



Identification and estimation of categorical random coefficient models

Zhan Gao¹ · M. Hashem Pesaran^{1,2}

Received: 25 April 2022 / Accepted: 27 February 2023 / Published online: 6 April 2023
© The Author(s) 2023

Abstract

This paper proposes a linear categorical random coefficient model, in which the random coefficients follow parametric categorical distributions. The distributional parameters are identified based on a linear recurrence structure of moments of the random coefficients. A generalized method of moments estimation procedure is proposed, also employed by Peter Schmidt and his coauthors to address heterogeneity in time effects in panel data models. Using Monte Carlo simulations, we find that moments of the random coefficients can be estimated reasonably accurately, but large samples are required for the estimation of the parameters of the underlying categorical distribution. The utility of the proposed estimator is illustrated by estimating the distribution of returns to education in the USA by gender and educational levels. We find that rising heterogeneity between educational groups is mainly due to the increasing returns to education for those with postsecondary education, whereas within-group heterogeneity has been rising mostly in the case of individuals with high school or less education.

Keywords Random coefficient models · Categorical distribution · Return to education

JEL Classification C01 · C21 · C13 · C46 · J30

✉ Zhan Gao
zhangao@usc.edu

M. Hashem Pesaran
pesaran@usc.edu

¹ Department of Economics, University of Southern California, 3620 South Vermont Avenue, Los Angeles, CA 90089, USA

² Trinity College, Cambridge, UK

1 Introduction

Random coefficient models have been used extensively in time series, cross-section and panel regressions. Nicholls and Pagan (1985) consider the estimation of first and second moments of the random coefficient β_i and the error term u_i , in a linear regression model. In a seminal paper, Beran and Hall (1992) establish conditions for identifying and estimating the distribution of β_i and u_i nonparametrically. The baseline linear univariate regression in Beran and Hall (1992) has been extended in nonparametric framework by Beran (1993), Beran and Millar (1994), Beran et al. (1996), Hoderlein et al. (2010), Hoderlein et al. (2017) and Breunig and Hoderlein (2018), to just name a few. Hsiao and Pesaran (2008) survey random coefficient models in linear panel data models.

In some econometric applications, Hausman (1981), Hausman and Newey (1995), Foster and Hahn (2000), for examples, the main interest is to estimate the consumer surplus distribution based on a linear demand system where the coefficient associated with the price is random. In such settings, the distribution of the random coefficients is needed when computing the consumer surplus function, and the nonparametric estimation is more general, flexible and suitable for the purpose. On the other hand, parametric models may be favored in applications in which the implied economic meaning of the distribution of the random coefficients is of interests. Examples include estimation of the return to education (Lemieux 2006b, c) and the labor supply equation (Bick et al. 2022).

In this paper, we consider a linear regression model with a random coefficient β_i that is assumed to follow a categorical distribution, i.e., β_i has a discrete support $\{b_1, b_2, \dots, b_K\}$, and $\beta_i = b_k$ with probability π_k . The discretization of the support of the random coefficient β_i naturally corresponds to the interpretation that each individual belongs to a certain category, or group, k with probability π_k . Compared to a nonparametric distribution with continuous support, assuming a categorical distribution allows us not only to model the heterogeneous responses across individuals but also to interpret the results with sharper economic meaning. As we will illustrate in the empirical application in Sect. 6, it is hard to clearly interpret the distribution of returns to education without imposing some form of parametric restrictions.

In addition, with the categorical distribution imposed, the identification and estimation of the distribution of β_i do not rely on identically distributed error terms u_i and regressors w_i , as shown in Sect. 2 and 3. Heterogeneously generated errors can be allowed, which is important in many empirical applications. To the best of our knowledge, this is the first identification result in linear random coefficient model without a strict IID setting.

The identification of the distribution of β_i is established in this paper based on the identification of the moments of β_i , which coincides with the identification condition in Beran and Hall (1992) that the distribution of β_i is uniquely determined by its moments, which is assumed to exist up to an arbitrary order. Since under our setup the distribution of β_i is parametrically specified, the moments of β_i exist and can be derived explicitly. The parameters of the assumed categorical distribution can then be uniquely determined by a system of equations in terms of the moments, as in Theorem 2. The parameters of the categorical distribution are then estimated consistently by the

generalized method of moments (GMM). The estimation procedure based on moment conditions shares similar spirits as in Ahn et al. (2001, 2013) in which Peter Schmidt and coauthors study panel data models with interactive effects where they allow for the time effects to vary across individual units. Compared to alternative nonparametric random coefficient models, the standard GMM estimation is easy to implement, and the identified categorical structure has a clear economic interpretation.

Using Monte Carlo (MC) simulations, we find that moments of the random coefficients can be estimated reasonably accurately, but large samples are required for estimation of the parameters of the underlying categorical distributions. Our theoretical and MC results also suggest that our method is suitable when the number of heterogeneous coefficients and the number of categories are small (2 or 3). With the number of categories rising the burden on identification from the moments to the parameters of the categorical distribution also rises rapidly. The quality of identification also deteriorates as we need to rely on higher and higher moments to identify a larger number of categories, since the information content of the moments tends to decline with their order.

The proposed method is also illustrated by providing estimates of the distribution of returns to education in the USA by gender and educational levels, using the May and Outgoing Rotation Group (ORG) supplements of the Current Population Survey (CPS) data. Comparing the estimates obtained over the sub-periods 1973–1975 and 2001–2003, we find that rising between group heterogeneity is largely due to rising returns to education in the case of individuals with postsecondary education, while within-group heterogeneity has been rising in the case of individuals with high school or less education.

Related Literature This paper draws mainly upon the literature of random coefficient models. As already mentioned, the main body of the recent literature is focused on nonparametric identification and estimation. Following Beran and Hall (1992), Beran (1993) and Beran and Millar (1994) extend the model to a linear semi-parametric model with a multivariate setup and propose a minimum distance estimator for the unknown distribution. Foster and Hahn (2000) extend the identification results in Beran and Hall (1992) and apply the minimum distance estimator to a gasoline consumption data to estimate the consumer surplus function. Beran et al. (1996) and Hoderlein et al. (2010) propose kernel density estimators based on the Radon inverse transformation in linear models.

In addition to linear models, Ichimura and Thompson (1998) and Gautier and Kitamura (2013) incorporate the random coefficients in binary choice models. Gautier and Hoderlein (2015) and Hoderlein et al. (2017) consider triangular models with random coefficients allowing for causal inference. Matzkin (2012) and Masten (2018) discuss the identification of random coefficients in simultaneous equation models. Breunig and Hoderlein (2018) propose a general specification test in a variety of random coefficient models. Random coefficients are also widely studied in panel data models, for example Hsiao and Pesaran (2008) and Arellano and Bonhomme (2012)

The rest of the paper is organized as follows: Sect. 2 establishes the main identification results. The GMM estimation procedure is proposed and discussed in Sect. 3. An extension to a multivariate setting is considered in Sect. 4. Small sample properties of the proposed estimator are investigated in Sect. 5, using Monte Carlo techniques

under different regressor and error distributions. Section 6 presents and discusses our empirical application to the return to education. Section 7 provides some concluding remarks and suggestions for future work. Technical proofs are given in “Appendix A.1.”

Notations Largest and smallest eigenvalues of the $p \times p$ matrix $\mathbf{A} = (a_{ij})$ are denoted by $\lambda_{\max}(\mathbf{A})$ and $\lambda_{\min}(\mathbf{A})$, respectively, its spectral norm by $\|\mathbf{A}\| = \lambda_{\max}^{1/2}(\mathbf{A}'\mathbf{A})$, $\mathbf{A} > 0$ means that \mathbf{A} is positive definite, $\text{vech}(\mathbf{A})$ denotes the vectorization of distinct elements of \mathbf{A} , $\mathbf{0}$ denotes zero matrix (or vector). For $\mathbf{a} \in \mathbb{R}^p$, $\text{diag}(\mathbf{a})$ represents the diagonal matrix with elements of a_1, a_2, \dots, a_p . For random variables (or vectors) u and v , $u \perp v$ represents u is independent of v . We use $c(C)$ to denote some small (large) positive constants. For a differentiable real-valued function $f(\boldsymbol{\theta})$, $\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta})$ denotes the gradient vector. Operator \rightarrow_p denotes convergence in probability, and \rightarrow_d convergence in distribution. The symbols $O(1)$, and $O_p(1)$ denote asymptotically bounded deterministic and random sequences, respectively.

2 Categorical random coefficient model

We suppose the single cross-section observations, $\{y_i, x_i, \mathbf{z}_i\}_{i=1}^n$, follow the categorical random coefficient model

$$y_i = x_i \beta_i + \mathbf{z}_i' \boldsymbol{\gamma} + u_i, \tag{2.1}$$

where $y_i, x_i \in \mathbb{R}$, $\mathbf{z}_i \in \mathbb{R}^{p_z}$, and $\beta_i \in \{b_1, b_2, \dots, b_K\}$ admits the following K -categorical distribution,

$$\beta_i = \begin{cases} b_1, & \text{w.p. } \pi_1, \\ b_2, & \text{w.p. } \pi_2, \\ \vdots & \vdots \\ b_K, & \text{w.p. } \pi_K, \end{cases} \tag{2.2}$$

w.p. denotes “with probability,” $\pi_k \in (0, 1)$, $\sum_{k=1}^K \pi_k = 1$, $b_1 < b_2 < \dots < b_K$, $\boldsymbol{\gamma} \in \mathbb{R}^{p_z}$ is homogeneous and \mathbf{z}_i could include an intercept term as its first element. It is assumed that $\beta_i \perp \mathbf{w}_i = (x_i, \mathbf{z}_i)'$, and the idiosyncratic errors u_i are independently distributed with mean 0.

Remark 1 The model can be extended to allow $\mathbf{x}_i, \boldsymbol{\beta}_i \in \mathbb{R}^p$, with $\boldsymbol{\beta}_i$ following a multivariate categorical distribution, though with more complicated notations. We will consider possible extensions in Sect. 4.

Remark 2 Since we consider a pure cross-sectional setting, the key assumption that β_i and x_i are independently distributed cannot be relaxed. Allowing β_i to vary with \mathbf{w}_i , without any further restrictions, is tantamount to assuming y_i is a general function of \mathbf{w}_i , in effect rendering a nonparametric specification.

Remark 3 The number of categories, K , is assumed to be fixed and known. Conditions $\sum_{k=1}^K \pi_k = 1$, $b_1 < b_2 < \dots < b_K$, and $\pi_k \in (0, 1)$ together are sufficient for the existence of K categories. For example, if $b_k = b_{k'}$, then we can merge categories k and k' , and the number of categories reduces to $K - 1$. Similarly, if $\pi_k = 0$ for some k , then category k can be deleted, and the number of categories is again reduced to $K - 1$. Information criteria can be used to determine K , but this will not be pursued in this paper. Model specification tests could also be considered. See, for examples, Andrews (2001) and Breunig and Hoderlein (2018).

In the rest of this section, we focus on the model (2.1) and establish the conditions under which the distribution of β_i is identified.

2.1 Identifying the moments of β_i

Assumption 1 (a) (i) u_i is distributed independently of $\mathbf{w}_i = (x_i, \mathbf{z}'_i)'$ and β_i . (ii) $\sup_i E(\|u_i^r\|) < C, r = 1, 2, \dots, 2K - 1$. (iii) $n^{-1} \sum_{i=1}^n u_i^4 = O_p(1)$.
 (b) (i) Let $\mathbf{Q}_{n,ww} = n^{-1} \sum_{i=1}^n \mathbf{w}_i \mathbf{w}'_i$, and $\mathbf{q}_{n,wy} = n^{-1} \sum_{i=1}^n \mathbf{w}_i y_i$. Then $\|E(\mathbf{Q}_{n,ww})\| < C < \infty$, and $\|E(\mathbf{q}_{n,wy})\| < C < \infty$, and there exists $n_0 \in \mathbb{N}$ such that for all $n \geq n_0$,

$$0 < c < \lambda_{\min}(\mathbf{Q}_{n,ww}) < \lambda_{\max}(\mathbf{Q}_{n,ww}) < C < \infty.$$

(ii) $\sup_i E(\|\mathbf{w}_i\|^r) < C < \infty, r = 1, 2, \dots, 4K - 2$.
 (iii) $n^{-1} \sum_{i=1}^n \|\mathbf{w}_i\|^4 = O_p(1)$.
 (c) $\|\mathbf{Q}_{n,ww} - E(\mathbf{Q}_{n,ww})\| = O_p(n^{-1/2}), \|\mathbf{q}_{n,wy} - E(\mathbf{q}_{n,wy})\| = O_p(n^{-1/2})$, and

$$E(\mathbf{Q}_{n,ww}) = n^{-1} \sum_{i=1}^n E(\mathbf{w}_i \mathbf{w}'_i) > 0.$$

(d) $\|E(\mathbf{Q}_{n,ww}) - \mathbf{Q}_{ww}\| = O(n^{-1/2}), \|E(\mathbf{q}_{n,wy}) - \mathbf{q}_{wy}\| = O(n^{-1/2})$, where $\mathbf{q}_{wy} = \lim_{n \rightarrow \infty} E(\mathbf{q}_{n,wy}), \mathbf{Q}_{ww} = \lim_{n \rightarrow \infty} E(\mathbf{Q}_{n,ww})$ and $\mathbf{Q}_{ww} > 0$.

