{k \in \{1, \dots, K\}} \sum_{t=1}^{T} \sum_{j=1}^{J} \mathbf{1}\{j_{it} = j\} \left( \ln \Pr\left(j_{it} = j \mid j_{i,t-1}, \mathcal{J}_{i,t-1}, \widehat{\alpha}(k, \mathcal{J}_{i,t-1}), \widehat{c}(k, \mathcal{J}_{i,t-1}), \widehat{\theta}\right) \\
+ \ln \phi(\ln W_{it}; \widehat{\alpha}(k, j), \widehat{\sigma}^{2}) \right),$$

where  $\phi$  denotes the normal density. Note that information on both wages and choices is used to reclassify individuals.

Given group assignments, we update parameters as:

$$\widehat{\alpha}^{(2)}(k,j) = \frac{\sum_{i=1}^{N} \sum_{t=1}^{T} \mathbf{1}\{\widehat{k}_{i}^{(2)} = k\} \mathbf{1}\{j_{it} = j\} \ln W_{it}}{\sum_{i=1}^{N} \sum_{t=1}^{T} \mathbf{1}\{\widehat{k}_{i}^{(2)} = k\} \mathbf{1}\{j_{it} = j\}}$$

with an update for  $\sigma^2$  analogous to (H3), and:

$$(\widehat{\theta}^{(2)}, \widehat{c}^{(2)}) = \underset{(\theta,c)}{\operatorname{argmax}} \sum_{i=1}^{N} \sum_{t=1}^{T} \sum_{j=1}^{J} \mathbf{1}\{j_{it} = j\} \ln \Pr\left(j_{it} = j \mid j_{i,t-1}, \mathcal{J}_{i,t-1}, \widehat{\alpha}^{(2)}(\widehat{k}_{i}^{(2)}, \mathcal{J}_{i,t-1}), c(\widehat{k}_{i}^{(2)}, \mathcal{J}_{i,t-1}), \theta\right).$$

This procedure may be iterated further. Note that in the update step we do not maximize the full likelihood as a function of parameters  $\alpha$ , c,  $\sigma^2$ ,  $\theta$ . Rather, we use a partial likelihood estimator by which we first estimate wage parameters  $\alpha$  and  $\sigma^2$ , and then estimate utility and cost parameters  $\theta$  and c. We use this approach for computational reasons; see Rust (1994) and Arcidiacono and Jones (2003) for related approaches.

Specification of cost heterogeneity. In the DGP with heterogeneous costs, we specify  $c_i(j) = \exp(c_0 + c_1\alpha_i(j))$ , where we estimate  $c_0$  and  $c_1$  using an exponential regression of the estimates  $\hat{c}(\hat{k}_i, j)$  on  $\hat{\alpha}(\hat{k}_i, j)$  and a constant, across both groups and regions. We find  $\hat{c}_0 = 0.85$  and  $\hat{c}_1 = 0.32$ . A positive relationship between returns and mobility costs could reflect a link between costs and investments. In the Monte Carlo, GFE recovers  $c_0$  quite well, although  $c_1$  estimates are somewhat biased downwards.

## H.2 Additional results

In this subsection we show additional estimation results for the illustration in Section 5.

**Parameters** a and b. In Figure S1 we report analogous Monte Carlo results as in Figures 4 and 5, for the parameters a and b that govern the probability of being a "mover type". In particular, we see that two-step GFE and its bias-reduced and iterated versions perform quite well.

Fit of the model. The model reproduces well the probability of moving, both unconditionally and conditional on past wages; in particular it reproduces the negative relationship between past wages and mobility. It also reproduces means and variances of log-wages by location. However, the model does not fit well average wages after mobility, as it tends to predict mean wage increases upon job move while the data do not show such a pattern.

**Fixed-***K* **GFE.** Until the end of this subsection we focus on the model with homogeneous costs. We start by reporting results based on fixed values of *K*, from K = 2 to K = 8, in Figure S2. We see that taking K = 2 yields imprecise estimates, in particular for  $\rho$ . In comparison, taking K = 4, K = 6 or K = 8 results in better performance. The most accurate results are obtained taking K = 6 or K = 8 and using bias reduction and one or three iterations. Those results are close to the ones using our method to select *K* (see Figure 4, where the average value for  $\hat{K}$  is 6.4). **Complements: fixed-effects estimates.** In the DGP with homogeneous costs, fixed-effects estimation is computationally tractable. This is due to the fact that the  $\alpha$ 's and the structural parameters can be estimated sequentially. One fixed-effects estimation of the structural parameters is about 2.5 times slower than one estimation of the model with 6 – 7 groups (an average value of  $\hat{K}$ ), although it becomes 5 times slower in a sample with 5 times as many individuals.

**Complements: random-effects estimates.** To compute random-effects estimators based a finite mixture with K = 2, K = 4, and K = 8 types, respectively, we use the EM algorithm of Arcidiacono and Jones (2003), where wage-specific parameters and structural parameters are estimated sequentially in each M-step of the algorithm. Setting a tolerance of  $10^{-6}$  on the change in the likelihood, the algorithm stops after 27, 80, and 339 iterations with K = 2, K = 4, and K = 8 types, respectively. Estimation is substantially more time-consuming than when using two-step grouped fixed-effects.

# I Empirical illustration: firm and worker heterogeneity

In this section of the appendix we present an empirical illustration, where we consider the question of assessing the sources of wage dispersion across workers and firms.

#### I.1 Setup and main results

We consider an additive model in worker and firm heterogeneity:

$$Y_{it} = \eta_i + \psi_{j(i,t)} + \varepsilon_{it},\tag{I1}$$

where  $Y_{it}$  denote log-wages, worker *i* works in firm j(i,t) at time *t*, and  $\eta_i$  and  $\psi_j$  denote unobserved attributes of worker *i* and firm *j*, respectively. Equation (I1) corresponds to the model of Abowd, Kramarz and Margolis (1999) for matched employer-employee data, where we abstract from observed covariates for simplicity. Our interest centers on the decomposition of the variance of log-wages into a worker component, a firm component, a component reflecting the sorting of workers into heterogeneous firms, and an idiosyncratic match component:

$$\operatorname{Var}(y_{i1}) = \operatorname{Var}(\eta_i) + \operatorname{Var}\left(\psi_{j(i,1)}\right) + 2\operatorname{Cov}\left(\eta_i, \psi_{j(i,1)}\right) + \operatorname{Var}(\varepsilon_{i1}).$$
(I2)

Identification of firm effects  $\psi_j$  comes from job movements. As an example, with two time periods the fixed-effects estimators of the  $\psi_j$ 's are obtained from:

$$Y_{i2} - Y_{i1} = \psi_{j(i,2)} - \psi_{j(i,1)} + \varepsilon_{i2} - \varepsilon_{i1},$$

which is uninformative for workers who remain in the same firm in the two periods. When the number of job movers into and out of firm j is low,  $\psi_j$  may be poorly estimated; see Andrews, Gill, Schank and Upward (2008) and Jochmans and Weidner (2016). This source of incidental parameter bias may be particularly severe in short panels.

To alleviate this "low-mobility bias", Bonhomme, Lamadon and Manresa (2019, BLM) propose to reduce the number of firm-specific parameters by grouping firms based on firm-level observables in a first step. Then, in a second step, the  $\psi_j$ 's are recovered at the group level, thus pooling information across job movers within firm groups. Specifically, given a kmeans-based classification  $\{\hat{k}_j\}$  of firms, the  $\psi(k)$ 's are estimated based on the following criterion:

$$\min_{(\psi(1),\dots,\psi(K))} \sum_{i=1}^{n} \left\| \widehat{\mathbb{E}} \left( Y_{i2} - Y_{i1} \,|\, \widehat{k}_{j(i,1)}, \widehat{k}_{j(i,2)} \right) - \psi \left( \widehat{k}_{j(i,2)} \right) + \psi \left( \widehat{k}_{j(i,1)} \right) \right\|^{2},$$

where  $\widehat{\mathbb{E}}$  denotes a group-pair average and *n* denotes the number of workers, subject to a single normalization (e.g.,  $\psi(K) = 0$ ).

This estimator is a two-step GFE estimator based on external moments. Here N is the number of firms, T is the number of available observations to estimate the firm-specific parameters  $\psi_j$  (that is, T is the number of job movers per firm), and S is the number of measurements on firm heterogeneity. In BLM, firms are classified based on their empirical wage distribution functions. Note that a difference with the setting of Theorem 1 is that here the likelihood function is not separable across firms, due to the fact that two firms are linked by the workers who move between them. We conjecture that Theorem 1 could be extended to such network settings, although formally developing this extension exceeds the scope of this paper.

Classifying firms based on wage distribution functions. Here we outline the distributionbased approach we use to classify firms. Let  $Y_{is}$  denote the wage of worker s in firm i. We denote  $\hat{F}_i(y) = \frac{1}{S} \sum_{s=1}^S \mathbb{1}\{Y_{is} \leq y\}$  the empirical cumulative distribution function of  $Y_{is}$ . We classify firms based on  $h_i = \hat{F}_i$ , using the norm  $\|g\|_{\omega}^2 = \int g(w)^2 \omega(w) dw$ , where  $\omega$  is an integrable function. The classification step then is:  $\min_{(\tilde{h}, \{k_i\})} \sum_{i=1}^N \|h_i - \tilde{h}(k_i)\|_{\omega}^2$ , where the  $\tilde{h}(k)$ 's are functions. In practice we discretize the integral, leading to a weighted kmeans objective function.

Simulation exercise. We now present simulation results based on model (I1). We focus on a two-period model, where  $\varepsilon_{it}$  are independent of j(i, 1), j(i, 2),  $\eta$ 's, and  $\psi$ 's, have zero means, and are i.i.d. across workers and time. Following BLM we adopt a correlated random-effects approach to model worker heterogeneity within firms. The parameters of the model are the firm fixed-effects  $\psi_j$ , the means and variances of worker effects in each firm  $\mu_j = \mathbb{E}(\eta_i | j(i, 1) = j)$  and  $\sigma_j^2 = \text{Var}(\eta_i | j(i, 1) = j)$ , and the variance of idiosyncratic errors  $s^2 = \text{Var}(\varepsilon_{i1})$ . We will be estimating the components of the variance decomposition in (I2). In addition we will report estimates of the correlation between worker and firm effects,  $\text{Corr}(\eta_i, \psi_{j(i,1)})$ , which is commonly interpreted as a measure of sorting.

In the baseline DGP firm heterogeneity is continuous and three-dimensional, and its underlying dimension equals one. Specifically, the vector of firm-specific parameters is:

$$\alpha_j = \left(\psi_j, \, \mu_j, \, \sigma_j^2\right) = \left(\psi_j, \, \mathbb{E}\left(\eta_i | \psi_{j(i,1)} = \psi_j\right), \, \operatorname{Var}\left(\eta_i | \psi_{j(i,1)} = \psi_j\right)\right),$$

so all firm-specific parameters are (nonlinear) functions of the scalar firm effects  $\psi_j$ . This specification is consistent with theoretical models of worker-firm sorting where firms are characterized by their scalar productivity level.<sup>2</sup> Below we report simulations using several alternative designs. We study cases where the underlying dimension of firm heterogeneity is equal to two, which allows for a second dimension of latent firm heterogeneity in addition to the wage effects  $\psi_j$  (below we provide evidence suggesting that the underlying dimension of firm heterogeneity is low in the data). We also consider a DGP where firm heterogeneity is discrete in the population.

We start by estimating model (11) on Swedish register data, following BLM. We select male workers full-year employed in 2002 and 2004, and define as job movers the workers whose firm IDs change between the two years. We focus on firms that are present throughout the period. There are about 20,000 job movers in the sample. We use two-step GFE with a classification based on the firms' empirical distributions of log-wages in 2002, evaluated at 20 percentiles of the overall log-wage distribution, with K = 10 groups. In the second step, we estimate the model's parameters  $\hat{\psi}(\hat{k}_j)$ ,  $\hat{\mu}(\hat{k}_j)$ ,  $\hat{\sigma}^2(\hat{k}_j)$ , and  $\hat{s}^2$ .