Remark 4 Part (a) of Assumption 1 relaxes the assumption that u_i is identically distributed, and allows for heterogeneously generated errors. For identification of the distribution of β_i , we require u_i to be distributed independently of \mathbf{w}_i and β_i , which rules out conditional heteroskedasticity. However, estimation and inference involving $E(\beta_i)$ and $\boldsymbol{\gamma}$ can be carried out in presence of conditionally error heteroskedastic, as shown in Theorem 3. Parts (c) and (d) of Assumption 1 relax the condition that \mathbf{w}_i is identically distributed across i . As we proceed, only β_i , whose distribution is of interest, is assumed to be IID across i , and it is not required for \mathbf{w}_i and u_i to be identically distributed over i .

Remark 5 The high-level conditions in Assumption 1, concerning the convergence in probability of averages such as $\mathbf{Q}_{n,ww} = n^{-1} \sum_{i=1}^n \mathbf{w}_i \mathbf{w}'_i$, can be verified under weak

cross-sectional dependence. Let $f_i = f(\mathbf{w}_i, \beta_i, u_i)$ be a generic function of \mathbf{w}_i, β_i and u_i .¹ Assume that $\sup_i E(f_i^2) < C$, and $\sup_j \sum_{i=1}^n |\text{cov}(f_i, f_j)| < C$, for some fixed $C < \infty$. Then,

$$\text{var} \left(\frac{1}{n} \sum_{i=1}^n f_i \right) \leq \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n |\text{cov}(f_i, f_j)| \leq \frac{1}{n} \sup_j \sum_{i=1}^n |\text{cov}(f_i, f_j)| \leq \frac{C}{n}.$$

By Chebyshev’s inequality, for any $\varepsilon > 0$, we have $M_\varepsilon > \sqrt{C/\varepsilon}$ such that

$$\Pr \left(\sqrt{n} \left| \frac{1}{n} \sum_{i=1}^n [f_i - E(f_i)] \right| > M_\varepsilon \right) \leq \frac{n \text{var} (n^{-1} \sum_{i=1}^n f_i)}{C} \varepsilon \leq \varepsilon,$$

i.e., $n^{-1} \sum_{i=1}^n [f_i - E(f_i)] = O_p(n^{-1/2})$.

Denote $\phi_i = (\beta_i, \boldsymbol{\gamma}')'$ and $\boldsymbol{\phi} = E(\phi_i) = (E(\beta_i), \boldsymbol{\gamma}')'$. Consider the moment condition,

$$E(\mathbf{w}_i y_i) = E(\mathbf{w}_i \mathbf{w}_i') \boldsymbol{\phi}, \tag{2.3}$$

and sum (2.3) over i

$$\frac{1}{n} \sum_{i=1}^n E(\mathbf{w}_i y_i) = \left[\frac{1}{n} \sum_{i=1}^n E(\mathbf{w}_i \mathbf{w}_i') \right] \boldsymbol{\phi}. \tag{2.4}$$

Let $n \rightarrow \infty$, then $\boldsymbol{\phi}$ is identified by

$$\boldsymbol{\phi} = \mathbf{Q}_{ww}^{-1} \mathbf{q}_{wy}, \tag{2.5}$$

under Assumption 1.

Assumption 2 Let $\tilde{y}_i = y_i - \mathbf{z}_i' \boldsymbol{\gamma}$.

- (a) $|n^{-1} \sum_{i=1}^n E(\tilde{y}_i^r x_i^s) - \rho_{r,s}| = O(n^{-1/2})$, and $|\rho_{r,s}| < \infty$, for $r, s = 0, 1, \dots, 2K - 1$.
- (b) $|n^{-1} \sum_{i=1}^n E(u_i^r) - \sigma_r| = O(n^{-1/2})$, and $|\sigma_r| < \infty$, for $r = 2, 3, \dots, 2K - 1$.
- (c) $n^{-1} \sum_{i=1}^n [\text{var}(x_i^r) - (\rho_{0,2r} - \rho_{0,r}^2)] = O(n^{-1/2})$ where $\rho_{0,2r} - \rho_{0,r}^2 > 0$, for $r = 2, 3, \dots, 2K - 1$.

Remark 6 The above assumption allows for a limited degree of heterogeneity of the moments. As an example, let $E(u_i^r) = \sigma_{ir}$ and denote the heterogeneity of the r^{th} moment of u_i by $e_{ir} = \sigma_{ir} - \sigma_r$. Then

$$\left| n^{-1} \sum_{i=1}^n E(u_i^r) - \sigma_r \right| \leq n^{-1} \sum_{i=1}^n |e_{ir}|,$$

¹ f_i is assumed to be a scalar, and we can apply the analysis element-by-element to a matrix, for example $\mathbf{w}_i \mathbf{w}_i'$.

and condition (b) of Assumption 2 is met if $\sum_{i=1}^n |e_{ir}| = O(n^{\alpha_r})$ with $\alpha_r < 1/2$. α_r measures the degree of heterogeneity with $\alpha_r = 1$ representing the highest degree of heterogeneity. A similar idea is used by Pesaran and Zhou (2018) in their analysis of poolability in panel data models.

Theorem 1 *Under Assumptions 1 and 2, $E(\beta_i^r)$ and $\sigma_r, r = 2, 3, \dots, 2K - 1$ are identified.*

Proof For $r = 2, \dots, 2K - 1$,

$$E(\tilde{y}_i^r) = E(x_i^r) E(\beta_i^r) + E(u_i^r) + \sum_{q=2}^{r-1} \binom{r}{q} E(x_i^{r-q}) E(u_i^q) E(\beta_i^{r-q}), \tag{2.6}$$

$$E(\tilde{y}_i^r x_i^r) = E(x_i^{2r}) E(\beta_i^r) + E(x_i^r) E(u_i^r) + \sum_{q=2}^{r-1} \binom{r}{q} E(x_i^{2r-q}) E(u_i^q) E(\beta_i^{r-q}). \tag{2.7}$$

where $\binom{r}{q} = \frac{r!}{q!(r-q)!}$ are binomial coefficients, for nonnegative integers $q \leq r$.
 Sum over i , then by parts (a) and (b) of Assumption 2,

$$\rho_{0,r} E(\beta_i^r) + \sigma_r = \rho_{r,0} - \sum_{q=2}^{r-1} \binom{r}{q} \rho_{0,r-q} \sigma_q E(\beta_i^{r-q}), \tag{2.8}$$

$$\rho_{0,2r} E(\beta_i^r) + \rho_{0,r} \sigma_r = \rho_{r,r} - \sum_{q=2}^{r-1} \binom{r}{q} \rho_{0,2r-q} \sigma_q E(\beta_i^{r-q}). \tag{2.9}$$

Derivation details are relegated to ‘‘Appendix A.1.’’ By part (c) of Assumption 2, the matrix $\begin{pmatrix} \rho_{0,r} & 1 \\ \rho_{0,2r} & \rho_{0,r} \end{pmatrix}$ is invertible for $r = 2, 3, \dots, 2K - 1$. As a result, we can sequentially solve (2.8) and (2.9) for $E(\beta_i^r)$ and σ_r , for $r = 2, 3, \dots, 2K - 1$. \square

2.2 Identifying the distribution of β_i

Beran and Hall (1992, Theorem 2.1, pp. 1972) prove the identification of the distribution of the random coefficient, β_i , in a canonical model without covariates, z_i , under the condition that the distribution of β_i is uniquely determined by its moments. We show the identification of moments of β_i holds more generally when x_i and u_i are not identically distributed and the distribution of β_i is identified if it follows a categorical distribution. Note that under (2.2),

$$E(\beta_i^r) = \sum_{k=1}^K \pi_k b_k^r, \quad r = 0, 1, 2, \dots, 2K - 1, \tag{2.10}$$

with $E(\beta_i^r)$ identified under Assumption 1. To identify $\pi = (\pi_1, \pi_2, \dots, \pi_K)'$ and $\mathbf{b} = (b_1, b_2, \dots, b_K)'$, we need to verify that the system of $2K$ equations in (2.10)

has a unique solution if $b_1 < b_2 < \dots < b_K$, and $\pi_k \in (0, 1)$. In the proof, we construct a linear recurrence relation and make use of the corresponding characteristic polynomial.

Theorem 2 Consider the random coefficient regression model (2.1), suppose that Assumptions 1 and 2 hold. Then $\theta = (\pi', \mathbf{b}')'$ is identified subject to $b_1 < b_2 < \dots < b_K$ and $\pi_k \in (0, 1)$, for all $k = 1, 2, \dots, K$.

Proof We motivate the key idea of the proof in the special case where $K = 2$, and relegate the proof of the general case to the ‘‘Appendix A.1.’’ Let $b_1 = \beta_L, b_2 = \beta_H, \pi_1 = \pi$ and $\pi_2 = 1 - \pi$. Note that

$$E(\beta_i) = \pi\beta_L + (1 - \pi)\beta_H, \tag{2.11}$$

$$E(\beta_i^2) = \pi\beta_L^2 + (1 - \pi)\beta_H^2, \tag{2.12}$$

$$E(\beta_i^3) = \pi\beta_L^3 + (1 - \pi)\beta_H^3, \tag{2.13}$$

and $E(\beta_i^k), k = 1, 2, 3$ are identified. (π, β_L, β_H) can be identified if the system of Eqs. (2.11)–(2.13), has a unique solution. By (2.11),

$$\pi = \frac{\beta_H - E(\beta_i)}{\beta_H - \beta_L}, \text{ and } 1 - \pi = \frac{E(\beta_i) - \beta_L}{\beta_H - \beta_L}. \tag{2.14}$$

Plug (2.14) into (2.12) and (2.13),

$$E(\beta_i)(\beta_L + \beta_H) - \beta_L\beta_H = E(\beta_i^2), \tag{2.15}$$

$$E(\beta_i^2)(\beta_L + \beta_H) - E(\beta_i)\beta_L\beta_H = E(\beta_i^3). \tag{2.16}$$

Denote $\beta_{L+H} = \beta_L + \beta_H$ and $\beta_{LH} = \beta_L\beta_H$, and write (2.15) and (2.16) in matrix form,

$$\mathbf{M}\mathbf{D}\mathbf{b}^* = \mathbf{m}, \tag{2.17}$$

where

$$\mathbf{M} = \begin{pmatrix} 1 & E(\beta_i) \\ E(\beta_i) & E(\beta_i^2) \end{pmatrix}, \mathbf{D} = \begin{pmatrix} -1 & 0 \\ 0 & 1 \end{pmatrix}, \mathbf{b}^* = \begin{pmatrix} \beta_{LH} \\ \beta_{L+H} \end{pmatrix}, \text{ and } \mathbf{m} = \begin{pmatrix} E(\beta_i^2) \\ E(\beta_i^3) \end{pmatrix}.$$

Under the conditions $0 < \pi < 1$ and $\beta_H > \beta_L$,

$$\det(\mathbf{M}) = \text{var}(\beta_i) = E(\beta_i^2) - E(\beta_i)^2 = \pi(1 - \pi)(\beta_H - \beta_L)^2 > 0.$$

As a result, we can solve (2.17) for β_{L+H} and β_{LH} as

$$\beta_{L+H} = \frac{E(\beta_i^3) - E(\beta_i) E(\beta_i^2)}{\text{var}(\beta_i)}, \tag{2.18}$$

$$\beta_{LH} = \frac{E(\beta_i) E(\beta_i^3) - E(\beta_i^2)^2}{\text{var}(\beta_i)}. \tag{2.19}$$

β_L and β_H are solutions to the quadratic equation,

$$\beta^2 - \beta_{L+H}\beta + \beta_{LH} = 0. \tag{2.20}$$

We can verify that $\Delta = \beta_{L+H}^2 - 4\beta_{LH} > 0$ by direct calculation using (2.18) and (2.19). Simplifying Δ in terms of $E(\beta_i^k)$ and then plugging in (2.11), (2.12) and (2.13),

$$\begin{aligned} \Delta &= \frac{[E(\beta_i^3) - E(\beta_i) E(\beta_i^2)]^2 - 4\text{var}(\beta_i) [E(\beta_i) E(\beta_i^3) - E(\beta_i^2)^2]}{[\text{var}(\beta_i)]^2} \\ &= (\beta_H - \beta_L)^2 > 0. \end{aligned}$$

Then, we obtain the unique solutions,

$$\beta_L = \frac{1}{2} \left(\beta_{L+H} - \sqrt{\beta_{L+H}^2 - 4\beta_{LH}} \right), \tag{2.21}$$

$$\beta_H = \frac{1}{2} \left(\beta_{L+H} + \sqrt{\beta_{L+H}^2 - 4\beta_{LH}} \right), \tag{2.22}$$

and π can be determined by (2.14) correspondingly. □

Remark 7 The key identifying assumption in (2) is the assumed existence of the strict ordinal relation $b_1 < b_2 < \dots < b_K$ so that b_k and $b_{k'}$ are not symmetric for $k \neq k'$, and $0 < \pi_k < 1$ so that the distribution of β_i does not degenerate. When $K = 2$, the conditions $b_1 < b_2 < \dots < b_K$, and $\pi_k \in (0, 1)$, are equivalent to $\text{var}(\beta_i) = \pi_1(1 - \pi_1)(b_2 - b_1)^2 > 0$. In other words, not surprisingly, the categorical distribution of β_i is identified only if $\text{var}(\beta_i) > 0$.

In practice, a test for $\mathbb{H}_0 : \text{var}(\beta_i) = 0$ is possible, by noting that $\text{var}(\beta_i) = 0$ is equivalent to

$$\kappa^2 = \frac{E(\beta_i)^2}{E(\beta_i^2)} = 1,$$

where κ^2 is well defined as long as $\beta_i \not\equiv 0$. One important advantage of basing the test of slope homogeneity on κ^2 rather than on $\text{var}(\beta_i) = 0$ is that κ^2 is scale-invariant. $E(\beta_i)$ and $E(\beta_i^2)$ are identified as in Sect. 2.1, whose consistent estimation does not require $\text{var}(\beta_i) > 0$. Consequently, in principle it is possible to test slope homogeneity by testing $\mathbb{H}_0 : \kappa^2 = 1$. However, the problem becomes much more

complicated when there are more than two categories and/or there are more than one regressor under consideration. A full treatment of testing slope homogeneity in such general settings is beyond the scope of the present paper.

Remark 8 Note that in the special case of the proof of Theorem 2 where $K = 2$, $\beta_{L+H} = \beta_L + \beta_H$ and $\beta_{LH} = \beta_L\beta_H$ corresponds to b_1^* and b_2^* , and (2.17) is the same as (A.1.6) when $K = 2$. This special case illustrates the procedure of identification: identify $(b_k^*)_{k=1}^K$ by the moments of β_i , then solve for $(b_k)_{k=1}^K$ and finally identify $(\pi_k)_{k=1}^K$.

3 Estimation

In this section, we propose a generalized method of moments estimator for the distributional parameters of β_i . To reduce the complexity of the moment equations, we first obtain a \sqrt{n} -consistent estimator of $\boldsymbol{\gamma}$ and consider the estimation of the distribution of β_i by replacing $\boldsymbol{\gamma}$ by $\hat{\boldsymbol{\gamma}}$.

3.1 Estimation of $\boldsymbol{\gamma}$

Let $\boldsymbol{\phi} = (E(\beta_i), \boldsymbol{\gamma}')'$, $v_i = \beta_i - E(\beta_i)$ and using the notation in Assumption 1, (2.1) can be written as

$$y_i = \mathbf{w}_i' \boldsymbol{\phi} + \xi_i, \tag{3.1}$$

where $\xi_i = u_i + x_i v_i$. Then, $\boldsymbol{\phi}$ can be estimated consistently by $\hat{\boldsymbol{\phi}} = \mathbf{Q}_{n,ww}^{-1} \mathbf{q}_{n,wy}$ where $\mathbf{Q}_{n,ww}$ and $\mathbf{q}_{n,wy}$ are defined in Assumption 1.

Assumption 3 $\|n^{-1} \sum_{i=1}^n E(\mathbf{w}_i \mathbf{w}_i' \xi_i^2) - \mathbf{V}_{w\xi}\| = O(n^{-1/2})$, $\mathbf{V}_{w\xi} \succ 0$, and

$$\left\| \frac{1}{n} \sum_{i=1}^n \mathbf{w}_i \mathbf{w}_i' \xi_i^2 - \frac{1}{n} \sum_{i=1}^n E(\mathbf{w}_i \mathbf{w}_i' \xi_i^2) \right\| = O_p(n^{-1/2}). \tag{3.2}$$

Remark 9 As in the case of Assumption 1, the high-level condition (3.2) can be shown to hold under weak cross-sectional dependence, assuming that elements of $\mathbf{w}_i \mathbf{w}_i' \xi_i^2$ are cross-sectionally weakly correlated over i . See Remark 5.

Theorem 3 Under Assumption 1, $\hat{\boldsymbol{\phi}}$ is a consistent estimator for $\boldsymbol{\phi}$. In addition, under Assumptions 1 and 3, as $n \rightarrow \infty$,

$$\sqrt{n}(\hat{\boldsymbol{\phi}} - \boldsymbol{\phi}) \rightarrow_d N(\mathbf{0}, \mathbf{V}_\phi), \tag{3.3}$$

where $\mathbf{V}_\phi = \mathbf{Q}_{ww}^{-1} \mathbf{V}_{w\xi} \mathbf{Q}_{ww}^{-1}$. \mathbf{V}_ϕ is consistently estimated by

$$\hat{\mathbf{V}}_\phi = \mathbf{Q}_{n,ww}^{-1} \hat{\mathbf{V}}_{w\xi} \mathbf{Q}_{n,ww}^{-1} \rightarrow_p \mathbf{V}_\phi,$$

as $n \rightarrow \infty$, where $\hat{V}_{w\xi} = n^{-1} \sum_{i=1}^n \mathbf{w}_i \mathbf{w}'_i \hat{\xi}_i^2$, and $\hat{\xi}_i = y_i - \mathbf{w}'_i \hat{\phi}$.

The proof of Theorem 3 is provided in Sect. S.2 in the online supplement.

3.2 Estimation of the distribution of β_i

Denote the moments of β_i on the right-hand side of (2.10) by

$$\begin{aligned} \mathbf{m}_\beta &= (m_1, m_2, \dots, m_{2K-1})' \\ &= [E(\beta_i^r)]_{r=1}^{2K-1} \in \Theta_m \subset \left\{ \mathbf{m}_\beta \in \mathbb{R}^{2K-1} : m_r \geq 0, r \text{ is even} \right\}, \end{aligned}$$

and note that

$$\mathbf{m}_\beta = \begin{pmatrix} m_1 \\ m_2 \\ \vdots \\ m_{2K-1} \end{pmatrix} = \begin{pmatrix} b_1 & b_2 & \dots & b_K \\ b_1^2 & b_2^2 & \dots & b_K^2 \\ \vdots & \vdots & \vdots & \vdots \\ b_1^{2K-1} & b_2^{2K-1} & \dots & b_K^{2K-1} \end{pmatrix} \begin{pmatrix} \pi_1 \\ \pi_2 \\ \vdots \\ \pi_K \end{pmatrix}, \tag{3.4}$$

so in general we can write $\mathbf{m}_\beta \triangleq h(\boldsymbol{\theta})$, where $\boldsymbol{\theta} = (\boldsymbol{\pi}', \mathbf{b}')' \in \Theta$, and $\boldsymbol{\theta}$ can be uniquely determined in terms of \mathbf{m}_β by Theorem 2. To estimate $\boldsymbol{\theta}$, we consider moment conditions following a similar procedure as in Sect. 2 and propose a generalized method of moments (GMM) estimator.

We consider the following moment conditions:

$$E(\tilde{y}_i^r) = \sum_{q=0}^r \binom{r}{q} E(x_i^{r-q}) E(u_i^q) m_{r-q},$$

and

$$E(\tilde{y}_i^r x_i^{s_r}) = \sum_{q=0}^r \binom{r}{q} E(x_i^{r-q+s_r}) E(u_i^q) m_{r-q}, \tag{3.5}$$

where $E(u_i) = 0$, $\tilde{y}_i = y_i - \mathbf{z}'_i \boldsymbol{\gamma}$, $r = 1, 2, \dots, 2K - 1$, and $s_r = 0, 1, \dots, S - r$, where S is a user-specific tuning parameter, chosen such that the highest order moments of x_i included is at most S , where $S > 2K - 1$.²

² For identification, we require the moments of x_i to exist up to order $4K - 2$. S can take values between $2K$ to $4K - 2$. In practice, the choice of S affects the trade-off between bias and efficiency.

Let $\sigma_0 = 1$ and $\sigma_1 = 0$ such that σ_r is well defined for $r = 0, 1, \dots, 2K - 1$. Sum (3.5) over i and rearrange terms,

$$\begin{aligned} 0 &= \sum_{q=0}^r \binom{r}{q} \left[\frac{1}{n} \sum_{i=1}^n E \left(x_i^{r-q+s_r} \right) E \left(u_i^q \right) \right] m_{r-q} - \frac{1}{n} \sum_{i=1}^n E \left(\tilde{y}_i^r x_i^{s_r} \right) \\ &= \sum_{q=0}^r \binom{r}{q} \left[\frac{1}{n} \sum_{i=1}^n E \left(x_i^{r-q+s_r} \right) \right] \sigma_q m_{r-q} - \frac{1}{n} \sum_{i=1}^n E \left(\tilde{y}_i^r x_i^{s_r} \right) + \delta_n^{(r,s_r)}, \end{aligned} \tag{3.6}$$

where

$$\delta_n^{(r,s_r)} = \sum_{q=0}^r \binom{r}{q} \left[\frac{1}{n} \sum_{i=1}^n E \left(x_i^{r-q+s_r} \right) \left[E \left(u_i^q \right) - \sigma_q \right] \right] m_{r-q} = O \left(n^{-1/2} \right),$$

as shown in the proof of Theorem 1.

Letting $n \rightarrow \infty$ in (3.6),

$$\sum_{q=0}^r \binom{r}{q} \rho_{0,r-q+s_r} \sigma_q m_{r-q} - \rho_{r,s_r} = 0, \tag{3.7}$$

by Assumption 2. We stack the left-hand side of (3.7) over $r = 1, 2, \dots, 2K - 1$, and $s_r = 0, 1, \dots, S - r$ and transform $\mathbf{m}_\beta = h(\boldsymbol{\theta})$ to obtain $\mathbf{g}_0(\boldsymbol{\theta}, \boldsymbol{\sigma}, \boldsymbol{\gamma})$.

To implement the GMM estimation, we replace \tilde{y}_i , by $\hat{\tilde{y}}_i = y_i - \mathbf{z}'_i \hat{\boldsymbol{\gamma}}$, and ρ_{r,s_r} by $n^{-1} \sum_{i=1}^n \hat{\tilde{y}}_i^r x_i^{s_r}$. Noting that $\mathbf{m}_\beta = h(\boldsymbol{\theta})$, denote the sample version of the left-hand side of (3.7) by

$$\hat{\mathbf{g}}_n^{(r,s_r)}(\boldsymbol{\theta}, \boldsymbol{\sigma}, \hat{\boldsymbol{\gamma}}) = \frac{1}{n} \sum_{i=1}^n \hat{\mathbf{g}}_i^{(r,s_r)}(\boldsymbol{\theta}, \boldsymbol{\sigma}, \hat{\boldsymbol{\gamma}}), \tag{3.8}$$

where

$$\hat{\mathbf{g}}_i^{(r,s_r)}(\boldsymbol{\theta}, \boldsymbol{\sigma}, \hat{\boldsymbol{\gamma}}) = \sum_{q=0}^r \binom{r}{q} x_i^{r-q+s_r} \sigma_q [h(\boldsymbol{\theta})]_{r-q} - \hat{\tilde{y}}_i^r x_i^{s_r},$$

and $\boldsymbol{\sigma} = (\sigma_2, \sigma_3, \dots, \sigma_{2K-1})'$. Stack the equations in (3.8), over $r = 0, 1, \dots, 2K - 1$ and $s_r = 0, 1, \dots, S - r$ ($S > 2K - 1$), in vector notations we have

$$\hat{\mathbf{g}}_n(\boldsymbol{\theta}, \boldsymbol{\sigma}, \hat{\boldsymbol{\gamma}}) = \frac{1}{n} \sum_{i=1}^n \hat{\mathbf{g}}_i(\boldsymbol{\theta}, \boldsymbol{\sigma}, \hat{\boldsymbol{\gamma}}). \tag{3.9}$$

Given $\hat{\boldsymbol{\gamma}}$, the GMM estimator of $(\boldsymbol{\theta}', \boldsymbol{\sigma}')$ is now computed as

$$(\hat{\boldsymbol{\theta}}', \hat{\boldsymbol{\sigma}}')' = \arg \min_{\boldsymbol{\theta} \in \Theta, \boldsymbol{\sigma} \in \mathcal{S}} \hat{\Phi}_n(\boldsymbol{\theta}, \boldsymbol{\sigma}, \hat{\boldsymbol{\gamma}}),$$

where $\hat{\Phi}_n = \hat{\mathbf{g}}_n(\boldsymbol{\theta}, \boldsymbol{\sigma}, \hat{\boldsymbol{\gamma}})' \mathbf{A}_n \hat{\mathbf{g}}_n(\boldsymbol{\theta}, \boldsymbol{\sigma}, \hat{\boldsymbol{\gamma}})$, and \mathbf{A}_n is a positive definite matrix. We follow the GMM literature using the following choice of \mathbf{A}_n ,

$$\hat{\mathbf{A}}_n = \left[\frac{1}{n} \sum_{i=1}^n \hat{\mathbf{g}}_i(\tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\sigma}}, \hat{\boldsymbol{\gamma}}) \hat{\mathbf{g}}_i(\tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\sigma}}, \hat{\boldsymbol{\gamma}})' - \bar{\mathbf{g}}_n \bar{\mathbf{g}}_n' \right]^{-1}, \tag{3.10}$$

where $\bar{\mathbf{g}}_n = \frac{1}{n} \sum_{i=1}^n \hat{\mathbf{g}}_i(\tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\sigma}}, \hat{\boldsymbol{\gamma}})$, and $\tilde{\boldsymbol{\theta}}$ and $\tilde{\boldsymbol{\sigma}}$ are preliminary estimators.

Assumption 4 Denote the true values of $\boldsymbol{\theta}$, $\boldsymbol{\sigma}$ and $\boldsymbol{\gamma}$ by $\boldsymbol{\theta}_0$, $\boldsymbol{\sigma}_0$ and $\boldsymbol{\gamma}_0$.

- (a) Θ and \mathcal{S} are compact. $\boldsymbol{\theta}_0 \in \text{int}(\Theta)$ and $\boldsymbol{\sigma}_0 \in \text{int}(\mathcal{S})$.
- (b) $\mathbf{A}_n \rightarrow_p \mathbf{A}$ as $n \rightarrow \infty$, where \mathbf{A} is some positive definite matrix.
- (c)

$$\frac{1}{n} \sum_{i=1}^n \left[\hat{y}_i^r x_i^{s_r} - E(\tilde{y}_i^r x_i^{s_r}) \right] = O_p(n^{-1/2}),$$

for $r = 0, 1, 2, \dots, 2K - 1$, $s_r = 0, 1, \dots, S - r$, and $S > 2K - 1$.