Given parameter estimates, we then simulate a two-period model where firm heterogeneity is continuously distributed. Specifically, the  $\psi_j$ 's are drawn from a normal distribution, calibrated to match the mean and variance of the  $\widehat{\psi}(\widehat{k}_j)$ 's. We draw 120,000 workers in the cross-section, including 20,000 job movers. We run simulations for different firm sizes, from 10 workers per firm to 200 workers per firm. The total number of job movers is kept constant, so the number of movers per firm increases with firm size.

In Table S9 and Figure S3 we report the mean and 95% confidence intervals of GFE and fixedeffects estimators of the components of the variance decomposition (I2), across 500 simulations. We estimate the number of groups in every simulation. We see that biases of two-step GFE estimators decrease quite rapidly when firm size grows, although biases are not negligible when firms are small. As an example, the variance of firm effects is two thirds of the true value on average when firm size equals 10, and 75% of the true value for a firm size of 20. Moreover, although we do not have a theoretical justification for it in this setting, bias correction using the half-panel jackknife method of Dhaene and Jochmans (2015) tends to provide performance improvements: for example, biases for the variance of firm effects become 25% and 5% for firm sizes of 10 and 20, respectively. Note that here bias correction is not associated with large increases in dispersion. To implement the bias-correction

<sup>&</sup>lt;sup>2</sup>Here Assumption 2 requires the mapping  $\psi_j \mapsto (\psi_j + \mathbb{E}(\eta_i | \psi_j), \text{Var}(\eta_i | \psi_j))$  to be injective. Firm-specific means of log-wages being monotone in  $\psi_j$  is sufficient, though not necessary, for this to hold.

method we select two halves within each firm at random, and we re-estimate the number of groups in each half-sample. In addition, the last column in the table shows that the estimated number of groups is rather small, and close to proportional to the square root of firm size (which is to be expected in this DGP with one-dimensional underlying heterogeneity).

Lastly, in the bottom panel of Table S9 we report the results for a fixed-effects estimator, which is computationally feasible in this linear setting. We see that fixed-effects is substantially biased. This shows that incidental parameter bias due to low mobility is particularly acute in this DGP. The contrast between fixed-effects and GFE is in line with Theorem 1, since here the number T of job movers per firm is small relative to the total number S of workers in the firm which we use to group firms. Hence, GFE, possibly combined with bias reduction, provides an effective regularization in this context.

#### I.2 Complements

Asymptotic properties of the distributional first step. Let  $F_i(y) = \Pr(Y_{is} \le y | \alpha_{i0}) = G(y; \alpha_{i0})$  denote the population cumulative distribution function of  $Y_{is}$ . Similarly to Lemma 1, the following convergence rate holds:

$$\frac{1}{N}\sum_{i=1}^{N}\left\|\widehat{h}(\widehat{k}_{i})-G(\cdot;\alpha_{i0})\right\|_{\omega}^{2}=O_{p}\left(\frac{1}{S}\right)+O_{p}\left(B_{\alpha}(K)\right),$$

provided (i)  $\frac{1}{N} \sum_{i=1}^{N} \|\widehat{F}_i - F_i\|_{\omega}^2 = O_p(S^{-1})$ , and (ii)  $G(\cdot; \alpha_{i0})$  is Lipschitz-continuous in its second argument. Here  $\alpha \mapsto G(\cdot; \alpha)$  maps individual-specific parameters to functions in  $L^2(\omega)$ . In order for (i) to hold, a functional central limit theorem on  $\widehat{F}_i$ , together with  $\omega$  being integrable, will suffice. The Lipschitz condition in (ii) will be satisfied provided  $\int \frac{\partial \ln f(y \mid \alpha_i)}{\partial \alpha} \frac{\partial \ln f(y \mid \alpha_i)}{\partial \alpha'} f(y \mid \alpha_i) dy dx$  is uniformly bounded.

For the second step to deliver estimators with similar properties as in Theorem 1, an injectivity condition is needed. When classifying individuals based on empirical distributions, this condition does not impose further restrictions other than  $\alpha_{i0}$  being identified. Indeed,  $\alpha \mapsto G(\cdot, \alpha)$  being injective is equivalent to  $[G(\cdot, \alpha_2) = G(\cdot, \alpha_1) \Rightarrow \alpha_2 = \alpha_1]$ , which in turn is equivalent to  $\alpha_{i0}$  being identified given knowledge of the function G.

**Details on estimation.** Following BLM we exploit the following restrictions, where we denote as  $m_i = \mathbf{1}\{j(i,1) \neq j(i,2)\}$  the job mobility indicator. For job movers, using the fact that mobility does not depend on  $\varepsilon$ 's, and that  $\varepsilon_{i1}$  is independent of  $\varepsilon_{i2}$ , we have:

$$\mathbb{E}\left(Y_{i2} - Y_{i1} \mid m_i = 1, j(i,1), j(i,2)\right) = \psi_{j(i,2)} - \psi_{j(i,1)},\tag{I3}$$

$$\operatorname{Var}(Y_{i2} - Y_{i1} | m_i = 1, j(i, 1), j(i, 2)) = \operatorname{Var}(\varepsilon_{i2}) + \operatorname{Var}(\varepsilon_{i1}) = 2s^2.$$
(I4)

Then, in the first cross-section we have:

$$\mathbb{E}(Y_{i1} | j(i,1)) = \psi_{j(i,1)} + \mathbb{E}(\eta_i | j(i,1)) = \psi_{j(i,1)} + \mu_{j(i,1)},$$
(I5)

$$\operatorname{Var}(Y_{i1} | j(i,1)) = \operatorname{Var}(\eta_i | j(i,1)) + \operatorname{Var}(\varepsilon_{i1}) = \sigma_{j(i,1)}^2 + s^2.$$
(I6)

In estimation, we first compute a firm partition  $\{\hat{k}_j\}$  into K groups based on firm-specific empirical distributions of log-wages (evaluated at 20 points). In the second step, we use the following algorithm:

- 1. Compute  $\widehat{\psi}(\widehat{k}_i)$  based on sample counterparts to (I3).
- 2. Compute  $\hat{s}^2$  based on (I4).
- 3. Given  $\widehat{\psi}(\widehat{k}_i)$ , compute  $\widehat{\mu}(\widehat{k}_i)$  based on (I5).
- 4. Given  $\hat{s}^2$ , compute  $\hat{\sigma}^2(\hat{k}_j)$  based on (16). In practice we impose non-negativity of the variances using a quadratic programming routine.

Given parameter estimates, we then estimates the variances and covariance in (I2) by aggregation across types.

The fixed-effects estimator in Table S9 is computed following the same algorithm, except that K is taken equal to N. Hence, the estimates of the firm effects  $\psi_j$  correspond to the estimator of Abowd, Kramarz and Margolis (1999). However, instead of relying on a fixed-effects approach on the worker side, in this two-period setting we rely on a correlated random-effects approach to deal with worker heterogeneity. In that specification, the mean and variance of worker effects  $\eta_i$  are firm-specific. We compute the connected set in an initial step, and use sparse matrix coding for efficient computation.

Monte Carlo designs. We consider four additional DGP, in addition to DGP 1 reported in Table S9. In Table S7 we show the sample sizes that we use in all designs, including the average number of job movers per firm. DGP 2 has one-dimensional underlying heterogeneity, with different parameter values: the variance of firm effects is larger than in DGP 1, while the correlation between firm effects and worker effects is smaller, the relative magnitudes being close to the estimates of Card, Heining and Kline (2013). DGP 3 and DGP 4 have two-dimensional underlying heterogeneity  $(\psi_j, V_j)$ , where  $\psi_j$  is the wage firm effect and  $V_j$  drives workers' firm choice.  $(\psi_j, V_j)$  are drawn from a bivariate normal distribution, and the mean and variance of worker effects in the firm are set to  $\mu_j = V_j$  and  $\sigma_j^2 = (a + bV_j)^2$  for some constants a, b which are calibrated to the Swedish sample. We interpret  $V_j$  as a present value driving workers' mobility decisions across firms, which may be only imperfectly correlated with  $\psi_j$  in the presence of non-pecuniary attributes valued by workers. As displayed in Table S8, the two-dimensional DGP differ in terms of parameter values. The last row of the table shows the correlation between the wage firm effect  $\psi_j$  and the present value  $V_j$  in all DGP. Lastly, DGP 5 has discrete heterogeneity. Specifically, there are  $K^* = 10$  "true" groups in the population. The groups are chosen by approximating the firm heterogeneity of DGP 1.

Alternative DGP with one-dimensional heterogeneity: results. In Table S10 we report the results of two-step GFE and its bias-corrected version, as well as fixed-effects, in DGP 2 with one-dimensional heterogeneity and a larger variance of firm effects than in Table S9. The performance of the estimators is comparable to Table S9.

**Bias-corrected fixed-effects.** In Table S11 we report the results of bias-corrected fixed-effects estimation in DGP 1 (top panel, see Table S9) and DGP 2 (bottom panel). In order to implement the bias correction we use the half-panel jackknife of Dhaene and Jochmans (2015), splitting all workers in every firm into two random halves, including job movers. We see that, although bias correction improves relative to fixed-effects, the bias-corrected estimator is still substantially biased, even for moderately large firms. Notice that some of the variance estimates are in fact negative. This is due to the fact that the additive bias correction method does not enforce non-negativity.

Inferring the underlying dimension of firm heterogeneity. As a motivation for considering DGP with an underlying dimension higher than one, but still relatively low, we now attempt to learn the underlying dimension of firm heterogeneity on the Swedish matched employer-employee data set. In statistics, the literature on manifold learning aims at inferring the low intrinsic dimension of large dimensional data; see for example Levina and Bickel (2004) and Raginsky and Lazebnik (2005). Motivated by our method for selecting the number of groups, the method we use here consists in comparing the length of the panel S with the number of groups  $\hat{K}$  estimated from (4). If the underlying dimension of  $\varphi(\xi_{i0})$  is d > 0, then we expect  $\hat{Q}(K)$  to decrease at a rate  $O_p(K^{-\frac{2}{d}}) + o_p(S^{-1})$ . This suggests that  $\hat{K}$  and  $S^{\frac{d}{2}}$  will have a similar order of magnitude. In such a case the underlying dimension may be inferred by plotting the relationship, for panels of different lengths, between  $\ln \hat{K}$ and  $\ln S$ , the slope of which is  $\hat{d}/2$ .

In Figure S4 we report the results of this exercise, taking firms with more than 50 employees, and then randomly selecting x% in each firm, where x varies between 5 and 100. The left graph shows the shape of the objective function  $\hat{Q}(K)$  as a function of K, in logs. In each sample the estimated number of groups  $\hat{K}$  lies at the intersection of that curve and the horizontal line  $\ln(\hat{V}_h)$ . We use a slight modification of the  $\hat{V}_h$  formula to deal with the fact that here the "panel" is unbalanced, since different firms may have different sizes. On the right graph we then plot  $\ln \hat{K}$  against the logarithm of the average firm size in each sample (here we report results based on empirical cdfs of log-wages evaluated at 20 points, but we checked that using 40 points instead did not affect the results). We see that the relationship is approximately linear and the slope is close to one, suggesting that the underlying dimension is around  $\hat{d} = 2$ .

**Two-dimensional heterogeneity: results.** In Table S12 we report the simulation results for DGP 3 with continuous two-dimensional firm heterogeneity. The results for DGP 4, with a smaller

variance of firm effects, are reported in Table S14. The results are shown graphically in Figures S5 and S6. Focusing on the first panel, which corresponds to our recommended choice for the selection rule of the number of groups (that is, taking  $\gamma = 1$  in (4)), we see that the two-step estimators show larger biases than in the one-dimensional case, especially for the variance of firm effects and the correlation parameter. Moreover, bias correction does not succeed at reducing the bias substantially. This suggests that, for the selected number of groups, the approximation error is still substantial. At the same time, as shown by the two bottom panels of the tables, taking  $\gamma = .5$  and  $\gamma = .25$  improves the performance of the two-step estimator. Notice that while the selected number of groups  $\hat{K}$  is monotone in firm size for  $\gamma = 1$  and  $\gamma = .5$ , it is not monotone for  $\gamma = .25$ . This is a finite sample issue: when taking  $\gamma = .25$  and focusing on large firms the number of groups is no longer negligible with respect to the number of firms in the sample. Lastly, performance is further improved when using the bias-reduced estimator.

As pointed out in Section 4, features of the model may be exploited to improve the classification. In the two-dimensional designs DGP 3 and DGP 4, we perform the following moment-based iteration. The two-step method delivers estimates of the mean and variance of worker effects  $\eta_i$  in firm group  $\hat{k}_j$ :  $\hat{\mu}(\hat{k}_j)$  and  $\hat{\sigma}^2(\hat{k}_j)$ , respectively. Regressing  $\sqrt{\hat{\sigma}^2(\hat{k}_j)}$  on  $\hat{\mu}(\hat{k}_j)$  and a constant then gives estimates  $\hat{b}$  and  $\hat{a}$ . Given those, we construct the (iterated) moments:

$$h_{1j} = \widehat{\mathbb{E}}(Y_{i1} \mid j) - \frac{\sqrt{\widehat{\operatorname{Var}}(Y_{i1} \mid j) - \widehat{s}^2} - \widehat{a}}{\widehat{b}}, \ h_{2j} = \frac{\sqrt{\widehat{\operatorname{Var}}(Y_{i1} \mid j) - \widehat{s}^2} - \widehat{a}}{\widehat{b}}$$

where  $\widehat{\mathbb{E}}$  and  $\widehat{\text{Var}}$  denote firm-specific means and variances. Those moments will be consistent for  $\psi_j$ and  $V_j$ , respectively, as S tends to infinity. We then apply two-step GFE to the moments  $h_{1j}$  and  $h_{2j}$ . In Tables S14 and S15 we report the results for the iterated estimator (only iterated once) and its bias-corrected version, for DGP 3 and DGP 4, respectively. We see that the iteration improves performance substantially for DGP 3, although it has small effects on performance in DGP 4.

Low mobility bias and regularization. As shown by Theorem 1, a benefit of discretizing unobserved heterogeneity is that it can reduce the incidental parameter bias of fixed-effects estimators. In the illustration on matched employer-employee data, fixed-effects estimators may be biased due to low rates of worker mobility between firms. In order to assess the impact of mobility rates on the performance of fixed-effects and GFE estimators, in Figures S7 and S8 we report the results of the estimated variance decomposition on 500 simulations, comparing fixed-effects, bias-corrected fixed-effects, two-step GFE with bias correction, and iterated two-step GFE with bias correction. We perform simulations for different number of job movers per firm, from 2 to 10 (shown on the x-axis), and a fixed firm size of 50. The two figures show the results for the two-dimensional DGP: DGP 3 and DGP 4, respectively. We see a striking difference between fixed-effects and GFE: while the former is very sensitive to the number of job movers, the latter is virtually insensitive. In particular, for low

numbers of job movers fixed-effects and its bias-corrected counterpart are severely biased, while the biases of GFE remain moderate. This is in line with Theorem 1. It is worth noting that the average number of job movers per firm is around 0.5 in the original Swedish sample. This suggests that, at least in short panels, the discrete regularization achieved in GFE may result in practical improvements relative to fixed-effects in data sets of realistic dimensions.

**Discrete heterogeneity: results.** Finally, in Table S16 we report results for a discrete DGP (DGP 5) where all firm population parameters are constant within groups  $\hat{k}_j$ , with  $K^* = 10$ . In this case the results of two-step GFE with  $K = K^*$  turn out to be quite similar to those obtained in the continuous DGP in Table S9. However, as the last column in the table shows, in this discrete DGP misclassification frequencies are sizable: 69% misclassification when firm size equals 10, and still 23% when size is 100. We computed misclassification frequencies by solving a linear assignment problem using the simplex algorithm in every simulation. This suggests that, for this DGP, an "oracle" asymptotic theory based on the premise that group misclassification is absent in the limit may not provide reliable guidance for finite sample inference, even when the true number of groups is known. Lastly, the table shows some evidence that bias correction (where the number of groups is estimated in every simulation) improves the performance of the estimator in this setting too.

# References

- Abowd, J., F. Kramarz, and D. Margolis (1999): "High Wage Workers and High Wage Firms", Econometrica, 67(2), 251–333.
- [2] Andrews, M. J., L. Gill, T. Schank, and R. Upward (2008): "High Wage Workers and Low Wage Firms: Negative Assortative Matching or Limited Mobility Bias?" *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 171(3), 673–697.
- [3] Arcidiacono, P., and J. B. Jones (2003): 'Finite Mixture Distributions, Sequential Likelihood and the EM Algorithm", *Econometrica*, 71(3), 933–946.
- [4] Arellano, M., and J. Hahn (2007): "Understanding Bias in Nonlinear Panel Models: Some Recent Developments,". In: R. Blundell, W. Newey, and T. Persson (eds.): Advances in Economics and Econometrics, Ninth World Congress, Cambridge University Press.
- [5] Arellano, M., and J. Hahn (2016): "A likelihood-Based Approximate Solution to the Incidental Parameter Problem in Dynamic Nonlinear Models with Multiple Effects," *Global Economic Review*, 45(3), 251–274.
- [6] Banna, M., F. Merlevède, and P. Youssef (2016): "Bernstein-Type Inequality for a Class of Dependent Random Matrices," *Random Matrices: Theory and Applications*, 5(02), 1650006.
- Bonhomme, S., and E. Manresa (2015): "Grouped Patterns of Heterogeneity in Panel Data," Econometrica, 83(3), 1147–1184.
- [8] Bonhomme, S., T. Lamadon, and E. Manresa (2019): "A Distributional Framework for Matched Employer-Employee Data," *Econometrica*, 87(3), 699–739.
- Bryant, P. and Williamson, J. A. (1978): "Asymptotic Behaviour of Classification Maximum Likelihood Estimates," *Biometrika*, 65, 273–281.
- [10] Card, D., J. Heining, and P. Kline (2013): "Workplace Heterogeneity and the Rise of West German Wage Inequality," *Quarterly Journal of Economics*, 128(3), 967–1015.
- [11] Dhaene, G. and K. Jochmans (2015): "Split Panel Jackknife Estimation," Review of Economic Studies, 82(3), 991–1030.
- [12] Hahn, J., and G. Kuersteiner (2011): "Bias Reduction for Dynamic Nonlinear Panel Models with Fixed Effects," *Econometric Theory*, 27(6), 1152–1191.
- [13] Hahn, J., and H. Moon (2010): "Panel Data Models with Finite Number of Multiple Equilibria," *Econometric Theory*, 26(3), 863–881.

- [14] Hahn, J. and W.K. Newey (2004): "Jackknife and Analytical Bias Reduction for Nonlinear Panel Models", *Econometrica*, 72, 1295–1319.
- [15] Hsu, D., S.M. Kakade, and T. Zhang (2012): "A Tail Inequality for Quadratic Forms of Subgaussian Random Vectors," *Electron. Commun. Probab.*, 17, 52, 1–6.
- [16] Jochmans, K., and M. Weidner (2016): "Fixed-Effect Regressions on Network Data", arXiv preprint arXiv:1608.01532.
- [17] Levina, E., and P. J. Bickel (2004): "Maximum Likelihood Estimation of Intrinsic Dimension," Advances in neural information processing systems, 777–784.
- [18] Newey, W. K. (1997): "Convergence Rates and Asymptotic Normality for Series Estimators," *Journal of Econometrics*, 79(1), 147–168.
- [19] Pollard, D. (1981): "Strong Consistency of K-means Clustering," Annals of Statistics, 9, 135–140.
- [20] Pollard, D. (1982): "A Central Limit Theorem for K-Means Clustering," Annals of Probability, 10, 919–926.
- [21] Raginsky, M., and Lazebnik, S. (2005): "Estimation of Intrinsic Dimensionality Using High-Rate Vector Quantization," In Advances in neural information processing systems, 1105–1112.
- [22] Rust, J. (1994): "Structural Estimation of Markov Decision Processes," Handbook of econometrics, 4(4), 3081–3143.
- [23] Tropp, J. A. (2012): "User-Friendly Tail Bounds for Sums of Random Matrices," Foundations of Computational Mathematics, 12(4), 389–434.
- [24] Vershynin, R. (2010): Introduction to the Non-Asymptotic Analysis of Random Matrices, in Y.
   C. Eldar and G. Kutyniok, ed., Compressed Sensing: Theory and Applications. Cambridge University Press.

	A. Probit with fixed effects											
	Two-step				Iterated				Fixed Effects			
	analy	ytical	boot	strap	analytical bootstrap		analy	rtical	boot	strap		
DGP	UC	BC	UC	BC	UC	BC	UC	BC	UC	BC	UC	BC
1	0.874	0.914	0.870	0.944	0.010	0.898	0.014	0.934	0.010	0.906	0.012	0.902
1*	0.864	0.948	0.878	0.944	0.006	0.932	0.008	0.952	0.010	0.906	0.012	0.902
2	0.000	0.000	0.000	0.000	0.022	0.858	0.030	0.952	0.010	0.920	0.014	0.922
$2^{*}$	0.092	0.544	0.142	0.852	0.006	0.866	0.020	0.930	0.010	0.920	0.014	0.922
3	0.000	0.000	0.000	0.000	0.002	0.700	0.004	0.884	0.000	0.848	0.002	0.858
4	0.000	0.000	0.000	0.000	0.000	0.634	0.008	0.816	0.002	0.758	0.002	0.796

Table S1: Coverage probabilities

		B. Probit with time-varying unobservables										
		Two-step				Iterated				Fixed Effects		
	analy	ytical	boot	strap	analy	ytical	boot	strap	analyt	cical	bootst	rap
DGP	UC	BC	UC	BC	UC	BC	UC	BC	UC	BC	UC	BC
1	0.896	0.846	0.732	0.876	0.000	0.518	0.000	1.000	_	_	_	_
$1^*$	0.876	0.730	0.680	0.670	0.000	0.226	0.000	0.952	_	_	_	_
2	0.000	0.000	0.000	0.000	0.018	0.268	0.074	0.994	_	_	_	_
$2^{*}$	0.002	0.234	0.082	0.970	0.000	0.030	0.000	0.520	_	_	_	_
3	0.000	0.000	0.000	0.000	0.000	0.326	0.000	1.000	_	—	_	_
4	0.000	0.000	0.000	0.000	0.000	0.388	0.000	1.000	_	_	_	_

	C. Probit with random coefficients											
	Two-step Iterated							Fixed	Effects			
	analy	vtical	boot	strap	analy	vtical	boot	strap	analy	vtical	boot	strap
DGP	UC	BC	UC	BC	UC	BC	UC	BC	UC	BC	UC	BC
3	0.000	0.000	0.000	0.000	0.002	0.398	0.002	0.876	0.000	0.160	0.000	0.464
4	0.000	0.000	0.000	0.000	0.020	0.380	0.042	0.928	0.000	0.172	0.000	0.812

Notes: See notes to Tables 1 and 2. The analytical coverage is based on an estimate of the asymptotic variance matrix. The bootstrap is clustered at the individual level, computed using 100 bootstrap simulations. Coverage probabilities for a nominal level of 95%. Results of 500 simulations.