Remark 10 Parts (a) and (b) of Assumption 4 are standard regularity conditions in the GMM literature. Part (c) together with Assumption 2 are high-level regularity conditions which allow us to generalize the usual IID assumption and nest the IID data generation process as a special case. The sample analog terms in (c) include $\hat{y}_i = y_i - \mathbf{z}_i' \hat{\boldsymbol{\gamma}}$, instead of the infeasible $\tilde{y}_i = y_i - \mathbf{z}_i' \boldsymbol{\gamma}$. The \sqrt{n} -consistency of $\hat{\boldsymbol{\gamma}}$ shown in Theorem 3 ensures that replacing \tilde{y}_i by \hat{y}_i does not alter the convergence rate.

Theorem 4 Let $\boldsymbol{\eta} = (\boldsymbol{\theta}', \boldsymbol{\sigma}')$ and $\boldsymbol{\eta}_0 = (\boldsymbol{\theta}'_0, \boldsymbol{\sigma}'_0)'$. Under Assumptions 1, 2, and 4, $\hat{\boldsymbol{\eta}} \rightarrow_p \boldsymbol{\eta}_0$ as $n \rightarrow \infty$.

The proof of Theorem 4 is provided in ‘‘Appendix A.1.’’

Assumption 5 Follow the notations as in Assumption 4 and in addition denote $\mathbf{G}(\boldsymbol{\theta}, \boldsymbol{\sigma}, \boldsymbol{\gamma}) = \nabla_{(\boldsymbol{\theta}', \boldsymbol{\sigma}')} \mathbf{g}_0(\boldsymbol{\theta}, \boldsymbol{\sigma}, \boldsymbol{\gamma})$, $\mathbf{G}_0 = \mathbf{G}(\boldsymbol{\theta}_0, \boldsymbol{\sigma}_0, \boldsymbol{\gamma}_0)$, $\mathbf{G}_\gamma(\boldsymbol{\theta}, \boldsymbol{\sigma}, \boldsymbol{\gamma}) = \nabla_\gamma \mathbf{g}_0(\boldsymbol{\theta}, \boldsymbol{\sigma}, \boldsymbol{\gamma})$, $\mathbf{G}_{0,\gamma} = \mathbf{G}_\gamma(\boldsymbol{\theta}_0, \boldsymbol{\sigma}_0, \boldsymbol{\gamma}_0)$.

- (a) $\sqrt{n} \hat{\mathbf{g}}_n(\boldsymbol{\theta}_0, \boldsymbol{\sigma}_0, \boldsymbol{\gamma}_0) \rightarrow_d \boldsymbol{\zeta} \sim N(0, \mathbf{V})$ as $n \rightarrow \infty$.
- (b) $\mathbf{G}'_0 \mathbf{A} \mathbf{G}_0 > 0$.

Remark 11 In Assumption 5, parts (a) is the high-level condition required to ensure the asymptotic normality of $\hat{\mathbf{g}}_n(\boldsymbol{\theta}_0, \boldsymbol{\sigma}_0, \boldsymbol{\gamma}_0)$, which can be verified by Lindeberg central limit theorem under low-level regularity conditions. Part (c) of Assumption 5 represents the full-rank condition on \mathbf{G}_0 , required for identification of $\boldsymbol{\theta}_0$ and $\boldsymbol{\sigma}_0$.

By Theorem 3, we have $\sqrt{n}(\hat{\boldsymbol{\gamma}} - \boldsymbol{\gamma}) \rightarrow_d \zeta_\gamma \sim N(0, V_\gamma)$. The following theorem shows the asymptotic normality of the GMM estimator $\hat{\boldsymbol{\eta}}$.

Theorem 5 *Under Assumptions 1, 3, 4 and 5,*

$$\sqrt{n}(\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}_0) \rightarrow_d (\mathbf{G}'_0 \mathbf{A} \mathbf{G}_0)^{-1} \mathbf{G}'_0 \mathbf{A} (\boldsymbol{\zeta} + \mathbf{G}_{0,\gamma} \zeta_\gamma),$$

as $n \rightarrow \infty$.

The proof of Theorem 5 is provided in ‘‘Appendix A.1.’’

Remark 12 In practice, we estimate the variance of the asymptotic distribution of $\hat{\boldsymbol{\eta}}$ by

$$\hat{\mathbf{V}}_\eta = (\hat{\mathbf{G}}' \hat{\mathbf{A}}_n \hat{\mathbf{G}})^{-1} \hat{\mathbf{G}}' \hat{\mathbf{A}}_n \hat{\mathbf{V}}_\zeta \hat{\mathbf{A}}'_n \hat{\mathbf{G}} (\hat{\mathbf{G}}' \hat{\mathbf{A}}_n \hat{\mathbf{G}})^{-1}, \tag{3.11}$$

where $\hat{\mathbf{G}} = \nabla_{(\boldsymbol{\sigma}', \boldsymbol{\theta}')'} \hat{\mathbf{g}}_n(\hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\sigma}}, \hat{\boldsymbol{\gamma}})$, $\hat{\mathbf{A}}_n$ is given by (3.10), and

$$\hat{\mathbf{V}}_\zeta = \frac{1}{n} \sum_{i=1}^n \boldsymbol{\psi}_{n,i} \boldsymbol{\psi}'_{n,i},$$

where

$$\boldsymbol{\psi}_{n,i} = \hat{\mathbf{g}}_i(\hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\sigma}}, \hat{\boldsymbol{\gamma}}) + \nabla_\gamma \hat{\mathbf{g}}_n(\hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\sigma}}, \hat{\boldsymbol{\gamma}}) \mathbf{L} \mathbf{Q}_{n,ww}^{-1}(\mathbf{w}_i \hat{\xi}_i),$$

and $\mathbf{L} = (\mathbf{0}_{p_z \times 1} \ \mathbf{I}_{p_z})$ is the loading matrix that selects $\boldsymbol{\gamma}$ out of $\boldsymbol{\phi}$.

4 Multiple regressors with random coefficients

One important extension of the regression model (2.1) is to allow for multiple regressors with random coefficients having categorical distribution. With this in mind consider

$$y_i = \mathbf{x}'_i \boldsymbol{\beta}_i + \mathbf{z}'_i \boldsymbol{\gamma} + u_i, \tag{4.1}$$

where the $p \times 1$ vector of random coefficients, $\boldsymbol{\beta}_i \in \mathbb{R}^p$ follows the multivariate distribution³

$$\Pr(\beta_{i1} = b_{1k_1}, \beta_{i2} = b_{2k_2}, \dots, \beta_{ip} = b_{pk_p}) = \pi_{k_1, k_2, \dots, k_p}, \tag{4.2}$$

³ We assume the number of categories K is homogeneous across $j = 1, 2, \dots, p$. This is for notational simplicity, and can be readily generalized to allow for $K_j \neq K_{j'}$ without affecting the main results.

Table 1 Distribution of β_i with $p = 2$ and $K = 2$

	$k_2 = L$	$k_2 = H$
$k_1 = L$	$\pi_{LL} = \Pr(\beta_{i1} = b_{1L}, \beta_{i2} = b_{2L})$	$\pi_{LH} = \Pr(\beta_{i1} = b_{1L}, \beta_{i2} = b_{2H})$
$k_1 = H$	$\pi_{HL} = \Pr(\beta_{i1} = b_{1H}, \beta_{i2} = b_{2L})$	$\pi_{HH} = \Pr(\beta_{i1} = b_{1H}, \beta_{i2} = b_{2H})$

with $k_j \in \{1, 2, \dots, K\}, b_{j1} < b_{j2} < \dots < b_{jK}$, and

$$\sum_{k_1, k_2, \dots, k_p \in \{1, 2, \dots, K\}} \pi_{k_1, k_2, \dots, k_p} = 1.$$

As in Sect. 2, $\boldsymbol{y} \in \mathbb{R}^{pz}$, $\mathbf{w}_i = (\mathbf{x}'_i, \mathbf{z}'_i)'$, $\boldsymbol{\beta}_i \perp \mathbf{w}_i$, $u_i \perp \mathbf{w}_i$, and u_i are independently distributed over i with mean 0.

Example 1 Consider the simple case with $p = 2$ and $K = 2$. For $j = 1, 2$, denote two categories as $\{L, H\}$. The probabilities of four possible combinations of realized $\boldsymbol{\beta}_i$ are summarized in Table 1, where $\pi_{LL} + \pi_{LH} + \pi_{HL} + \pi_{HH} = 1$.

We first identify the moments of $\boldsymbol{\beta}_i$. As in Sect. 2, $\boldsymbol{\phi} = (E(\boldsymbol{\beta}_i)', \boldsymbol{\gamma}')'$ is identified by

$$\boldsymbol{\phi} = \mathbf{Q}_{ww}^{-1} \mathbf{q}_{wy}, \tag{4.3}$$

under Assumption 1. We now consider the identification of the higher-order moments of $\boldsymbol{\beta}_i$ up to the finite order $2K - 1$.

Since \boldsymbol{y} is identified as in (4.3), we treat it as known and let $\tilde{y}_i^r = y_i - \mathbf{z}'_i \boldsymbol{\gamma}$. For $r = 2, 3, \dots, 2K - 1$, consider the moment conditions

$$\begin{aligned} E(\tilde{y}_i^r) &= E[(\mathbf{x}'_i \boldsymbol{\beta}_i + u_i)^r] \\ &= E[(\mathbf{x}'_i \boldsymbol{\beta}_i)^r] + E(u_i^r) + \sum_{s=2}^{r-1} \binom{r}{s} E[(\mathbf{x}'_i \boldsymbol{\beta}_i)^{r-s}] E(u_i^s). \end{aligned} \tag{4.4}$$

Note that $\mathbf{x}'_i \boldsymbol{\beta}_i = \sum_{j=1}^p \beta_{ij} x_{ij}$, and

$$E \left[\left(\sum_{j=1}^p \beta_{ij} x_{ij} \right)^r \right] = \sum_{\sum_{j=1}^p q_j = r} \binom{r}{\mathbf{q}} E \left(\prod_{j=1}^p x_{ij}^{q_j} \right) E \left(\prod_{j=1}^p \beta_{ij}^{q_j} \right),$$

where $\binom{r}{\mathbf{q}} = \frac{r!}{q_1! q_2! \dots q_p!}$, for nonnegative integers r, q_1, \dots, q_p with $r = \sum_{j=1}^p q_j$, denotes the multinomial coefficients. We stack $\prod_{j=1}^p x_{ij}^{q_j}$ with $\mathbf{q} \in$

$\{\mathbf{q} \in \{0, 1, \dots, r\}^p : \sum_{j=1}^p q_j = r\}$ in a vector form by denoting⁴

$$\boldsymbol{\tau}_r(\mathbf{x}_i) = [\varphi(\mathbf{x}_i, \mathbf{q}_1), \varphi(\mathbf{x}_i, \mathbf{q}_2), \dots, \varphi(\mathbf{x}_i, \mathbf{q}_{v_r})]'$$

where $\varphi(\mathbf{x}_i, \mathbf{q}) = \prod_{j=1}^p x_{ij}^{q_j}$ and $v_r = \binom{r+p-1}{p-1}$ is the number of distinct monomials of degree r on the variables $x_{i1}, x_{i2}, \dots, x_{ip}$. Similarly,

$$\boldsymbol{\tau}_r(\boldsymbol{\beta}_i) = [\varphi(\boldsymbol{\beta}_i, \mathbf{q}_1), \varphi(\boldsymbol{\beta}_i, \mathbf{q}_2), \dots, \varphi(\boldsymbol{\beta}_i, \mathbf{q}_{v_r})]'$$

where $\varphi(\boldsymbol{\beta}_i, \mathbf{q}) = \prod_{j=1}^p \beta_{ij}^{q_j}$.

Example 2 Consider $p = 2$ and $r = 2$, we have

$$\boldsymbol{\tau}_2(\mathbf{x}_i) = (x_{i1}^2, x_{i1}x_{i2}, x_{i2}^2)'$$

$$\boldsymbol{\tau}_2(\boldsymbol{\beta}_i) = (\beta_{i1}^2, \beta_{i1}\beta_{i2}, \beta_{i2}^2)'$$

and

$$\begin{aligned} E[(x_{i1}\beta_{i1} + x_{i2}\beta_{i2})^2] &= E(x_{i1}^2) E(\beta_{i1}^2) + 2E(x_{i1}x_{i2}) E(\beta_{i1}\beta_{i2}) + E(x_{i2}^2) E(\beta_{i2}^2) \\ &= [E(x_{i1}^2), E(x_{i1}x_{i2}), E(x_{i2}^2)] \text{diag}[(1, 2, 1)'] [E(\beta_{i1}^2), E(\beta_{i1}\beta_{i2}), E(\beta_{i2}^2)]' \\ &= E[\boldsymbol{\tau}_2(\mathbf{x}_i)]' \boldsymbol{\Lambda}_2 E[\boldsymbol{\tau}_2(\boldsymbol{\beta}_i)], \end{aligned}$$

where $\boldsymbol{\Lambda}_2 = \text{diag}[(1, 2, 1)']$.

Then, the moment condition (4.4) can be written as

$$\begin{aligned} E(\tilde{y}_i^r) &= E[\boldsymbol{\tau}_r(\mathbf{x}_i)]' \boldsymbol{\Lambda}_r E[\boldsymbol{\tau}_r(\boldsymbol{\beta}_i)] + E(u_i^r) \\ &\quad + \sum_{s=2}^{r-1} \binom{r}{s} E[\boldsymbol{\tau}_{r-s}(\mathbf{x}_i)]' \boldsymbol{\Lambda}_{r-s} E[\boldsymbol{\tau}_{r-s}(\boldsymbol{\beta}_i)] E(u_i^s), \end{aligned} \tag{4.5}$$

where $\boldsymbol{\Lambda}_r = \text{diag} \left[\left[\binom{r}{\mathbf{q}} \right]_{\sum_{j=1}^p q_j=r} \right]$ is the $v_r \times v_r$ diagonal matrix of multinomial coefficients. We further consider the moment conditions

$$\begin{aligned} E(\tilde{y}_i^r \boldsymbol{\tau}_r(\mathbf{x}_i)) &= E[\boldsymbol{\tau}_r(\mathbf{x}_i) \boldsymbol{\tau}_r(\mathbf{x}_i)'] \boldsymbol{\Lambda}_r E[\boldsymbol{\tau}_r(\boldsymbol{\beta}_i)] + E[\boldsymbol{\tau}_r(\mathbf{x}_i)] E(u_i^r) \\ &\quad + \sum_{s=2}^{r-1} \binom{r}{s} E[\boldsymbol{\tau}_r(\mathbf{x}_i) \boldsymbol{\tau}_{r-s}(\mathbf{x}_i)'] \boldsymbol{\Lambda}_{r-s} E[\boldsymbol{\tau}_{r-s}(\boldsymbol{\beta}_i)] E(u_i^s), \end{aligned} \tag{4.6}$$

$r = 2, 3, \dots, 2K - 1$. (4.5) and (4.6) reduce to (2.6) and (2.7) when $p = 1$.

⁴ For $\mathbf{x} \in \mathbb{R}^p$, note that $\boldsymbol{\tau}_0(\mathbf{x}) = 1$, $\boldsymbol{\tau}_1(\mathbf{x}) = \mathbf{x}$ and $\boldsymbol{\tau}_2(\mathbf{x}) = \text{vech}(\mathbf{x}\mathbf{x}')$.

- Assumption 6** (a) $\|n^{-1} \sum_{i=1}^n E(\tilde{y}_i^r \boldsymbol{\tau}_s(\mathbf{x}_i)) - \boldsymbol{\rho}_{r,s}\| = O(n^{-1/2})$, and $\|\boldsymbol{\rho}_{r,s}\| < \infty, r, s = 0, 1, \dots, 2K - 1$.
 (b) $\|n^{-1} \sum_{i=1}^n E[\boldsymbol{\tau}_r(\mathbf{x}_i) \boldsymbol{\tau}_s(\mathbf{x}_i)'] - \boldsymbol{\Xi}_{r,s}\| = O(n^{-1/2})$, and $\|\boldsymbol{\Xi}_{r,s}\| < \infty, r, s = 0, 1, \dots, 2K - 1$.
 (c) $|n^{-1} \sum_{i=1}^n E(u_i^r) - \sigma_r| = O(n^{-1/2})$, and $|\sigma_r| < \infty$ for $r = 2, 3, \dots, 2K - 1$.
 (d) $\|n^{-1} \sum_{i=1}^n [\text{var}(\boldsymbol{\tau}_r(\mathbf{x}_i)) - (\boldsymbol{\Xi}_{r,r} - \boldsymbol{\rho}_{0,r} \boldsymbol{\rho}'_{0,r})]\| = O(n^{-1/2})$, where $\boldsymbol{\Xi}_{r,r} - \boldsymbol{\rho}_{0,r} \boldsymbol{\rho}'_{0,r} \succ 0$ for $r = 2, 3, \dots, 2K - 1$.

Theorem 6 For any $\mathbf{q} \in \{\mathbf{q} \in \{0, 1, \dots, r\}^p : \sum_{j=1}^p q_j = r\}$ and $r = 2, 3, \dots, 2K - 1$, $E(\prod_{j=1}^p \beta_{ij}^{q_j})$ and σ_r are identified under Assumptions 1 and 6.

Proof For $r = 2, 3, \dots, 2K - 1$, sum (4.5) and (4.6) over i , go through the same steps as in the proof of Theorem 1, then by Assumptions 6(a) to (c), we have (for $n \rightarrow \infty$)

$$\boldsymbol{\rho}'_{r,0} \boldsymbol{\Lambda}_r E[\boldsymbol{\tau}_r(\boldsymbol{\beta}_i)] + \sigma_r = \boldsymbol{\rho}_{r,0} - \sum_{s=2}^{r-1} \binom{r}{s} \boldsymbol{\rho}_{0,r-s} \boldsymbol{\Lambda}_{r-s} E[\boldsymbol{\tau}_{r-s}(\boldsymbol{\beta}_i)] \sigma_s, \tag{4.7}$$

$$\boldsymbol{\Xi}_{r,r} \boldsymbol{\Lambda}_r E[\boldsymbol{\tau}_r(\boldsymbol{\beta}_i)] + \boldsymbol{\rho}_{0,r} \sigma_r = \boldsymbol{\rho}_{r,r} - \sum_{s=2}^{r-1} \binom{r}{s} \boldsymbol{\Xi}_{r,r-s} \boldsymbol{\Lambda}_{r-s} E[\boldsymbol{\tau}_{r-s}(\boldsymbol{\beta}_i)] \sigma_s. \tag{4.8}$$

Note that

$$\mathbf{M}_r = \begin{pmatrix} \boldsymbol{\Xi}_{r,r} & \boldsymbol{\rho}_{0,r} \\ \boldsymbol{\rho}'_{0,r} & 1 \end{pmatrix} \begin{pmatrix} \boldsymbol{\Lambda}_r & \mathbf{0} \\ \mathbf{0} & 1 \end{pmatrix},$$

is invertible since $\det(\mathbf{M}_r) = \det(\boldsymbol{\Xi}_{r,r} - \boldsymbol{\rho}_{0,r} \boldsymbol{\rho}'_{0,r}) \det(\boldsymbol{\Lambda}_r) > 0$, for $r = 2, 3, \dots, R$, by Assumption 6(d). As a result, we can sequentially solve (4.7) and (4.8) for $E[\boldsymbol{\tau}_r(\boldsymbol{\beta}_i)]$ and σ_r , for $r = 2, 3, \dots, 2K - 1$. \square

We now move from the moments of $\boldsymbol{\beta}_i$ to the distribution of $\boldsymbol{\beta}_i$. We first focus on the identification of the marginal probabilities obtained from (4.2) by averaging out the effects of the other coefficients except for β_{ij} , namely we initially focus on identification of $\lambda_{jk} = \Pr(\boldsymbol{\beta}_i = b_{jk})$, for $k = 1, 2, \dots, K$, and $j = 1, 2, \dots, p$.

Remark 13 Focusing on the marginal distribution of $\boldsymbol{\beta}_i$ is similar to focusing on estimation of partial derivatives in the context of nonparametric estimation, where the curse of dimensionality applies. Consider the estimation of regressing y_i on $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{ip})'$,

$$y_i = F(x_{i1}, x_{i2}, \dots, x_{ip}) + u_i.$$

Then if $F(x_1, x_{i2}, \dots, x_{ip})$ is a homogeneous function (of degree $1/\mu$), then

$$y_i = \sum_{j=1}^p \left(\mu \frac{\partial F(\cdot)}{\partial x_{ij}} \right) x_{ij} + u_i,$$

and under certain conditions we can treat $\mu \frac{\partial F(\cdot)}{\partial x_{ij}} \equiv \beta_{ij}$.

By Theorem 6, $E(\beta_{ij}^r)$ is identified for $r = 1, 2, \dots, 2K - 1$ under Assumptions 1 and 6. By (4.2), we have equations

$$E(\beta_{ij}^r) = \sum_{k=1}^K \lambda_{jk} b_{jk}^r, \tag{4.9}$$

$r = 0, 1, \dots, 2K - 1$, which is of the same form as (2.10) and (3.4). To identify $\lambda_j = (\lambda_{j1}, \lambda_{j2}, \dots, \lambda_{jK})'$ and $\mathbf{b}_j = (b_{j1}, b_{j2}, \dots, b_{jK})'$, we can verify the system of $2K$ equations in (4.9) has a unique solution if $b_{j1} < b_{j2} < \dots < b_{jK}$ and $\lambda_{jk} \in (0, 1)$. The following corollary is a direct application of Theorem 2.

Corollary 7 Consider the model (4.1) and suppose that Assumptions 1 and 6 hold. Then, the parameters $\theta_j = (\lambda'_j, \mathbf{b}'_j)'$ of the marginal distribution of β_i with respect to β_{ij} is identified subject to $b_{j1} < b_{j2} < \dots < b_{jK}$ and $\lambda_{jk} \in (0, 1)$ for $j = 1, 2, \dots, p$.

The problem of identification and estimation of the joint distribution of β_i is subject to the curse of dimensionality. We have $K^p - 1$ probability weights, $\pi_{k_1, k_2, \dots, k_p}$, to be identified in addition to the pK categorical coefficients b_{ij} that are identified by Corollary 7. The number of parameters increases rapidly with p . Even in the simplest case with $K = 2$, the total number of unknown parameters is $2p + 2^p - 1$, which grows exponentially.

Note that the marginal probabilities λ_{jk} are related to the joint distribution by

$$\lambda_{jk} = \sum_{k_1, \dots, k_{j-1}, k_{j+1}, \dots, k_p \in \{1, 2, \dots, K\}} \pi_{k_1, k_2, \dots, k_{j-1}, k, k_{j+1}, \dots, k_p}, \tag{4.10}$$

$k = 1, 2, \dots, K$ and $j = 1, 2, \dots, p$. The number of linearly independent equations in (4.10) is $pK - (p - 1)$.

Example 3 Consider the same setup as in Example 1 with $p = 2$ and $K = 2$. The marginal probabilities are obtained by

$$\begin{aligned} \lambda_{1L} &= \Pr(\beta_{i1} = b_{1L}) = \pi_{LL} + \pi_{LH}, \\ \lambda_{1H} &= \Pr(\beta_{i1} = b_{1H}) = 1 - \lambda_{1L} = \pi_{HL} + \pi_{HH}, \\ \lambda_{2L} &= \Pr(\beta_{i2} = b_{2L}) = \pi_{LL} + \pi_{HL}, \\ \lambda_{2H} &= \Pr(\beta_{i2} = b_{2H}) = 1 - \lambda_{2L} = \pi_{LH} + \pi_{HH}. \end{aligned} \tag{4.11}$$

Note that any equation in (4.11) can be expressed as a linear combination of other three equations, for example $\lambda_{2H} = \lambda_{1L} + \lambda_{1H} - \lambda_{2L}$.

The equations corresponding to the cross-moments, $E\left(\prod_{j=1}^p \beta_{ij}^{q_j}\right)$, are

$$E\left(\prod_{j=1}^p \beta_{ij}^{q_j}\right) = \sum_{k_1, k_2, \dots, k_p \in \{1, 2, \dots, K\}} \left(\prod_{j=1}^p b_{jk_j}^{q_j}\right) \pi_{k_1, k_2, \dots, k_p}, \tag{4.12}$$

for $\mathbf{q} \in \left\{ \mathbf{q} \in \{0, 1, \dots, r-1\}^p : \sum_{j=1}^p q_j = r \right\}$, $r = 2, \dots, 2K - 1$. The linear system (4.12) has

$$\sum_{r=1}^{2K-1} \binom{r+p-1}{p-1} - p(2K-1)$$

equations. Then the total number of equations in (4.10) and (4.12) that can be utilized to identify joint probabilities is $C_r = \sum_{r=1}^{2K-1} \binom{r+p-1}{p-1} - pK$, which is smaller than the number of joint probabilities $K^p - 1$ for large p . When $K = 2$, $C_r < K^p - 1$ for $p \geq 7$.

Identification and estimation of the joint distribution of β_i in the general setting will not be pursued in this paper due to the curse of dimensionality. Instead, we consider special cases, that are empirically relevant, in which identification of the joint distribution of β_i can be readily established. We first consider small p and K , in particular $p = 2$ and $K = 2$ as in Example 1.

Example 4 Consider the same setup as in Example 1 with $p = 2$ and $K = 2$. In addition to (4.11), consider the cross-moment,

$$E(\beta_{i1}\beta_{i2}) = b_{1L}b_{2L}\pi_{LL} + b_{1L}b_{2H}\pi_{LH} + b_{1H}b_{2L}\pi_{HL} + b_{1H}b_{2H}\pi_{HH}. \tag{4.13}$$

Writing (4.11) and (4.13) in matrix form, we have

$$\mathbf{B}\boldsymbol{\pi} = \boldsymbol{\lambda},$$

where

$$\mathbf{B} = \begin{pmatrix} 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 1 & 0 & 1 & 0 \\ b_{1L}b_{2L} & b_{1L}b_{2H} & b_{1H}b_{2L} & b_{1H}b_{2H} \end{pmatrix}, \boldsymbol{\pi} = \begin{pmatrix} \pi_{LL} \\ \pi_{LH} \\ \pi_{HL} \\ \pi_{HH} \end{pmatrix}, \boldsymbol{\lambda} = \begin{pmatrix} \lambda_{1L} \\ \lambda_{1H} \\ \lambda_{2L} \\ E(\beta_{i1}\beta_{i2}) \end{pmatrix}.$$

Note that $E(\beta_{i1}\beta_{i2})$ is identified by Theorem 6, and b_{jk_j} and λ_{jk_j} are identified by Corollary 7, and matrix \mathbf{B} is invertible given that $b_{1L} < b_{1H}$ and $b_{2L} < b_{2H}$ (see ‘‘Appendix A.1’’). As a result, the joint probabilities, $\boldsymbol{\pi}$, are identified.