		А	. Probit with fix	xed effects					
		$\widehat{K}$	Two	-step	Iter	ated			
DGP	full	half	UC	BC	UC	BC			
1	11.8 (0.41)	8.45 (0.50)	0.985 (0.019)	1.002 (0.019)	1.087 (0.021)	0.999(0.021)			
2	$20.31 \ (0.69)$	12.48(0.50)	0.940 (0.018)	0.970(0.024)	$1.088 \ (0.022)$	$0.998\ (0.022)$			
			Conditional M	nal Methods					
		Linear	Cu	bic	Neural	network			
DGP	UC	BC	UC	BC	UC	BC			
1	1.043 (0.019)	1.002 (0.021)	1.059 (0.020)	1.014 (0.022)	1.079 (0.020)	1.027(0.024)			
2	$1.001 \ (0.020)$	$0.971 \ (0.024)$	$1.017 \ (0.021)$	$0.991 \ (0.026)$	$1.035\ (0.023)$	$1.016\ (0.029)$			
		B Probit	with time-vary	ing heterogenei	tv				
		$\hat{v}$	T T	, ing neterogener	т,	4 1			
DCP	full	n half	I WO	-step BC	UC	BC			
1	11.02 (0.00)			1.000 (0.005)		1.002 (0.070)			
1	11.92 (0.38)	12.0(0.4), 8.0(0.2)	1.015 (0.021)	1.028 (0.025)	1.319 (0.039)	1.093 (0.070)			
2	21.69(0.78)	20.5 (0.7), 11.1 (0.4)	0.889(0.023)	0.935(0.037)	1.278 (0.044)	1.196(0.080)			
			Conditional M	lethods					
		Linear	Cu	bic	Neural network				
DGP	UC	BC	UC	BC	UC	BC			
1	1.110 (0.023)	1.027 (0.032)	1.186 (0.027)	1.087 (0.040)	1.195(0.029)	1.143(0.045)			
2	$1.041 \ (0.028)$	$1.056\ (0.046)$	$1.096\ (0.035)$	1.098(0.060)	$1.085\ (0.035)$	1.144(0.062)			

# Table S2: Estimated K, unconditional rule

Notes: See notes to Tables 1 and 2. We estimate K using (4), with  $\gamma = 1$ .

			A. Probit wit	h fixed effects			
	ĺ	Ŕ	Two	-step	Iterated		
DGP	full	half	UC	BC	UC	BC	
1	4.98 (0.13)	4.00 (0.00)	0.999 (0.018)	1.013 (0.019)	1.094 (0.021)	1.008 (0.022)	
2	5.90(0.30)	5.00(0.06)	0.744(0.034)	0.774(0.066)	1.070(0.022)	$1.001 \ (0.025)$	
3	5.00(0.00)	4.83(0.38)	0.757(0.021)	0.763(0.028)	1.078(0.026)	0.969(0.033)	
4	4.02(0.15)	4.00(0.00)	0.749(0.022)	$0.754\ (0.028)$	1.069(0.029)	$0.949\ (0.038)$	
			Conditiona	al Methods			
	Lin	near	Cu	lbic	Neural network		
DGP	UC	BC	UC	BC	UC	BC	
1	1.041 (0.019)	0.999 (0.020)	1.066 (0.020)	1.011 (0.022)	1.080 (0.021)	1.030 (0.024)	
2	$0.944 \ (0.019)$	0.910(0.024)	1.005(0.020)	$0.980 \ (0.025)$	$1.021 \ (0.022)$	1.003(0.028)	
3	0.923(0.023)	0.903(0.035)	0.924(0.026)	0.923 (0.040)	1.004(0.027)	0.992(0.044)	

Table S3: Estimated K, conditional rule

Notes: See notes to Table 1. We implement the choice of K described in Appendix B.2, based on a linear conditional kmeans specification, and  $K_{\text{max}} = 30$ .

0.897(0.041)

0.996(0.034)

0.982(0.059)

0.893(0.026)

4

0.894(0.024)

0.875(0.038)

		B. Probit with time-varying heterogeneity								
		$\widehat{K}$	Two	-step	Iter	ated				
DGP	full	half	UC	BC	UC	BC				
1	5.04(0.19)	4.5 (0.5), 4.0 (0.1)	1.029 (0.019)	$1.020 \ (0.025)$	1.191 (0.031)	1.119(0.055)				
2	7.25(0.51)	6.0 (0.4), 5.0 (0.1)	$0.631 \ (0.044)$	0.692(0.086)	1.082(0.038)	1.159(0.070)				
3	5.05(0.21)	5.0(0.1), 4.0(0.2)	$0.562 \ (0.027)$	$0.547 \ (0.039)$	1.082(0.040)	$1.060 \ (0.062)$				
4	4.17(0.38)	4.3 (0.5), 3.1 (0.3)	0.559(0.025)	$0.542 \ (0.035)$	1.046(0.042)	$1.026\ (0.070)$				
			Conditional I	Methods						
	]	Linear	Cu	bic	Neural	network				
DGP	UC	BC	UC	BC	UC	BC				
1	1.065(0.020)	1.012(0.0269)	1.125 (0.022)	1.079(0.034)	1.136 (0.027)	1.092(0.046)				
2	$0.906\ (0.028)$	0.922(0.046)	$0.955\ (0.031)$	$0.956\ (0.057)$	$0.979\ (0.035)$	1.020(0.063)				
3	$0.845\ (0.035)$	$0.851 \ (0.064)$	0.760(0.045)	$0.771 \ (0.080)$	0.948(0.041)	0.994(0.083)				
4	0.802(0.032)	$0.834\ (0.056)$	$0.753 \ (0.037)$	0.782(0.066)	$0.923 \ (0.054)$	0.981 (0.104)				
		C. P	robit with rand	om coefficients						
			Two	-sten	Iter	ated				
DGP	full	half	UC	BC	UC	BC				
3	5.08 (0.28)	4.86 (0.40)	0.690 (0.023)	0.699 (0.028)	1.035(0.025)	0.949 (0.037)				
4	4.31(0.46)	4.11(0.32)	$0.594\ (0.029)$	$0.602 \ (0.034)$	$0.976\ (0.034)$	$0.917 \ (0.058)$				
			Conditional I	Methods						
		Linear	Cu	bic	Neural	network				
DGP	UC	BC	UC	BC	UC	BC				
3	0.879 (0.028)	0.864 (0.038)	0.909 (0.026	0.889 (0.039)	0.982 (0.025)	0.947 (0.040)				
4	0.765(0.032)	0.756 (0.049)	0.798(0.045)	0.805 (0.073)	0.923 (0.044)	0.893 (0.084)				

Table S4: Estimated K, conditional rule (cont.)

Notes: See notes to Table 2. We implement the choice of K described in Appendix B.2, based on a linear conditional kmeans specification, and  $K_{\text{max}} = 30$ .

	Table 5	5. Starti	ing value	s in kine	ans			
				DGP 1				
$\sharp$ starting values	1	5	10	50	100	500	1000	
within	0.027	0.007	0.004	0.002	0.001	0.000	0.000	
between	3.110	3.118	3.121	3.127	3.130	3.132	3.133	
total	3.137	3.125	3.125	3.128	3.131	3.133	3.134	
within (% total)	0.9%	0.2%	0.1%	0.0%	0.0%	0.0%	0.0%	
	DGP 2							
$\sharp$ starting values	1	5	10	50	100	500	1000	
within	1.730	0.634	0.384	0.121	0.068	0.018	0.010	
between	2.815	3.192	3.387	3.568	3.589	3.603	3.606	
total	4.522	3.814	3.761	3.681	3.650	3.614	3.609	
within (% total)	38.3%	16.6%	10.2%	3.3%	1.9%	0.5%	0.3%	
				DGP 3				
# starting values	1	5	10	50	100	500	1000	
$\sharp$ starting values within	1 0.894	$5 \\ 0.675$	$\begin{array}{c} 10\\ 0.581 \end{array}$	$50 \\ 0.354$	$100 \\ 0.281$	$500 \\ 0.148$	$1000 \\ 0.112$	
<pre># starting values within between</pre>	1 0.894 2.705	$5 \\ 0.675 \\ 2.866$	$10 \\ 0.581 \\ 2.959$	$50 \\ 0.354 \\ 3.185$	100 0.281 3.271	$500 \\ 0.148 \\ 3.450$	1000 0.112 3.486	
<pre># starting values within between total</pre>	1 0.894 2.705 3.585	5 0.675 2.866 3.529	10 0.581 2.959 3.529	50 0.354 3.185 3.529	100 0.281 3.271 3.543	500 0.148 3.450 3.589	$   \begin{array}{r}     1000 \\     0.112 \\     3.486 \\     3.590   \end{array} $	
<pre># starting values within between total within (% total)</pre>	$     1 \\     0.894 \\     2.705 \\     3.585 \\     24.9\% $	5 0.675 2.866 3.529 19.1%	$10 \\ 0.581 \\ 2.959 \\ 3.529 \\ 16.5\%$	50 0.354 3.185 3.529 10.0%	100 0.281 3.271 3.543 7.9%	500 0.148 3.450 3.589 4.1%	1000 0.112 3.486 3.590 3.1%	
<pre># starting values within between total within (% total)</pre>	$     1 \\     0.894 \\     2.705 \\     3.585 \\     24.9\% $	5 0.675 2.866 3.529 19.1%	$10 \\ 0.581 \\ 2.959 \\ 3.529 \\ 16.5\%$	50 0.354 3.185 3.529 10.0% DGP 4	100 0.281 3.271 3.543 7.9%	500 0.148 3.450 3.589 4.1%	1000 0.112 3.486 3.590 3.1%	
<pre># starting values within between total within (% total) # starting values</pre>	1 0.894 2.705 3.585 24.9%	5 0.675 2.866 3.529 19.1%	10 0.581 2.959 3.529 16.5%	50 0.354 3.185 3.529 10.0% DGP 4 50	100 0.281 3.271 3.543 7.9%	500 0.148 3.450 3.589 4.1% 500	1000 0.112 3.486 3.590 3.1% 1000	
<pre># starting values within between total within (% total) # starting values within</pre>	1 0.894 2.705 3.585 24.9% 1 1.023	5 0.675 2.866 3.529 19.1% 5 0.847	10 0.581 2.959 3.529 16.5% 10 0.781	50 0.354 3.185 3.529 10.0% DGP 4 50 0.621	100 0.281 3.271 3.543 7.9% 100 0.549	500 0.148 3.450 3.589 4.1% 500 0.395	1000 0.112 3.486 3.590 3.1% 1000 0.338	
<pre># starting values within between total within (% total) # starting values within between</pre>	1 0.894 2.705 3.585 24.9% 1 1.023 3.553	5 0.675 2.866 3.529 19.1% 5 0.847 3.618	10 0.581 2.959 3.529 16.5% 10 0.781 3.689	50 0.354 3.185 3.529 10.0% DGP 4 50 0.621 3.853	100 0.281 3.271 3.543 7.9% 100 0.549 3.909	500 0.148 3.450 3.589 4.1% 500 0.395 4.090	1000 0.112 3.486 3.590 3.1% 1000 0.338 4.158	
<pre># starting values within between total within (% total) # starting values within between total</pre>	$ \begin{array}{c} 1\\ 0.894\\ 2.705\\ 3.585\\ 24.9\%\\ \hline 1\\ 1.023\\ 3.553\\ 4.560\\ \end{array} $	5 0.675 2.866 3.529 19.1% 5 0.847 3.618 4.449	10 0.581 2.959 3.529 16.5% 10 0.781 3.689 4.455	50 0.354 3.185 3.529 10.0% DGP 4 50 0.621 3.853 4.460	100 0.281 3.271 3.543 7.9% 100 0.549 3.909 4.445	500 0.148 3.450 3.589 4.1% 500 0.395 4.090 4.473	1000 0.112 3.486 3.590 3.1% 1000 0.338 4.158 4.485	

Notes: See notes to Table 1. Let n be a number of starting values. In each of the 500 simulated samples, we compute 100 kmeans partitions starting at n randomly generated values and the resulting two-step GFE estimators. We then report the total variance across samples and kmeans runs, as well as the between-sample and within-sample variances. All variances are multiplied by 10000.

Ν	100	500	1000	5000	10000
Kmeans	0.11	0.12	0.16	0.66	1.51
Two-step $\operatorname{GFE}$	0.12	0.15	0.21	1.11	2.71
Iterated GFE	0.28	0.58	0.94	7.07	20.1
Fixed-effects	0.02	0.16	0.50	11.7	93.5

0.02 0.16

Table S6: Computation Time (in seconds)

Notes: Computation time for various estimators in the random coefficients probit model of DGP 3. See notes to Table 2. We use sparse matrix computation to compute the fixed-effects estimator, and we use parallel computing for kmeans. Computations are performed on a 20 CPU core machine.