Remark 14 The argument in Example 4 is applicable for identification of the joint distribution of $(\beta_{ij}, \beta_{i,j'})'$ for $j \neq j'$ when $p > 2$ and $K = 2$.

5 Finite sample properties using Monte Carlo experiments

We examine the finite sample performance of the categorical coefficient estimator proposed in Sect. 3 by Monte Carlo experiments.

5.1 Data generating processes

we generate y_i as

$$y_i = \alpha + x_i\beta_i + z_{i1}\gamma_1 + z_{i2}\gamma_2 + u_i, \text{ for } i = 1, 2, \dots, n, \tag{5.1}$$

with β_i distributed as in (2.2) with $K = 2$, and the parameters π, β_L and β_H .⁵

We draw β_i for each individual i independently by setting $\beta_i = \beta_L$ with probability π and $\beta_i = \beta_H$ with probability $1 - \pi$, through a sequence of independent Bernoulli draws. We consider two sets of parameters in all DGPs, denoted as *high variance* and *low variance* parametrization, respectively,

$$(\pi, \beta_L, \beta_H, E(\beta_i), \text{var}(\beta_i)) = \begin{cases} (0.5, 1, 2, 1.5, 0.25) & (\text{high variance}) \\ (0.3, 0.5, 1.345, 1.0915, 0.15) & (\text{low variance}) \end{cases} \tag{5.2}$$

$\beta_H/\beta_L = 2$ for the *high variance* parametrization, and $\beta_H/\beta_L = 2.69$, for the *low variance* parametrization, which is motivated by the estimates in our empirical illustration in Sect. 6.⁶ The values of $E(\beta_i)$ and $\text{var}(\beta_i)$ are obtained noting that $E(\beta_i) = \pi\beta_L + (1 - \pi)\beta_H$, and $\text{var}(\beta_i) = \pi(1 - \pi)(\beta_H - \beta_L)^2$. The remaining parameters are set as $\alpha = 0.25$, and $\gamma = (1, 1)'$, across DGPs.

We generate the regressors and the error terms as follows.

DGP 1 (Baseline) We first generate $\tilde{x}_i \sim \text{IID}\chi^2(2)$, and then set $x_i = (\tilde{x}_i - 2)/2$ so that x_i has 0 mean and unit variance. The additional regressors, z_{ij} , for $j = 1, 2$ with homogeneous slopes are generated as

$$z_{i1} = x_i + v_{i1} \text{ and } z_{i2} = z_{i1} + v_{i2},$$

with $v_{ij} \sim \text{IID } N(0, 1)$, for $j = 1, 2$. This ensures that the regressors are sufficiently correlated. The error term, u_i , is generated as $u_i = \sigma_i\varepsilon_i$, where σ_i^2 are generated as $0.5(1 + \text{IID}\chi^2(1))$, and $\varepsilon_i \sim \text{IID}N(0, 1)$. Note that ε_i and σ_i^2 are generated independently, and $E(u_i^2) = 1$.

DGP 2 (Categorical x) This setup deviates from the baseline DGP, and allows the distribution of x_i to differ across i . Accordingly, we generate $x_i = (\tilde{x}_{1i} - 2)/2$ where $\tilde{x}_{1i} \sim \text{IID}\chi^2(2)$ for $i = 1, 2, \dots, \lfloor n/2 \rfloor$, and $x_i = (\tilde{x}_{2i} - 2)/4$ where $\tilde{x}_{2i} \sim \text{IID}\chi^2(4)$, for $i = \lfloor n/2 \rfloor + 1, \dots, n$. The additional regressors, z_{ij} , for $j = 1, 2$ with homogeneous slopes are generated as

$$z_{i1} = x_i + v_{i1} \text{ and } z_{i2} = z_{i1} + v_{i2},$$

⁵ A Monte Carlo experiment with $K = 3$ is relegated to Sect. S.3.5 in the online supplement.

⁶ The estimates for β_H/β_L in our empirical analysis range from 1.50 to 2.79.

with $v_{ij} \sim \text{IID } N(0, 1)$, for $j = 1, 2$. The error term u_i is generated the same as in DGP 1.

DGP 3 (Categorical u) We generate x_i and \mathbf{z}_i the same as in DGP 1, but allow the error term u_i to have a heterogeneous distribution over i . For $i = 1, 2, \dots, \lfloor n/2 \rfloor$, we set $u_i = \sigma_i \varepsilon_i$, where $\sigma_i^2 \sim \text{IID } \chi^2(2)$ and $\varepsilon_i \sim \text{IID } N(0, 1)$, and for $i = \lfloor n/2 \rfloor + 1, \dots, n$, we set $u_i = (\tilde{u}_i - 2) / 2$, where $\tilde{u}_i \sim \text{IID } \chi^2(2)$.

We investigate the finite sample performance of the estimator proposed in Sect. 3 across DGP 1 to 3 with *low variance* and *high variance* scenarios.⁷ Details of the computational algorithm used to carry out the Monte Carlo experiments (and the empirical results that follow) are given in Sect. S.5 of the online supplement. An accompanying R package is available at <https://github.com/zhan-gao/ccrm>.

5.2 Summary of the MC results

For each sample size $n = 100, 1000, 2000, 5000, 10,000$ and $100,000$ we run 5000 replications of experiments for DGP 1 (baseline), DGP 2 (categorical x) and DGP 3 (categorical u) with *high variance* and *low variance* parametrization, as set out in (5.2).

We first investigate the finite sample performance of $\hat{\phi}$, as an estimator of $\phi = (E(\beta_i), \boldsymbol{\gamma}')'$. Bias, root mean squared errors (RMSE) for estimation of $E(\beta_i)$, γ_1 and γ_2 , as well as the size of testing of the null values at the 5 percent nominal value are reported in Table 2. In addition, we plot the associated empirical power functions in Figs. 1 and 2, for cases of high and low $\text{var}(\beta_i)$. The results show that $\hat{\phi}$ has very good small sample properties with small bias and RMSEs, with size very close to the nominal value of 5 percent across all DGPs and parametrization, even when sample size is relatively small. The power of the test increases steadily as the sample size increases.

Then, we turn to the GMM estimator for the distributional parameters of β_i proposed in Sect. 3.2. The bias, RMSE and the test size based on the asymptotic distribution given in Theorem 5, for π , β_L and β_H , are reported in Table 3. The empirical power functions are reported in Figs. 3 and 4. The reported results are based on $S = 4$, where $S (> 2K - 1 = 3)$ denotes the highest order of moments of x_i included in estimation.⁸

The upper panel of this table reports the results of the high variance and the lower panel for the low variance parametrization, as set out in (5.2). For all parameters and under all DGPs, the bias and RMSE decline steadily with the sample size as predicted by Theorem 4, and confirm the robustness of the GMM estimates to the heterogeneity in the regressor and the error processes. But for a given sample size, the relative

⁷ We can consider a DGP with conditional heteroskedasticity, in which we follow the baseline DGP and generate the error term as $u_i = x_i \varepsilon_i$, where $\varepsilon_i \sim N(0, 1)$. The least square estimator for ϕ is valid in this setup in terms of estimation and inference, whereas the GMM estimator for the distributional parameters θ breaks down, which is to be expected since we can only identify the first moment of β_i under conditional heteroskedasticity. The results are available on request.

⁸ We also tried estimation based on a larger number of moments (using $S = 5$ and $S = 6$). In the case of current Monte Carlo results, adding more moments does not seem to add much to the precision of the estimates and could be counterproductive when n is not sufficiently large. The results are available in Sect. S.3.1 in the online supplement.

Table 2 Bias, RMSE and size of the least square estimator $\hat{\phi}$

DGP	Sample size n	Baseline		Categorical x		Categorical u		
		Bias	RMSE	Bias	RMSE	Bias	RMSE	Size
<i>High variance: var(β_i) = 0.25</i>								
<i>E(β_i) = 1.5</i>								
	100	-0.0024	0.2035	0.0966	0.2035	-0.0037	0.2035	0.0858
	1000	-0.0017	0.0669	0.0568	0.0657	-0.0002	0.0657	0.0540
	2000	-0.0008	0.0463	0.0512	0.0475	-0.0015	0.0475	0.0534
	5000	-0.0004	0.0301	0.0540	0.0300	-0.0008	0.0300	0.0546
	10,000	0.0002	0.0214	0.0508	0.0212	0.0000	0.0212	0.0510
	100,000	-0.0001	0.0066	0.0472	0.0066	0.0000	0.0066	0.0460
$\gamma_1 = 1$								
	100	-0.0022	0.1571	0.0604	0.1598	-0.0006	0.1598	0.0666
	1000	0.0004	0.0501	0.0496	0.0496	-0.0005	0.0496	0.0508
	2000	0.0003	0.0352	0.0530	0.0350	-0.0004	0.0350	0.0544
	5000	-0.0001	0.0222	0.0470	0.0225	0.0005	0.0225	0.0548
	10,000	-0.0004	0.0157	0.0470	0.0157	0.0002	0.0157	0.0512
	100,000	-0.0001	0.0049	0.0494	0.0049	0.0000	0.0049	0.0468
$\gamma_2 = 1$								
	100	0.0011	0.1115	0.0616	0.1121	0.0016	0.1121	0.0654
	1000	-0.0003	0.0358	0.0558	0.0354	0.0001	0.0354	0.0550
	2000	-0.0001	0.0253	0.0522	0.0246	0.0006	0.0246	0.0502
	5000	0.0000	0.0158	0.0480	0.0159	0.0000	0.0159	0.0570
	10,000	0.0002	0.0111	0.0494	0.0111	-0.0002	0.0111	0.0530
	100,000	0.0001	0.0035	0.0488	0.0034	0.0000	0.0034	0.0446

Table 2 continued

DGP	Sample size n	Baseline			Categorical x			Categorical u		
		Bias	RMSE	Size	Bias	RMSE	Size	Bias	RMSE	Size
<i>Low variance: var(β_i) = 0.15</i>										
$E(\beta_i) = 1.0915$	100	-0.0006	0.1829	0.0810	-0.0023	0.1855	0.0766	-0.0025	0.2094	0.0828
	1000	-0.0005	0.0597	0.0610	0.0005	0.0590	0.0478	-0.0006	0.0670	0.0542
	2000	-0.0002	0.0408	0.0516	-0.0007	0.0427	0.0606	-0.0004	0.0475	0.0544
	5000	-0.0002	0.0264	0.0530	-0.0006	0.0266	0.0480	-0.0005	0.0302	0.0538
	10,000	0.0000	0.0189	0.0546	-0.0002	0.0188	0.0486	-0.0002	0.0208	0.0482
	100,000	-0.0001	0.0059	0.0474	0.0000	0.0059	0.0494	0.0000	0.0068	0.0508
$\gamma_1 = 1$	100	-0.0027	0.1521	0.0614	-0.0001	0.1538	0.0622	0.0014	0.1847	0.0624
	1000	0.0001	0.0480	0.0520	-0.0007	0.0481	0.0542	-0.0003	0.0584	0.0570
	2000	0.0002	0.0338	0.0514	-0.0006	0.0334	0.0512	0.0001	0.0417	0.0572
	5000	-0.0002	0.0213	0.0474	0.0003	0.0216	0.0532	0.0007	0.0257	0.0498
	10,000	-0.0003	0.0150	0.0466	0.0002	0.0152	0.0542	0.0001	0.0183	0.0518
	100,000	-0.0001	0.0047	0.0482	0.0000	0.0047	0.0474	0.0000	0.0057	0.0500
$\gamma_2 = 1$	100	0.0011	0.1081	0.0592	0.0013	0.1079	0.0622	-0.0002	0.1323	0.0674
	1000	-0.0003	0.0345	0.0594	0.0003	0.0342	0.0556	0.0006	0.0409	0.0500
	2000	0.0000	0.0243	0.0534	0.0006	0.0235	0.0450	-0.0001	0.0292	0.0576
	5000	0.0001	0.0152	0.0490	0.0001	0.0152	0.0552	-0.0002	0.0179	0.0470
	10,000	0.0002	0.0106	0.0454	-0.0002	0.0107	0.0528	-0.0002	0.0131	0.0526
	100,000	0.0001	0.0033	0.0442	0.0000	0.0033	0.0448	0.0000	0.0040	0.0486

The data generating process is (5.1). *high variance* and *low variance* parametrization are described in (5.2). “Baseline,” “Categorical x ” and “Categorical u ” refer to DGP 1 to 3 as in Sect. 5.1. Generically, bias, RMSE and size are calculated by $R^{-1} \sum_{r=1}^R (\hat{\theta}^{(r)} - \theta_0)$, $\sqrt{R^{-1} \sum_{r=1}^R (\hat{\theta}^{(r)} - \theta_0)^2}$, and $R^{-1} \sum_{r=1}^R \mathbf{1} \left[\left| \hat{\theta}^{(r)} - \theta_0 \right| / \hat{\sigma}_{\theta}^{(r)} > cv_{0.05} \right]$, respectively, for true parameter θ_0 , its estimate $\hat{\theta}^{(r)}$, the estimated standard error of $\hat{\theta}^{(r)}$, $\hat{\sigma}_{\theta}^{(r)}$, and the critical value $cv_{0.05} = \Phi^{-1}(0.975)$ across $R = 5000$ replications, where $\Phi(\cdot)$ is the cumulative distribution function of standard normal distribution