0.50

93.5

Firm size	Number firms	Number job movers
		per firm
10	10000	2
20	5000	4
50	2000	10
100	1000	20
200	500	40

Table S7: Firms and workers, sample sizes

Notes: Sample sizes for different firm sizes, all DGP.

	small	$Var(\psi)$	large	$Var(\psi)$
	1D	2D	1D	2D
$Var(\psi)$	$0.0017 \\ 2.0\%$	$0.0017 \\ 2.0\%$	0.0204 21.2%	$0.0204 \\ 21.2\%$
$Var(\eta)$	$0.0758 \\ 85.2\%$	$0.0758 \\ 85.2\%$	$0.0660 \\ 68.4\%$	$0.0660 \\ 68.4\%$
$2Cov(\psi,\eta)$	$0.0057 \\ 12.8\%$	$0.0057 \\ 12.8\%$	$0.0050 \\ 10.4\%$	$0.0050 \\ 10.4\%$
$\begin{array}{c} Corr(\psi,\eta) \\ Var(\varepsilon) \end{array}$	$0.4963 \\ 0.0341$	$0.4963 \\ 0.0341$	$0.1373 \\ 0.0341$	$\begin{array}{c} 0.1373 \\ 0.0341 \end{array}$
$Corr(V,\psi)$	1.0000	0.7130	1.0000	0.2540

Table S8: Firms and workers, different DGP

Notes: The four columns show the parameter values and overall shares of variance in DGP 1, DGP 4, DGP 2, and DGP 3, respectively.

Firm size	$\operatorname{Var}\left(\eta_{i}\right)$	$\operatorname{Var}\left(\psi_{j}\right)$	$\operatorname{Cov}\left(\eta_{i},\psi_{j}\right)$	$\operatorname{Corr}\left(\eta_{i},\psi_{j}\right)$	$\operatorname{Var}\left(\varepsilon_{i1}\right)$	$\hat{K}$
			true va	lues		
-	0.0758	0.0017	0.0057	0.4963	0.0341	
			two-step es	stimator		
10	0.0775	0.0011	0.0048	0 5281	0.0348	3.0
10	[0.076, 0.079]	[0.001, 0.001]	[0.005, 0.005]	[0.519, 0.537]	[0.034, 0.035]	[3,3]
20	0.0769	0.0013	0.0051	0.5091	0.0345	4.0
	$\left[0.076, 0.078 ight]$	[0.001, 0.002]	[0.005, 0.005]	$\left[0.500, 0.518 ight]$	[0.034, 0.035]	[4, 4]
50	0.0764	0.0015	0.0054	0.4986	0.0343	6.0
	[0.075, 0.078]	[0.001,0.002]	0.005,0.006	[0.490,0.507]	[0.034,0.035]	[6,6]
100	0.0761	0.0016	0.0055	0.4955	0.0342	8.4
000	[0.075, 0.077]	[0.001,0.002]		[0.487, 0.504]	[0.034,0.035]	[8,9]
200	0.0760	0.0017		0.4930	0.0342	11.3
	[0.075,0.077]	[0.001,0.002]	[0.003,0.000]	[0.403,0.003]	[0.034,0.033]	[11,12]
		two-	step estimator	, bias-corrected	1	
10	0.0778	0.0013	0.0047	0.4511	0.0346	
	$\left[0.076, 0.079 ight]$	[0.001, 0.002]	[0.004, 0.005]	[0.439, 0.463]	[0.034, 0.035]	
20	0.0763	0.0016	0.0055	0.4902	0.0343	
	$\left[0.075, 0.078 ight]$	[0.001, 0.002]	$\left[0.005, 0.006 ight]$	$\left[0.479, 0.501 ight]$	$\left[0.034,\! 0.035 ight]$	
50	0.0762	0.0017	0.0055	0.4876	0.0342	
	[0.075, 0.078]	[0.001,0.002]	0.005,0.006	[0.476,0.499]	[0.034,0.035]	
100	0.0759	0.0017	0.0056	0.4923	0.0341	
000	[0.075,0.077]	[0.002, 0.002]	[0.005,0.006]	[0.481, 0.502]	[0.034,0.035]	
200	0.0759	1100.0	0.0056	0.4909	0.0341	
	[0.074,0.077]	[0.002,0.002]	[0.003,0.000]	[0.480,0.505]	[0.035,0.055]	
			fixed-effects	estimator		
10	0.1342	0.0342	-0.0267	-0.3949	0.0173	
	[0.132, 0.136]	$\left[0.033, 0.036 ight]$	[-0.028, -0.025]	[-0.409, -0.382]	[0.017, 0.018]	
20	0.1002	0.0130	-0.0056	-0.1548	0.0256	
	$\left[0.099, 0.102 ight]$	$\left[0.012, 0.014 ight]$	[-0.006, -0.005]	[-0.169, -0.139]	$\left[0.025, 0.026 ight]$	
50	0.0848	0.0055	0.0019	0.0895	0.0307	
	[0.083, 0.086]	[0.005, 0.006]	[0.002, 0.002]	[0.072, 0.107]	[0.030, 0.031]	
100	0.0802	0.0035	0.0039	0.2311	0.0324	
200	[0.079, 0.082]	[0.003, 0.004]	[0.004, 0.004]	[0.212, 0.250]	[0.032, 0.033]	
200	0.0780	0.0026	0.0048	0.3359	0.0333	
	[0.077, 0.079]	[0.002, 0.003]	[0.004, 0.005]	[0.319, 0.352]	[0.033,0.034]	

Table S9: Estimates of firm and worker heterogeneity across simulations

Notes: Means and 95% confidence intervals. Unobserved heterogeneity is continuously distributed in the DGP. The number of groups  $\hat{K}$  is estimated in every replication, using (4) with  $\gamma = 1$ , and it is reported in the last column of the first panel. We use the kmeans routine from R, with 100 starting values. 500 simulations.

Firm size	$\operatorname{Var}\left(\eta_{i}\right)$	$\operatorname{Var}\left(\psi_{j}\right)$	$\operatorname{Cov}\left(\eta_{i},\psi_{j}\right)$	$\operatorname{Corr}\left(\eta_{i},\psi_{j}\right)$	$\operatorname{Var}\left(\varepsilon_{i1}\right)$	$\hat{K}$
			true va	lues		
-	0.0660	0.0204	0.0050	0.1373	0.0341	
			two-step es	stimator		
10	0.0605 $[0.059, 0.062]$	0.0124 [0.012,0.013]	0.0078 [0.008,0.008]	0.2868 $[0.275, 0.300]$	0.0422 [0.041,0.043]	$3.0$ $_{[3,3]}$
20	0.0626 $[0.061, 0.064]$	0.0155 [0.015,0.016]	0.0068 $[0.006, 0.007]$	$\begin{array}{c} 0.2178 \\ \left[ 0.205, 0.230  ight] \end{array}$	0.0392 [0.038, 0.040]	$4.0 \\ [4,4]$
50	0.0645 [0.063,0.066]	0.0180 [0.017,0.019]	0.0058 [ $0.005, 0.006$ ]	0.1714 [0.158,0.183]	0.0365 [0.036,0.037]	6.0 [6,6]
100	0.0653 [0.064,0.066]	0.0191 [0.018,0.020]	0.0054 [0.005,0.006]	0.1542 [0.141,0.166]	0.0354 [0.035,0.036]	8.0 [8,8]
200	$\begin{array}{c} 0.0657 \\ [0.065, 0.067] \end{array}$	$\begin{array}{c} 0.0198\\ [0.019, 0.021]\end{array}$	$\begin{bmatrix} 0.0052\\ [0.005, 0.006] \end{bmatrix}$	$\begin{array}{c} 0.1448\\ [0.132, 0.157]\end{array}$	$\begin{array}{c} 0.0348\\ [0.034, 0.035]\end{array}$	$     \begin{array}{c}       10.9 \\       [10,12]     \end{array}   $
		two-	step estimator	, bias-corrected	1	
10	0.0650 [0.064,0.066]	0.0149 $[0.014, 0.016]$	0.0056 $[0.005, 0.006]$	0.1445 [0.127,0.163]	0.0397 $[0.039, 0.041]$	
20	0.0647 [0.063,0.066]	0.0185 [0.018,0.019]	$\begin{array}{c} 0.0057 \\ [0.005, 0.006] \end{array}$	0.1499 [0.133,0.167]	$\begin{array}{c} 0.0361 \\ [0.035, 0.037] \end{array}$	
50	0.0656 [0.064,0.067]	0.0202 [0.019,0.021]	0.0053 [0.005,0.006]	0.1416 [0.126,0.155]	0.0344 [0.034,0.035]	
100	0.0661 [0.065,0.067]	0.0202 [0.019,0.021]	0.0050 [ $0.005, 0.005$ ]	0.1371 [0.122,0.150]	0.0344 [0.034,0.035]	
200	$\begin{array}{c} 0.0661 \\ [0.065, 0.067] \end{array}$	$\begin{array}{c} 0.0204 \\ [0.020, 0.021] \end{array}$	$\begin{array}{c} 0.0050 \\ [0.005, 0.005] \end{array}$	$\begin{array}{c} 0.1361 \\ [0.123, 0.149] \end{array}$	$\begin{array}{c} 0.0342 \\ [0.033, 0.035] \end{array}$	
			fixed-effects	estimator		
10	0.1252 [0.123,0.127]	0.0528 [0.051,0.055]	-0.0273 [-0.029,-0.026]	-0.3357 [-0.346,-0.324]	0.0173 [0.017,0.018]	
20	0.0908	0.0318 [0.031,0.033]	-0.0063 [-0.007,-0.006]	-0.1165 [-0.127,-0.105]	0.0256 [0.025,0.026]	
50	0.0752 [0.074,0.076]	0.0242 [0.023,0.025]	$\begin{array}{c} 0.0013 \\ [0.001, 0.002] \end{array}$	$\begin{array}{c} 0.0301 \\ [0.019, 0.041] \end{array}$	0.0307 [0.030,0.031]	
100	0.0705 [0.069,0.072]	0.0222 [0.021,0.023]	0.0033 [0.003,0.004]	0.0827 [0.071,0.095]	0.0324 [0.032,0.033]	
200	0.0683 [0.067,0.069]	$\begin{array}{c} 0.0213 \\ [0.021, 0.022] \end{array}$	$\begin{array}{c} 0.0041 \\ [0.004, 0.005] \end{array}$	$\begin{array}{c} 0.1085 \\ [0.096, 0.120] \end{array}$	$\begin{array}{c} 0.0333 \\ [0.033, 0.034] \end{array}$	

Table S10: Estimates of firm and worker heterogeneity across simulations, one-dimensional firm heterogeneity, large variance of firm effects

Notes: See notes to Table S9. Results for DGP 2.

Firm size	$\mathrm{Var}\left(\eta_{i}\right)$	$\operatorname{Var}\left(\psi_{j}\right)$	$\operatorname{Cov}\left(\eta_{i},\psi_{j}\right)$	$\operatorname{Corr}\left(\eta_{i},\psi_{j}\right)$	$\operatorname{Var}\left(\varepsilon_{i1}\right)$					
	one-dimensional, small firm effect									
-	0.0758	0.0017	0.0057	0.4963	0.0341					
		fixed-	effects, bias-co	orrected						
10	0.0065 [-0.004,0.016]	-0.0717 [-0.082,-0.064]	0.0791 [0.071,0.089]	-0.0976 [-0.125, -0.072]	0.0300 [0.029,0.031]					
20	$\begin{array}{c} 0.0645 \\ [0.062, 0.067] \end{array}$	-0.0098 [-0.011,-0.008]	$\begin{array}{c} 0.0172 \\ [0.016, 0.019] \end{array}$	$\begin{array}{c} 0.0973 \\ [0.073, 0.125] \end{array}$	$\begin{bmatrix} 0.0339 \\ [0.033, 0.035] \end{bmatrix}$					
50	0.0733 $[0.072, 0.075]$	-0.0007 [-0.001,-0.000]	0.0082 [0.008,0.009]	0.3069 $[0.279, 0.335]$	$\begin{array}{c} 0.0341 \\ \left[ 0.033, 0.035  ight] \end{array}$					
100	0.0748 [0.073,0.076]	0.0007 $[0.000, 0.001]$	0.0067 $[0.006, 0.007]$	0.4173 [0.388, 0.447]	0.0341 [0.033,0.035]					
200	0.0753 [0.074,0.077]	$\begin{array}{c} 0.0012 \\ [0.001, 0.002] \end{array}$	$\begin{array}{c} 0.0062 \\ [0.006, 0.007] \end{array}$	$\begin{array}{c} 0.4822 \\ [0.451, 0.512] \end{array}$	0.0341 [0.033,0.035]					
	one-dimensional, large firm effect									
-	0.0660	0.0204	0.0050	0.1373	0.0341					
		fixed-	effects, bias-co	orrected						
10	-0.0036 [-0.013,0.006]	-0.0533	0.0788 [0.070,0.088]	-0.0077 [-0.034, 0.019]	0.0301 [0.029,0.031]					
20	0.0547 [0.053,0.057]	0.0089	0.0166 [0.015,0.018]	0.1163 [0.096,0.137]	0.0339 [0.033, 0.035]					
50	0.0636 [0.062,0.065]	0.0180 [0.017,0.019]	0.0075 [0.007,0.008]	0.1561 [0.139,0.173]	0.0341 [0.033,0.035]					
100	0.0650 [0.064,0.066]	0.0194 [0.019,0.020]	0.0061 [0.006,0.007]	0.1554 [0.139,0.171]	0.0341 [0.033,0.035]					
200	$\begin{array}{c} 0.0656 \\ [0.064, 0.067] \end{array}$	$\begin{array}{c} 0.0199 \\ [0.019, 0.021] \end{array}$	$\begin{bmatrix} 0.0055\\ [0.005, 0.006] \end{bmatrix}$	$\begin{array}{c} 0.1487 \\ [0.133, 0.164] \end{array}$	$\begin{bmatrix} 0.0341 \\ [0.033, 0.035] \end{bmatrix}$					

Table S11: Bias-corrected fixed-effects estimators, one-dimensional firm heterogeneity

Notes: Means and 95% confidence intervals. Unobserved heterogeneity is continuously distributed in the DGP. Bias correction is based on splitting both job movers and job stayers into two sub-samples. The top panel shows the results on DGP 1, with a small variance of firm effects, while the bottom panel shows the results for DGP 2, with a larger variance of firm effects. 500 simulations.

	two-step estimator						two-step estimator, bias corrected						
Firm size	$\mathrm{Var}\left(\eta_{i}\right)$	$\operatorname{Var}\left(\psi_{j}\right)$	$\operatorname{Cov}\left(\eta_{i},\psi_{j}\right)$	$\operatorname{Corr}\left(\eta_{i},\psi_{j}\right)$	$\operatorname{Var}\left(\varepsilon_{i1}\right)$	$\hat{K}$	$\mathrm{Var}\left(\eta_{i}\right)$	$\operatorname{Var}\left(\psi_{j}\right)$	$\operatorname{Cov}\left(\eta_{i},\psi_{j}\right)$	$\operatorname{Corr}\left(\eta_{i},\psi_{j}\right)$	$\operatorname{Var}\left(\varepsilon_{i1}\right)$	$\hat{K}$	
	true values							true values					
-	0.0660	0.0204	0.0050	0.1373	0.0341		0.0660	0.0204	0.0050	0.1373	0.0341		
			$\gamma =$	1.0			$\gamma = 1.0$						
10	0.0513	0.0098	0.0124	0.5500	0.0448	4.0	0.0529	0.0112	0.0115	0.4574	0.0434	4.2	
20	0.0515	0.0112	0.0124	[0.539,0.565] 0.5180 [0.498,0.536]	0.0433 [0.042,0.044]	5.7 [5,6]	0.0514	0.0126	0.0125	[0.437, 0.486] [0.4856] [0.454, 0.509]	[0.042, 0.044] [0.0420] [0.041, 0.043]	[4,5] 7.4 [6,8]	
50	0.0514	0.0123	0.0124	0.4939 [0.471.0.512]	0.0423	8.9 [8.9]	0.0513	0.0131	0.0125	0.4797	0.0415	11.8 [10.12]	
100	0.0519	0.0128	0.0124	0.4796	0.0416	13.3 [13,14]	0.0520	0.0133	0.0123	0.4664	0.0411	17.9 [17,20]	
200	0.0548 [0.051,0.058]	$\begin{array}{c} 0.0147 \\ [0.014, 0.016] \end{array}$	0.0104 [0.009,0.012]	0.3683 [0.303,0.426]	0.0399 [0.039,0.041]	21.4 [20,23]	0.0579 [0.053,0.062]	0.0165 [0.015, 0.018]	0.0089 [0.006,0.011]	0.2713 [0.168,0.374]	0.0381 [0.037,0.040]	30.1 [28,33]	
	$\gamma = 0.5$						$\gamma = 0.5$						
10	0.0498	0.0110	0.0134	0.5730	0.0435	12.8 [12,13]	0.0386	0.0126	0.0123	0.5385 [0.494,0.578]	0.0419	14.0 [12,15]	
20	0.0510	0.0123	0.0125	0.4997	0.0423	16.1 [16.17]	0.0520	0.0136	0.0116	0.4297	0.0410	19.7 [19.22]	
50	0.0536	0.0140	0.0113	0.4134	0.0407	26.0 [24.28]	0.0556	0.0152	0.0104	0.3484	0.0394	34.4 [30,38]	
100	0.0563	0.0153	0.0099	0.3371	0.0392	38.8 [36.41]	0.0589	0.0168	0.0086	0.2635	0.0377	52.9 [48,57]	
200	0.0596	0.0171	0.0085	0.2662	0.0375	53.2 [49,57]	0.0626	0.0187	0.0070	0.1948	0.0359	71.8 [65,78]	
			$\gamma = 0$	).25					$\gamma = 0$	).25			
10	0.0595	0.0119	0.0132	0.4988	0.0428	125.6	0.0504	0.0136	0.0117	0.4379	0.0411	152.6	
20	0.0567	0.0134	0.0119	0.4318	0.0412	138.3 [132.143]	0.0536	0.0149	0.0106	0.3677	0.0397	163.7	
50	0.0574	0.0155	0.0099	0.3316	0.0391	154.0	0.0582	0.0170	0.0085	0.2622	0.0376	190.1	
100	0.0601	0.0171	0.0083	0.2598	0.0374	151.1	0.0624	0.0186	0.0069	0.1932	0.0359	186.1	
200	$\begin{bmatrix} 0.057, 0.063 \end{bmatrix}$ 0.0626 $\begin{bmatrix} 0.058, 0.066 \end{bmatrix}$	[0.016,0.018] 0.0186 [0.018,0.020]	[0.007,0.010] 0.0069 [0.005,0.009]	[0.222,0.303] 0.2027 [0.156,0.251]	[0.037,0.038] 0.0361 [0.035,0.037]	$\begin{bmatrix} 142,163 \end{bmatrix}$ 133.7 $\begin{bmatrix} 127,141 \end{bmatrix}$	$\begin{bmatrix} 0.059, 0.066 \end{bmatrix}$ 0.0649 $\begin{bmatrix} 0.060, 0.069 \end{bmatrix}$	[0.018,0.019] 0.0199 [0.019,0.021]	[0.006,0.008] 0.0056 [0.004,0.007]	$\begin{array}{c} [0.151, 0.240] \\ 0.1512 \\ [0.095, 0.205] \end{array}$	[0.035,0.037] 0.0348 [0.034,0.036]	$\begin{bmatrix} 172,202 \end{bmatrix}$ 162.3 $\begin{bmatrix} 151,176 \end{bmatrix}$	

Table S12: Firm and worker effects, two-dimensional firm heterogeneity, large  $Var(\psi)$ , different choices of  $\gamma$ 

Notes: Means and 95% confidence intervals. Unobserved heterogeneity is continuously distributed in the DGP, with underlying dimension equal to 2. The number of groups K is estimated in every replication, with different choices for  $\gamma$  in (4). 500 simulations. Results for DGP 3.