Table 3 Bias, RMSE and size of the GMM estimator for distributional parameters of β

DGP	Sample size n	Baseline		Categorical x		Categorical u		
		Bias	RMSE	Size	RMSE	Bias	RMSE	Size
<i>High variance: var(β_i) = 0.25</i>								
$\pi = 0.5$								
	100	0.0457	0.2291	0.1737	0.0363	0.2410	0.0235	0.2231
	1000	0.0018	0.1019	0.1308	0.0033	0.1178	-0.0270	0.2033
	2000	0.0017	0.0688	0.1084	0.0015	0.0826	-0.0174	0.1545
	5000	-0.0003	0.0416	0.0936	-0.0015	0.0495	-0.0089	0.1048
	10,000	0.0002	0.0301	0.0774	-0.0006	0.0351	-0.0052	0.0864
	100,000	-0.0001	0.0096	0.0550	0.0002	0.0114	-0.0009	0.0582
$\beta_L = 1$								
	100	0.1415	0.4749	0.2472	0.1099	0.5110	0.1151	0.1820
	1000	0.0207	0.1242	0.1501	0.0200	0.1454	-0.0256	0.1225
	2000	0.0129	0.0819	0.1344	0.0116	0.1007	-0.0094	0.1094
	5000	0.0048	0.0512	0.1052	0.0027	0.0607	-0.0053	0.0850
	10,000	0.0031	0.0365	0.0854	0.0021	0.0428	-0.0020	0.0714
	100,000	0.0002	0.0112	0.0534	0.0007	0.0135	-0.0002	0.0574
$\beta_H = 2$								
	100	-0.0996	0.5609	0.2014	-0.0873	0.6154	-0.1071	0.1866
	1000	-0.0193	0.1407	0.1864	-0.0128	0.1581	-0.0319	0.2093
	2000	-0.0099	0.0893	0.1486	-0.0099	0.1094	-0.0239	0.1673
	5000	-0.0053	0.0519	0.1092	-0.0072	0.0622	-0.0127	0.1156
	10,000	-0.0020	0.0362	0.0878	-0.0033	0.0430	-0.0080	0.0986
	100,000	-0.0005	0.0114	0.0530	-0.0003	0.0134	-0.0017	0.0646

Table 3 continued

DGP	Sample size n	Baseline			Categorical x			Categorical μ		
		Bias	RMSE	Size	Bias	RMSE	Size	Bias	RMSE	Size
<i>Low variance: var(β_i) = 0.15</i>										
$\pi = 0.3$	100	0.2175	0.3084	0.2183	0.2227	0.3187	0.2464	0.2294	0.3157	0.2500
	1000	0.0170	0.1536	0.1873	0.0307	0.1837	0.2063	0.0511	0.2295	0.2493
	2000	0.0014	0.1010	0.1426	0.0105	0.1290	0.1601	0.0181	0.1815	0.2102
	5000	-0.0002	0.0590	0.1084	0.0010	0.0737	0.1158	0.0085	0.1232	0.1468
	10,000	-0.0001	0.0415	0.0894	0.0005	0.0515	0.0928	0.0067	0.0906	0.1046
	100,000	-0.0001	0.0129	0.0594	0.0003	0.0158	0.0536	0.0108	0.0349	0.0776
$\beta_L = 0.5$	100	0.3365	0.5905	0.2426	0.3153	0.6042	0.2432	0.3384	0.6746	0.2005
	1000	0.0352	0.2334	0.1560	0.0290	0.2813	0.1544	0.0131	0.4141	0.1233
	2000	0.0175	0.1414	0.1310	0.0131	0.1835	0.1382	-0.0157	0.2988	0.1037
	5000	0.0085	0.0830	0.1082	0.0041	0.1052	0.1118	-0.0057	0.1798	0.0928
	10,000	0.0055	0.0577	0.0966	0.0031	0.0730	0.0934	0.0019	0.1231	0.0760
	100,000	0.0005	0.0180	0.0596	0.0011	0.0222	0.0582	0.0130	0.0443	0.0962
$\beta_H = 1.345$	100	0.0023	0.4727	0.1377	0.0238	0.5290	0.1453	0.0185	0.6500	0.1461
	1000	-0.0081	0.1265	0.1737	0.0042	0.1621	0.1655	0.0120	0.2353	0.1738
	2000	-0.0092	0.0828	0.1428	-0.0026	0.1045	0.1475	0.0029	0.1607	0.1710
	5000	-0.0048	0.0489	0.1028	-0.0041	0.0586	0.1034	0.0006	0.0970	0.1172
	10,000	-0.0025	0.0340	0.0808	-0.0024	0.0412	0.0942	0.0019	0.0706	0.0958
	100,000	-0.0004	0.0105	0.0486	-0.0002	0.0125	0.0548	0.0073	0.0262	0.0696

The data generating process is (5.1). *high variance* and *low variance* parametrization are described in (5.2). “Baseline,” “Categorical x ” and “Categorical μ ” refer to DGP 1 to 3 as in Sect. 5.1. Generically, bias, RMSE and size are calculated by $R^{-1} \sum_{r=1}^R (\hat{\theta}^{(r)} - \theta_0)$, $\sqrt{R^{-1} \sum_{r=1}^R (\hat{\theta}^{(r)} - \theta_0)^2}$, and $R^{-1} \sum_{r=1}^R \mathbf{1} \left[|\hat{\theta}^{(r)} - \theta_0| / \hat{\sigma}_{\theta}^{(r)} > cv_{0.05} \right]$, respectively, for true parameter θ_0 , its estimate $\hat{\theta}^{(r)}$, the estimated standard error of $\hat{\theta}^{(r)}$, $\hat{\sigma}_{\theta}^{(r)}$, and the critical value $cv_{0.05} = \Phi^{-1}(0.975)$ across $R = 5000$ replications, where $\Phi(\cdot)$ is the cumulative distribution function of standard normal distribution

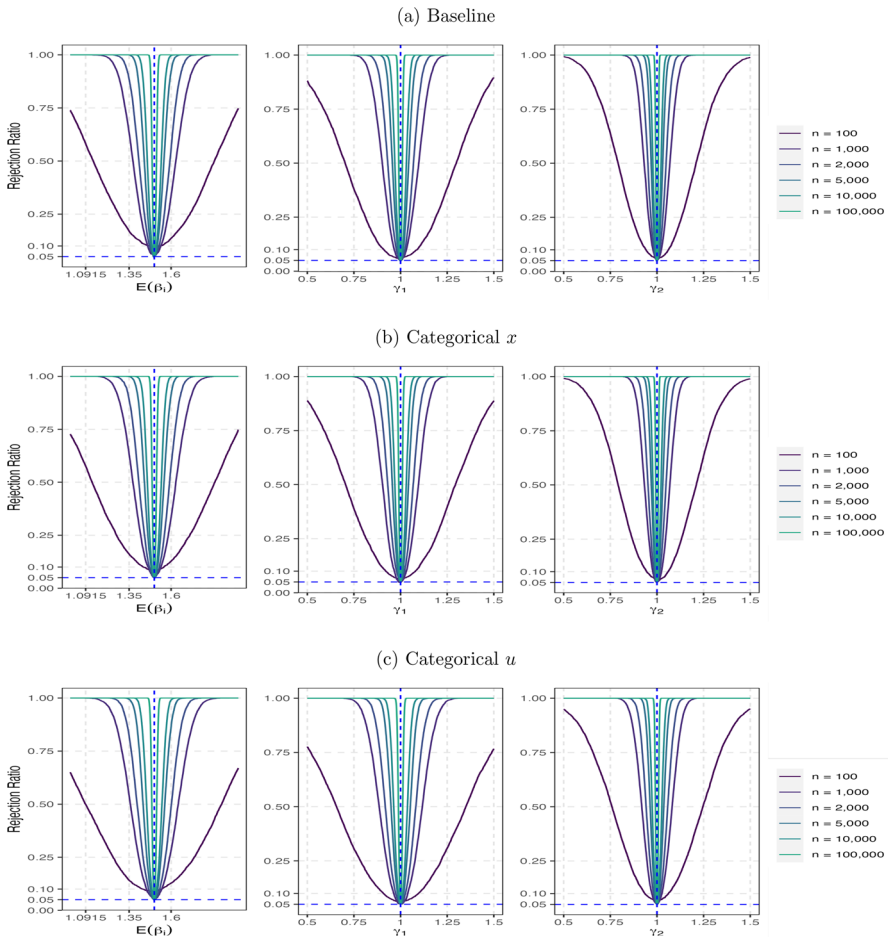


Fig. 1 Empirical power functions for the least square estimator $\hat{\phi}$ with the *high variance* parametrization ($\text{var}(\beta_i) = 0.25$). *Notes:* The data generating process is (5.1) with *high variance* parametrization that is described in (5.2). “Baseline,” “Categorical x ” and “Categorical u ” refer to DGP 1 to 3 as in Sect. 5.1. Generically, power is calculated by $R^{-1} \sum_{r=1}^R \mathbf{1} \left[\left| \hat{\theta}^{(r)} - \theta_{\delta} \right| / \hat{\sigma}_{\hat{\theta}}^{(r)} > cv_{0.05} \right]$, for θ_{δ} in a symmetric neighborhood of the true parameter θ_0 , the estimate at the r -th replication, $\hat{\theta}^{(r)}$, the estimated standard error of $\hat{\theta}^{(r)}$, $\hat{\sigma}_{\hat{\theta}}^{(r)}$, and the critical value $cv_{0.05} = \Phi^{-1}(0.975)$ across $R = 5000$ replications, where $\Phi(\cdot)$ is the cumulative distribution function of standard normal distribution

precision of the estimates depends on the variability of β_i , as characterized by the true value of $\text{var}(\beta_i)$. The precision of the estimates with *high variance* parametrization is relatively higher than that with *low variance* parametrization. This is to be expected since, unlike $E(\beta_i)$, the distributional parameters are only identified if $\text{var}(\beta_i) > 0$. As shown in (2.18) and (2.19) for the current case of $K = 2$, $\text{var}(\beta_i)$ is in the denominator when we recover the distributional parameters from the moments of β_i . When $\text{var}(\beta_i)$ is small, estimation errors in the moments of β_i can be amplified in the estimation of π , β_L and β_H . On the other hand, the larger the variance the more precisely π ,

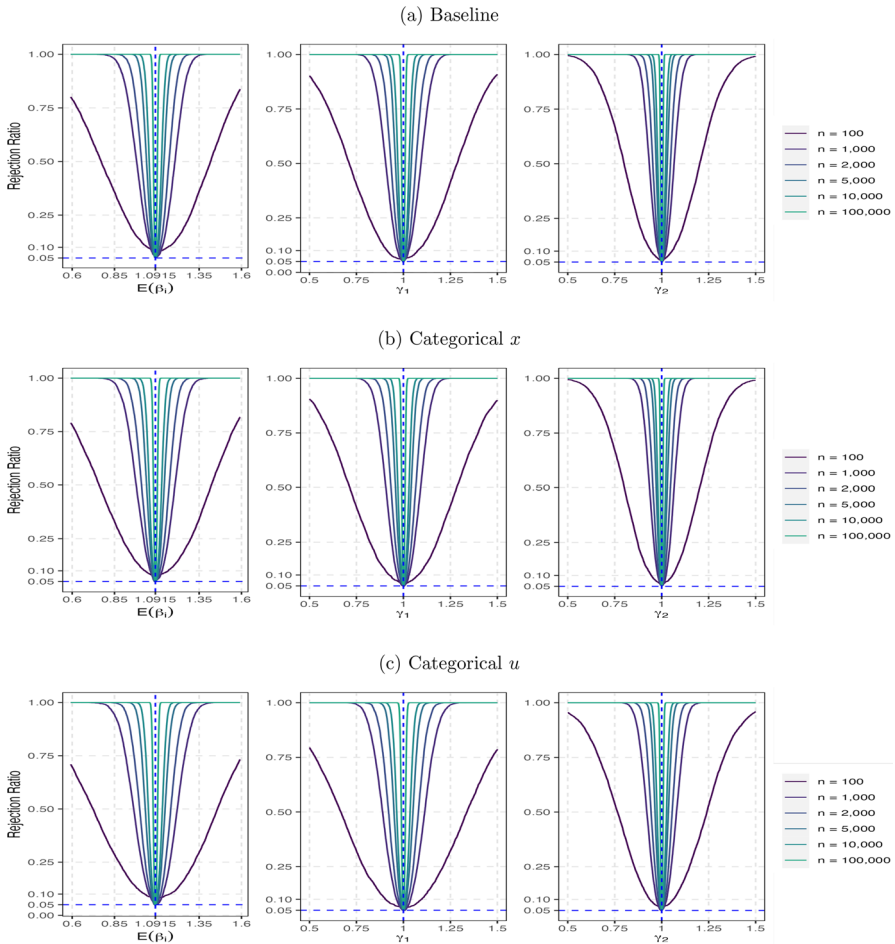


Fig. 2 Empirical power functions for the least square estimator $\hat{\phi}$ with the *low variance* parametrization ($\text{var}(\beta_i) = 0.15$). *Notes:* The data generating process is (5.1) with *low variance* parametrization that is described in (5.2). “Baseline,” “Categorical x ” and “Categorical u ” refer to DGP 1 to 3 as in Sect. 5.1. Generally, power is calculated by $R^{-1} \sum_{r=1}^R \mathbf{1} \left[\left| \hat{\theta}^{(r)} - \theta_{\delta} \right| / \hat{\sigma}_{\hat{\theta}}^{(r)} > \text{cv}_{0.05} \right]$, for θ_{δ} in a symmetric neighborhood of the true parameter θ_0 , the estimate at the r -th replication, $\hat{\theta}^{(r)}$, the estimated standard error of $\hat{\theta}^{(r)}$, $\hat{\sigma}_{\hat{\theta}}^{(r)}$, and the critical value $\text{cv}_{0.05} = \Phi^{-1}(0.975)$ across $R = 5000$ replications, where $\Phi(\cdot)$ is the cumulative distribution function of standard normal distribution

β_H and β_L can be estimated for a given n .⁹ The size and power also depends on the parametrization. With both *high variance* and *low variance* parametrization, we can achieve correct size and reasonable power when n is quite large ($n = 100,000$). We plot the empirical power functions for $n \geq 5000$ for π , β_H and β_L since the size is far

⁹ Section S.3.4 in the online supplement presents parametrization with $\text{var}(\beta_i) = 6.35$ and 18.95 , which further confirms the pattern that the larger the variance the more precisely π , β_H and β_L can be estimated for a given n .