63

two-step estimator							two-step estimator, bias corrected					
$\mathrm{Var}\left(\eta_{i}\right)$	$\operatorname{Var}\left(\psi_{j}\right)$	$\operatorname{Cov}\left(\eta_{i},\psi_{j}\right)$	$\operatorname{Corr}\left(\eta_{i},\psi_{j}\right)$	$\operatorname{Var}\left(\varepsilon_{i1}\right)$	$\hat{K}$	$\mathrm{Var}\left(\eta_{i}\right)$	$\operatorname{Var}\left(\psi_{j}\right)$	$\operatorname{Cov}\left(\eta_{i},\psi_{j}\right)$	$\operatorname{Corr}\left(\eta_{i},\psi_{j}\right)$	$\operatorname{Var}\left(\varepsilon_{i1}\right)$	$\hat{K}$	
		true va	alues		true values							
0.0758	0.0017	0.0057	0.4963	0.0341		0.0758	0.0017	0.0057	0.4963	0.0341		
$\gamma = 1.0$							$\gamma = 1.0$					
0.0759	0.0008	0.0056	0.7010	0.0350	4.0	0.0760	0.0009	0.0056	0.6487	0.0349	4.0	
[0.074,0.078]	[0.001,0.001]	[0.005,0.006] 0.0059	[0.691, 0.709] 0.6927	[0.034,0.036]	[4,4] 5.5	[0.074,0.078] 0.0749	[0.001,0.001]	[0.005,0.006]	[0.636,0.660]	[0.034,0.036]	[4,4] 6 9	
[0.073,0.077]	[0.001,0.001]	[0.005,0.006]	[0.680,0.707]	[0.034, 0.036]	[5,6]	[0.073,0.077]	[0.001,0.001]	[0.006,0.007]	[0.666,0.705]	[0.034, 0.035]	[6,8]	
0.0750	0.0011	0.0061	0.6877	0.0348	8.0	0.0747	0.0011	0.0062	0.6841	0.0347	10.0	
0.0752	0.0011	0.0062	0.6848	$\begin{bmatrix} 0.034, 0.035 \end{bmatrix}$ 0.0347	[8,8] 11 1	0.072,0.078]	0.0011	0.0063	0.6816	0.034,0.035	14.2	
[0.072, 0.079]	[0.001,0.001]	[0.006,0.007]	[0.668,0.701]	[0.034, 0.035]	[11,12]	[0.072,0.078]	[0.001,0.001]	[0.006,0.007]	[0.660,0.702]	[0.034, 0.035]	[14,16]	
0.0746	0.0011	0.0062	0.6765	0.0347	15.2	0.0745	0.0012	0.0062	0.6720	0.0347	19.3	
[0.069,0.079]	[0.001,0.001]	[0.006,0.007]	[0.654,0.697]	[0.034,0.035]	[14,16]	[0.069,0.079]	[0.001,0.001]	[0.006,0.007]	[0.647,0.694]	[0.034,0.035]	[17,21]	
		$\gamma =$	0.5			$\gamma = 0.5$						
0.0748	0.0010	0.0062	0.7333	0.0349	12.2 [12.13]	0.0731	0.0011	0.0062	0.6905	0.0348	12.7 [12.14]	
0.0747	0.0010	0.0062	0.7076	0.0348	15.1	0.0746	0.0011	0.0063	0.6814	0.0347	18.0	
[0.073, 0.076]	[0.001, 0.001]	[0.006, 0.007]	[0.696, 0.717]	[0.034, 0.035]	[15, 16]	[0.073, 0.076]	[0.001, 0.001]	[0.006, 0.007]	[0.664, 0.695]	[0.034, 0.035]	[18, 20]	
0.0744	0.0011	0.0062	0.6858	0.0347	21.6 [20.23]	0.0744	0.0012	0.0062	0.6717	0.0347	27.0 [24.30]	
0.0743	0.0011	0.0062	0.6709	0.0347	28.2	0.0743	0.0012	0.0062	0.6584	0.0347	35.7	
[0.071, 0.078]	[0.001, 0.001]	[0.006, 0.007]	[0.649, 0.690]	[0.034, 0.035]	[26, 31]	[0.071, 0.078]	[0.001, 0.001]	[0.006, 0.007]	[0.626, 0.684]	[0.034, 0.035]	[32, 40]	
0.0751	0.0012	0.0062	0.6542	0.0347	35.0	0.0751	0.0012	0.0062	0.6409	0.0346	44.0	
[0.071,0.079]	[0.001,0.001]	[0.006,0.007]	[0.631,0.683]	[0.034,0.035]	[31,39]	[0.071,0.080]	[0.001,0.001]	[0.006,0.007]	[0.603,0.681]	[0.034,0.035]	[38,50]	
		$\gamma = 0$	0.25					$\gamma = 0$	).25			
0.0796	0.0012	0.0062	0.6355	0.0346	124.3	0.0699	0.0013	0.0061	0.6190	0.0345	148.7	
[0.078,0.082]	0.0013	0.0061	[0.605,0.657]	[0.034,0.035]	[121, 127] 131.0	[0.066,0.073]	[0.001,0.002]	0.0060	0.6086	[0.034,0.035]	[142, 154] 150 3	
[0.074,0.078]	[0.001, 0.001]	[0.006,0.007]	[0.602,0.643]	[0.034, 0.035]	[125, 137]	[0.070,0.074]	[0.001, 0.002]	[0.006,0.006]	[0.567,0.647]	[0.034, 0.035]	[140, 160]	
0.0752	0.0014	0.0061	0.6020	0.0345	134.7	0.0749	0.0014	0.0060	0.5800	0.0344	159.0	
[0.072,0.077]	[0.001,0.002]	[0.006,0.007]	[0.572,0.624]	[0.034,0.035]	[127,142]	[0.072,0.077]	[0.001,0.002]	[0.006,0.006]	[0.529,0.622]	[0.034,0.035]	[146,171]	
0.0752	0.0014	[0.000.0.00	0.5879	0.0345	125.1 [116.132]	0.0752	0.0015	0.0060	0.5651	0.0344	140.7 [132.158]	
0.0754	0.0014	0.0060	0.5781	0.0344	105.4	0.0755	0.0015	0.0060	0.5569	0.0344	121.7	
[0.071, 0.079]	[0.001, 0.002]	[0.005, 0.007]	[0.545, 0.606]	[0.034, 0.035]	[99, 113]	[0.071, 0.079]	[0.001, 0.002]	[0.005, 0.006]	[0.498, 0.604]	[0.034, 0.035]	[107, 134]	
	$\begin{array}{c} {\rm Var}\left(\eta_i\right)\\ \hline\\ 0.0758\\ \hline\\ 0.0754\\ 0.0754\\ \hline\\ 0.073.0078\\ 0.0752\\ \hline\\ 0.073.0079\\ 0.0752\\ \hline\\ 0.0752\\ \hline\\ 0.072.0079\\ \hline\\ 0.0746\\ \hline\\ 0.073.0076\\ \hline\\ 0.0747\\ \hline\\ 0.0747\\ \hline\\ 0.0743\\ \hline\\ 0.0752\\ \hline\\ 0.0754\\ \hline\\ 0.0752\\ \hline\\ $	Var $(\eta_i)$ Var $(\psi_j)$ 0.0758         0.0017           0.0759         0.0008           0.0759         0.0009           0.0750         0.0009           0.0751         0.0009           0.07520         0.0011           0.0753         0.0011           0.0750         0.0011           0.0752         0.0011           0.0753         0.0011           0.0754         0.0011           0.0755         0.0011           0.0750         0.0011           0.0747         0.0010           0.0747         0.0010           0.0747         0.0011           0.0743         0.0011           0.0743         0.0011           0.0751         0.0012           0.0753         0.0012           0.0754         0.0012           0.0758         0.0013           0.0752         0.0014           0.0752         0.0014           0.0754         0.0014           0.0754         0.0014           0.0754         0.0014	$\begin{array}{c c} & \text{two-step e} \\ \text{Var}\left(\eta_i\right) & \text{Var}\left(\psi_j\right) & \text{Cov}\left(\eta_i,\psi_j\right) \\ & \text{true va} \\ 0.0758 & 0.0017 & 0.0057 \\ \hline \\ $	$\begin{array}{c c c c c c c c c c c c c c c c c c c $	$\begin{array}{c c c c c c c c c c c c c c c c c c c $	$\begin{array}{c c c c c c c c c c c c c c c c c c c $	$\begin{array}{c c c c c c c c c c c c c c c c c c c $	$\begin{array}{c c c c c c c c c c c c c c c c c c c $	$\begin{array}{c c c c c c c c c c c c c c c c c c c $	$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	

Table S13: Firm and worker effects, two-dimensional firm heterogeneity, small  $Var(\psi)$ , different choices of  $\gamma$ 

Notes: Means and 95% confidence intervals. Unobserved heterogeneity is continuously distributed in the DGP, with underlying dimension equal to 2. The number of groups K is estimated in every replication, with different choices for  $\gamma$  in (4). 500 simulations. Results for DGP 4.

	iterated estimator						iterated estimator, bias corrected					
Firm size	$\operatorname{Var}\left(\eta_{i}\right)$	$\operatorname{Var}\left(\psi_{j}\right)$	$\operatorname{Cov}\left(\eta_{i},\psi_{j}\right)$	$\operatorname{Corr}\left(\eta_{i},\psi_{j}\right)$	$\operatorname{Var}\left(\varepsilon_{i1}\right)$	$\hat{K}$	$\operatorname{Var}\left(\eta_{i}\right)$	$\operatorname{Var}\left(\psi_{j}\right)$	$\operatorname{Cov}\left(\eta_{i},\psi_{j}\right)$	$\operatorname{Corr}\left(\eta_{i},\psi_{j}\right)$	$\operatorname{Var}\left(\varepsilon_{i1}\right)$	$\hat{K}$
	true values								true va	alues		
-	0.0660	0.0204	0.0050	0.1373	0.0341		0.0660	0.0204	0.0050	0.1373	0.0341	
			$\gamma =$	1.0					$\gamma =$	1.0		
10	0.0661	0.0045	0.0050	0.2847	0.0501	4.0	0.0592	0.0079	0.0084	0.4058	0.0468	4.2
20	0.0632	0.0080	0.0066	0.2928	0.0466	[4,4] 5.7 [r.c]	0.0592	0.0117	0.0085	0.3192	0.0429	[4,5] 7.4 [c.e]
50	0.0608	0.0127	0.0077	0.2785	0.0420	[5,6] 8.9	0.0589	0.0163	0.0086	0.2702	0.0383	[0,6] 11.8
100	0.0617	0.0152	0.0074	0.2424	0.0392	[8,9] 13.3	0.0624	0.0179	0.0071	0.2044	0.0366	17.9
200	0.0628	0.0174	0.0064	0.1932 [0.139,0.242]	0.0371	$\begin{bmatrix} 13, 14 \\ 21.4 \\ 20, 23 \end{bmatrix}$	0.0642	0.0196	0.0057	0.1552	0.0350	30.1 [28,33]
	$\gamma = 0.5$						$\gamma = 0.5$					
10	0.0549	0.0093	0.0106	0.4708	0.0452	12.8	0.0473	0.0141	0.0144	0.5393	0.0404	14.0
20	0.0555	[0.009,0.010] 0.0117	0.010,0.011	[0.451,0.489] 0 4027	[0.044,0.046]	[12,13] 16.1	[0.038,0.053] 0.0560	[0.012,0.017]	[0.012,0.018]	[0.433,0.684] 0.3396	[0.037,0.042] 0.0405	[12,15] 19.7
20	[0.054, 0.057]	[0.011,0.012]	[0.010,0.011]	[0.385,0.424]	[0.042, 0.044]	[16,17]	[0.054, 0.058]	[0.013,0.015]	[0.009,0.011]	[0.311,0.375]	[0.039, 0.041]	[19,22]
50	0.0579	0.0148	0.0092	0.3135	0.0399	26.0	0.0599	0.0170	0.0082	0.2436	0.0376	34.4
100	0.0605	0.0167	0.0077	0.2437	0.0378	38.8	0.0630	0.0188	0.0065	0.1780	0.0358	52.9
	[0.057, 0.063]	[0.016, 0.018]	[0.006, 0.009]	[0.206, 0.281]	[0.037, 0.039]	[36, 41]	[0.060, 0.066]	[0.018, 0.020]	[0.005, 0.008]	[0.135, 0.221]	[0.035, 0.037]	[48, 57]
200	0.0631 [0.059, 0.067]	0.0185 [0.018,0.019]	0.0067 [0.005,0.008]	0.1976 [0.147,0.248]	0.0362 [0.036,0.037]	53.2 [49,57]	0.0653 [0.061,0.069]	0.0200	0.0056 [0.004,0.008]	0.1505 [0.097,0.209]	0.0346 [0.034,0.035]	71.8 [65,78]
			$\gamma = 0$	0.25					$\gamma = 0$	).25		
10	0.0653	0.0118	0.0129	0.4641	0.0428	125.6	0.0391	0.0141	0.0121	0.4642	0.0405	152.6
20	[0.063, 0.068] 0.0574	0.0137	0.0114	[0.443, 0.481] 0.4078	0.042,0.044]	[121,130] 138.3	0.034,0.045	[0.013, 0.015] 0.0155	[0.011, 0.013] 0.0100	[0.432, 0.494] 0.3537	[0.040, 0.041] 0.0391	[145,160] 163.7
50	0.056,0.060	[0.013,0.014]	0.011,0.012	0.391,0.425	0.040,0.042	[132,143]	0.047,0.053	0.015,0.016	[0.009,0.011]	0.326,0.380	0.038,0.040	[153,172]
50	0.0583	0.0162	0.0090	0.2939	0.0384	154.0 [146_163]	0.0599	0.0180	0.0073	0.2117	0.0365	190.1 [177-204]
100	0.0618	0.0179	0.0074	0.2222	0.0366	151.1	0.0646	0.0195	0.0058	0.1549	0.0350	186.1
	[0.059, 0.065]	[0.017, 0.019]	[0.006, 0.009]	[0.186, 0.270]	[0.036, 0.037]	[142, 163]	[0.062, 0.068]	[0.019, 0.020]	[0.004, 0.008]	[0.115, 0.205]	[0.034, 0.036]	[172, 202]
200	0.0642 [0.060, 0.068]	0.0193 [0.018, 0.020]	0.0061 [0.004,0.008]	0.1734 [0.127,0.225]	0.0354 [0.035,0.036]	133.7 [127,141]	0.0663 [0.062,0.070]	0.0205 [0.020, 0.022]	0.0050 [0.003,0.007]	0.1306 [0.082,0.186]	0.0342 [0.033, 0.035]	162.3 [151,176]

Table S14: Firm and worker effects, two-dimensional firm heterogeneity, large  $Var(\psi)$ , different choices of  $\gamma$ , iterated estimators

Notes: Means and 95% confidence intervals. Unobserved heterogeneity is continuously distributed in the DGP, with underlying dimension equal to 2. The number of groups K is estimated in every replication, with different choices for  $\gamma$  in (4). 500 simulations. Results for DGP 3.

65

	iterated estimator							iterated estimator, bias corrected					
Firm size	$\operatorname{Var}\left(\eta_{i}\right)$	$\operatorname{Var}\left(\psi_{j}\right)$	$\operatorname{Cov}\left(\eta_{i},\psi_{j}\right)$	$\operatorname{Corr}\left(\eta_{i},\psi_{j}\right)$	$\operatorname{Var}\left(\varepsilon_{i1}\right)$	$\hat{K}$	$\mathrm{Var}\left(\eta_{i}\right)$	$\operatorname{Var}\left(\psi_{j}\right)$	$\operatorname{Cov}\left(\eta_{i},\psi_{j}\right)$	$\operatorname{Corr}\left(\eta_{i},\psi_{j}\right)$	$\operatorname{Var}\left(\varepsilon_{i1}\right)$	$\hat{K}$	
	true values							true values					
-	0.0758	0.0017	0.0057	0.4963	0.0341		0.0758	0.0017	0.0057	0.4963	0.0341		
	$\gamma = 1.0$								$\gamma = 1$	1.0			
10	0.0866	0.0000	0.0003	0.1290 [ $0.108, 0.150$ ]	0.0358	4.0 [4,4]	0.0867	0.0000	0.0002	0.1169 [0.075,0.153]	0.0358	4.0 [4,4]	
20	0.0845	0.0002	0.0013	0.2921	0.0356	5.5 [5.6]	0.0823	0.0004	0.0024	0.4556	0.0354	6.9	
50	0.0791	0.0007	0.0041	0.5444 [0.516,0.573]	0.0352	[8,0] [8,8]	0.0761	0.0010	0.0056	0.6584	0.0349	10.0 [10,10]	
100	0.0775	0.0009	0.0050	0.6035	0.0349	11.1 [11.12]	0.0756	0.0011	0.0060	0.6600	0.0347	14.2 [14.16]	
200	0.0759 [0.071,0.080]	$\begin{array}{c} 0.0011 \\ [0.001, 0.001] \end{array}$	0.0055 [0.005,0.006]	0.6174 [0.589,0.648]	0.0348 [0.034,0.035]	15.2 [14,16]	0.0750 [0.070,0.079]	$\begin{array}{c} 0.0012\\ [0.001, 0.001] \end{array}$	0.0060 [0.005,0.007]	0.6345 [0.593,0.674]	$\begin{array}{c} 0.0347 \\ [0.034, 0.035] \end{array}$	19.3 [17,21]	
	$\gamma = 0.5$						$\gamma = 0.5$						
10	0.0799	0.0006	0.0037	0.5408	0.0353	12.2 [12,13]	0.0777	0.0008	0.0047	0.6157	0.0351	12.7 [12,14]	
20	0.0780	0.0008	0.0046	0.5940	0.0351	15.1 [15,16]	0.0760	0.0010	0.0056	0.6482	0.0349	18.0 [18,20]	
50	0.0760	0.0010	0.0054	0.6299	0.0349	21.6 [20,23]	0.0748	0.0011	0.0060	0.6568	0.0347	27.0 [24,30]	
100	0.0753	0.0011	0.0057	0.6364	0.0348	28.2	0.0746	0.0012	0.0061	0.6457	0.0347	35.7	
200	0.0756 [0.071,0.080]	0.0012 [0.001, 0.001]	0.0059 [0.005,0.007]	0.6231 [0.594,0.656]	0.0347 [0.034,0.035]	35.0 [31,39]	0.0753 [0.071,0.080]	0.0013 [0.001,0.002]	0.0061 [0.006,0.007]	0.6137 [0.569,0.657]	0.0346 [0.034,0.035]	44.0 [38,50]	
			$\gamma = 0$	0.25			$\gamma = 0.25$						
10	0.0781	0.0011	0.0059	0.6217 [ $0.602, 0.642$ ]	0.0347	124.3 [121,127]	0.0650	0.0014	0.0063	0.6481 [0.602,0.692]	0.0345	148.7 [142,154]	
20	0.0755	0.0012	0.0059	0.6122	0.0346	131.0	0.0731	0.0014	0.0060	0.5983	0.0345	150.3	
50	0.0754	0.0013	0.0060	0.5928	0.0345	134.7	0.0752	0.0014	0.0059	0.5713	0.0344	159.0	
100	0.0754	0.0014	0.0060	0.5806	0.0345	125.1	0.0753	0.0015	0.0060	0.5597	0.0344	146.7	
200	[0.072,0.078] 0.0755 [0.071,0.080]	$\begin{array}{c} [0.001, 0.002] \\ 0.0014 \\ [0.001, 0.002] \end{array}$	[0.006,0.006] 0.0060 [0.005,0.006]	$\begin{array}{c} [0.547, 0.609] \\ 0.5704 \\ [0.532, 0.602] \end{array}$	$\begin{bmatrix} 0.034, 0.035 \end{bmatrix}$ 0.0344 $\begin{bmatrix} 0.034, 0.035 \end{bmatrix}$	$\begin{bmatrix} 116, 132 \end{bmatrix}$ 105.4 $\begin{bmatrix} 99, 113 \end{bmatrix}$	[0.072,0.078] 0.0756 [0.071,0.080]	$\begin{bmatrix} 0.001, 0.002 \end{bmatrix}$ 0.0015 $\begin{bmatrix} 0.001, 0.002 \end{bmatrix}$	[0.006,0.006] 0.0059 [0.005,0.007]	$\begin{array}{c} [0.503, 0.609] \\ 0.5499 \\ [0.490, 0.604] \end{array}$	[0.034,0.035] 0.0343 [0.034,0.035]	$\begin{bmatrix} 132,158 \end{bmatrix}$ 121.7 $\begin{bmatrix} 107,134 \end{bmatrix}$	

Table S15: Firm and worker effects, two-dimensional firm heterogeneity, small  $Var(\psi)$ , different choices of  $\gamma$ , iterated estimators

Notes: Means and 95% confidence intervals. Unobserved heterogeneity is continuously distributed in the DGP, with underlying dimension equal to 2. The number of groups K is estimated in every replication, with different choices for  $\gamma$  in (4). 500 simulations. Results for DGP 4.

66

Firm size	$\mathrm{Var}\left(\eta_{i}\right)$	$\operatorname{Var}\left(\psi_{j}\right)$	$\operatorname{Cov}\left(\eta_{i},\psi_{j}\right)$	$\operatorname{Corr}\left(\eta_{i},\psi_{j}\right)$	$\operatorname{Var}\left(\varepsilon_{i1}\right)$	% misclass.				
	true values									
-	0.0758	0.0017	0.0057	0.4963	0.0341					
			two-step wit	h $K = K^* = 10$	)					
10	0.0758 [0.074,0.077]	0.0013 [0.001,0.001]	0.0057 $[0.005, 0.006]$	0.5770 [0.566,0.586]	0.0346 [0.034,0.035]	69.0%				
20	0.0758 [0.074,0.077]	$\begin{array}{c} 0.0015 \\ [0.001, 0.002] \end{array}$	$\begin{array}{c} 0.0057 \\ [0.005, 0.006] \end{array}$	$\begin{array}{c} 0.5355 \\ [0.525, 0.546] \end{array}$	$\begin{array}{c} 0.0344 \\ [0.034, 0.035] \end{array}$	58.5% [0.560,0.614]				
50	0.0759 [0.075,0.077]	$\begin{array}{c} 0.0016 \\ [0.001, 0.002] \end{array}$	0.0056 $[0.005, 0.006]$	0.5083 $[0.499, 0.517]$	$\begin{array}{c} 0.0342 \\ [0.034, 0.035] \end{array}$	39.3% [0.338,0.476]				
100	0.0759 $[0.075, 0.077]$	$\begin{array}{c} 0.0017 \\ [0.001, 0.002] \end{array}$	0.0056 $[0.005, 0.006]$	0.4981 [0.489,0.507]	$\begin{array}{c} 0.0342 \\ [0.033, 0.035] \end{array}$	22.6% [0.171,0.359]				
200	0.0759 [0.074,0.077]	$\begin{array}{c} 0.0017 \\ [0.002, 0.002] \end{array}$	0.0056 $[0.005, 0.006]$	$\begin{array}{c} 0.4945 \\ [0.484, 0.504] \end{array}$	$\begin{array}{c} 0.0341 \\ \left[ 0.033, 0.035  ight] \end{array}$	7.5% [0.050,0.115]				
		ł	bias corrected	with estimated	Κ					
10	0.0778 $[0.076, 0.079]$	0.0013 $[0.001, 0.002]$	0.0047 $[0.004, 0.005]$	0.4527 $[0.441, 0.465]$	0.0346 $[0.034, 0.035]$					
20	0.0762 [0.075,0.078]	$\begin{array}{c} 0.0016 \\ [0.001, 0.002] \end{array}$	0.0055 [0.005,0.006]	0.4917 [0.478, 0.502]	0.0342 [0.034,0.035]					
50	0.0760 [0.075,0.077]	$\begin{array}{c} 0.0017 \\ [0.001, 0.002] \end{array}$	0.0056 [0.005,0.006]	0.4906 [0.478,0.503]	0.0342 [0.033,0.035]					
100	0.0759 [0.074,0.077]	$\begin{array}{c} 0.0017 \\ [0.002, 0.002] \end{array}$	0.0057 $[0.005, 0.006]$	0.4909 [0.480, 0.501]	0.0341 [0.033,0.035]					
200	0.0757 [0.074,0.077]	$\begin{array}{c} 0.0018\\ [0.002, 0.002]\end{array}$	0.0057 $[0.005, 0.006]$	$\begin{array}{c} 0.4930 \\ [0.483, 0.503] \end{array}$	$\begin{array}{c} 0.0341 \\ [0.033, 0.035] \end{array}$					

Table S16: Firm and worker effects, discrete firm heterogeneity  $(K^* = 10)$ 

Notes: Means and 95% confidence intervals. Unobserved heterogeneity is discretely distributed in the DGP, with  $K^* = 10$  groups. In the top panel the true number of groups is used. The last column shows frequencies of misclassification. In the bottom panel the number of groups is estimated in every replication. 500 simulations. Results for DGP 5.





Notes: See notes to Figures 4 and 5. a and b are the intercept and slope in the probability of being a mover type.



Figure S2: Parameter estimates across simulations, fixed K

Notes: See note to Figure 4. K is kept fixed, and we focus on the model with homogeneous mobility costs. 500 replications.



Figure S3: Estimates of firm and worker heterogeneity across simulations

Notes: Means (solid line) and 95% confidence intervals. The dashed line indicates the true parameter value. Unobserved heterogeneity is continuously distributed in the DGP. The number of groups K is estimated in every replication. 500 replications.



Figure S4: Dimension of firm heterogeneity

Notes: Source Swedish administrative data. Left graph shows the logarithm of  $\widehat{Q}(K)$  as a function of K, for different average firm sizes S. Horizontal lines show the corresponding value of  $\ln(\widehat{V}_h)$ . The right graph shows the relationship between the log of  $\widehat{K}$  and the log of the average firm size in the sample, across samples.

Figure S5: Estimates of firm and worker heterogeneity across simulations, two-dimensional firm heterogeneity, large variance of firm effects



Notes: Means (solid line) and 95% confidence intervals.  $\blacksquare$  indicates the two-step bias-corrected GFE estimator and  $\blacktriangle$  indicates the iterated bias-corrected GFE estimator. The different columns represent different values of  $\gamma$  (that is, different selection rules for the number of groups). Unobserved heterogeneity is continuously distributed in the DGP. The number of groups K is estimated in every replication. 500 replications. Results for DGP 3.


Figure S6: Two-dimensional firm heterogeneity, small variance of firm effects

Notes: See the notes to Figure S5. Results for DGP 4.



Figure S7: Estimates of firm and worker heterogeneity across simulations, two-dimensional firm heterogeneity, large variance of firm effects, different number of job movers per firm

Notes: Means (solid line) and 95% confidence intervals. 
indicates the two-step bias-corrected GFE estimator,  $\blacktriangle$  the iterated bias-corrected GFE estimator,  $\bigcirc$  the fixed-effects estimator, and  $\bullet$  the biascorrected fixed-effects estimator. The different columns represent different values of  $\gamma$  (that is, different selection rules for the number of groups). Unobserved heterogeneity is continuously distributed in the DGP. The number of groups K is estimated in every replication. 500 replications. Results for DGP743.



Figure S8: Estimates of firm and worker heterogeneity across simulations, two-dimensional firm heterogeneity, small variance of firm effects, different number of job movers per firm

Notes: See the notes to Figure S7. Results for DGP 4.