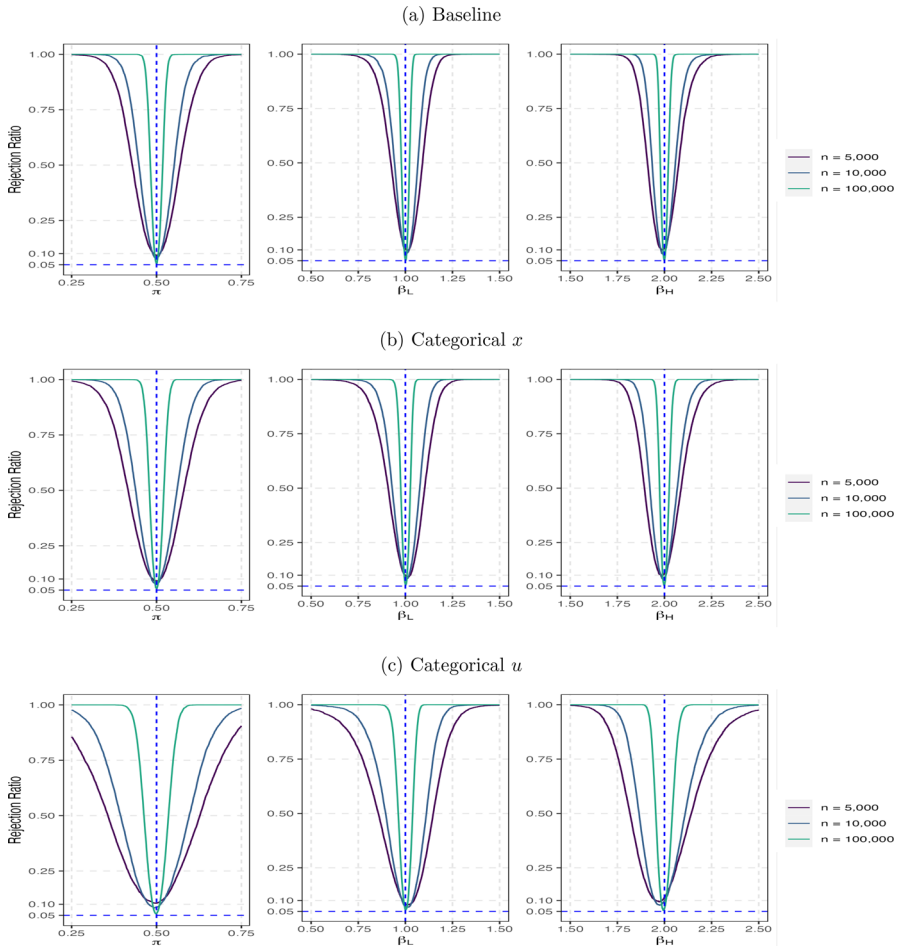


Fig. 3 Empirical power functions for the GMM estimator of distributional parameters of β with the *high variance* parametrization ($\text{var}(\beta_i) = 0.25$). *Notes:* The data generating process is (5.1) with *high variance* parametrization that is described in (5.2). “Baseline,” “Categorical x ” and “Categorical u ” refer to DGP 1 to 3 as in Sect. 5.1. The model is estimated with $S = 4$, the highest order of moments of x_i used in estimation. Generically, power is calculated by $R^{-1} \sum_{r=1}^R \mathbf{1} \left[\left| \hat{\theta}^{(r)} - \theta_\delta \right| / \hat{\sigma}_{\hat{\theta}}^{(r)} > cv_{0.05} \right]$, for θ_δ in a symmetric neighborhood of the true parameter θ_0 , the estimate at the r -th replication, $\hat{\theta}^{(r)}$, the estimated standard error of $\hat{\theta}^{(r)}$, $\hat{\sigma}_{\hat{\theta}}^{(r)}$, and the critical value $cv_{0.05} = \Phi^{-1}(0.975)$ across $R = 5000$ replications, where $\Phi(\cdot)$ is the cumulative distribution function of standard normal distribution

above 5 percent for smaller values of n , and power comparisons are not meaningful in such cases.

Remark 15 Note that GMM estimators of moments of β_i , namely \mathbf{m}_β , can be obtained using the moment conditions in (3.7), and the transformations $\mathbf{m}_\beta = h(\theta)$ in (3.4) are required only to derive the estimators of θ , the parameters of the underlying categorical distribution. The Monte Carlo results in Sect. S.3.2 in the online supplement show that

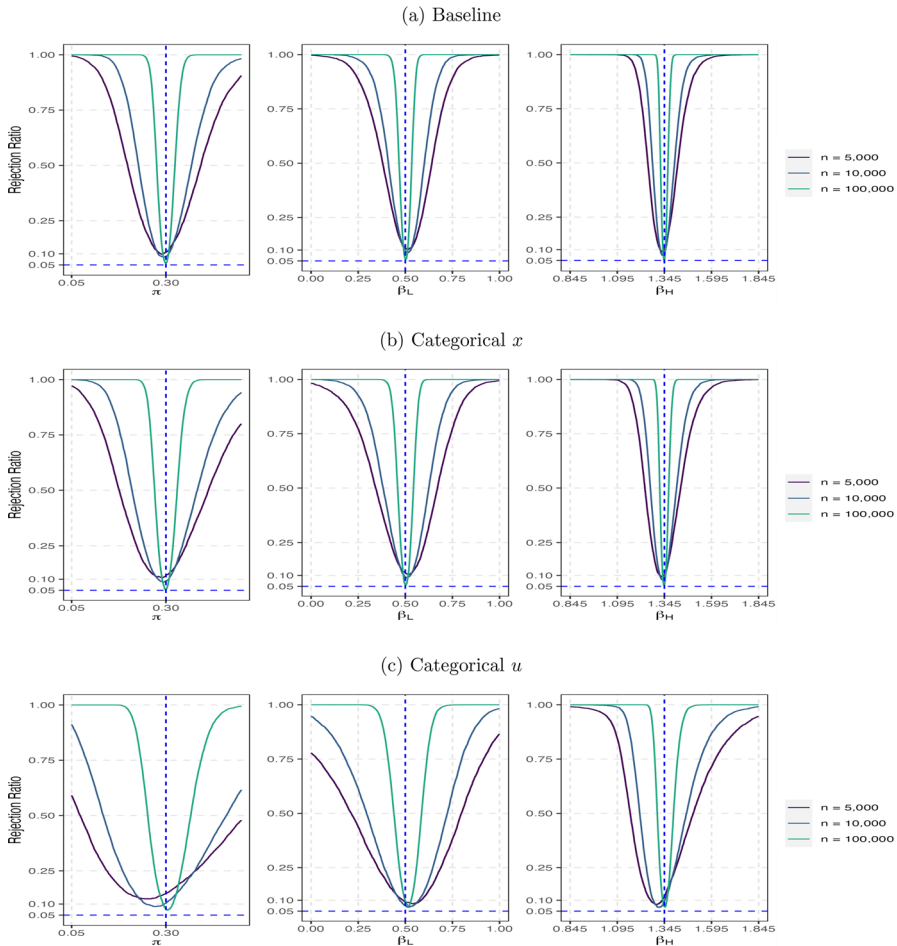


Fig. 4 Empirical power functions for the GMM estimator of distributional parameters of β with the low variance parametrization ($\text{var}(\beta_i) = 0.15$). *Notes:* The data generating process is (5.1) with low variance parametrization that is described in (5.2). “Baseline,” “Categorical x ” and “Categorical u ” refer to DGP 1 to 3 as in Sect. 5.1. The model is estimated with $S = 4$, the highest order of moments of x_i used in estimation. Generically, power is calculated by $R^{-1} \sum_{r=1}^R \mathbf{1} \left[\left| \hat{\theta}^{(r)} - \theta_\delta \right| / \hat{\sigma}_{\hat{\theta}}^{(r)} > \text{cv}_{0.05} \right]$, for θ_δ in a symmetric neighborhood of the true parameter θ_0 , the estimate at the r -th replication, $\hat{\theta}^{(r)}$, the standard error of $\hat{\theta}^{(r)}$, $\hat{\sigma}_{\hat{\theta}}^{(r)}$, and the critical value $\text{cv}_{0.05} = \Phi^{-1}(0.975)$ across $R = 5,000$ replications, where $\Phi(\cdot)$ is the cumulative distribution function of standard normal distribution

\mathbf{m}_β can be accurately estimated with relatively small sample sizes. In the estimation of both \mathbf{m}_β and θ , the same set of moment conditions are included, so the estimation of distributional parameters θ essentially relies on the relation $\theta = h^{-1}(\mathbf{m}_\beta)$. Sampling uncertainties in the estimation of \mathbf{m}_β , particularly in higher-order moments, are potentially amplified through the inverse transformation h^{-1} that involves matrix inversion, which causes the difficulties in estimation and inference of θ when sample sizes are small. This is analogous to the problem of precision matrix estimation from

an estimated covariance matrix. In practice, estimation of the categorical parameters is recommended for applications where the sample size is relatively large, otherwise it is advisable to focus on estimates of the lower-order moments of β_i .

6 Heterogeneous return to education: an empirical application

Since the pioneering work by Becker (1962, 1964) on the effects of investments in human capital, estimating returns to education has been one of the focal points of labor economics research. In his pioneering contribution Mincer (1974) models the logarithm of earnings as a function of years of education and years of potential labor market experience (age minus years of education minus six), which can be written in a generic form:

$$\log \text{wage}_i = \alpha_i + \beta_i \text{edu}_i + \phi(\mathbf{z}_i) + \varepsilon_i, \quad (6.1)$$

as in Heckman et al. (2018, Eq. (1)), where \mathbf{z}_i includes the labor market experience and other relevant control variables. The above wage equation, also known as the ‘‘Mincer equation’’, has become of the workhorse of the empirical works on estimating the return to education. In the most widely used specification of the Mincer equation (6.1),

$$\phi(\mathbf{z}_i) = \rho_1 \text{exper}_i + \rho_2 \text{exper}_i^2 + \tilde{\mathbf{z}}_i' \tilde{\boldsymbol{\gamma}},$$

where $\tilde{\mathbf{z}}_i$ is the vector of control variables other than potential labor market experience.

Along with the advancement of empirical research on this topic, there has been a growing awareness of the importance of heterogeneity in individual cognitive and non-cognitive abilities (Heckman 2001) and their significance for explaining the observed heterogeneity in return to education. Accordingly, it is important to allow the parameters of the wage equation to differ across individuals. In Eq. (6.1), we allow α_i and β_i to differ across individuals, but assume that $\phi(\mathbf{z}_i)$ can be approximated as nonlinear functions of experience and other control variables with homogeneous coefficients.

Specifically, following Lemieux (2006b, c) we also allow for time variations in the parameters of the wage equation and consider the following categorical coefficient model over a given cross-section sample indexed by t ¹⁰:

$$\log \text{wage}_{it} = \alpha_{it} + \beta_{it} \text{edu}_{it} + \rho_{1t} \text{exper}_{it} + \rho_{2t} \text{exper}_{it}^2 + \tilde{\mathbf{z}}_{it}' \tilde{\boldsymbol{\gamma}}_t + \varepsilon_{it}, \quad (6.2)$$

where the return to education follows the categorical distribution,

$$\beta_{it} = \begin{cases} b_{tL} & \text{w.p. } \pi_t, \\ b_{tH} & \text{w.p. } 1 - \pi_t, \end{cases}$$

¹⁰ Some investigators have suggested including higher powers of the experience variable in the wage equation. Lemieux (2006a), for example, proposes using a quartic rather than a quadratic function. As a robustness check we also provide estimation results with quartic experience specification in Sect. S.4 in the online supplement.

and $\tilde{\mathbf{z}}_{it}$ includes gender, marital status and race. $\alpha_{it} = \alpha_t + \delta_{it}$ where δ_{it} is mean 0 random variable assumed to be distributed independently of edu_{it} and $\mathbf{z}_{it} = (\text{exper}_{it}, \text{exper}_{it}^2, \tilde{\mathbf{z}}_t)'$. Let $u_{it} = \varepsilon_{it} + \delta_{it}$, and write (6.2) as

$$\log \text{wage}_{it} = \alpha_t + \beta_{it} \text{edu}_{it} + \rho_{1t} \text{exper}_{it} + \rho_{2t} \text{exper}_{it}^2 + \tilde{\mathbf{z}}_{it}' \tilde{\boldsymbol{\gamma}}_t + u_{it}. \tag{6.3}$$

The correlation between α_{it} and edu_{it} in (6.1) is the source of “ability bias” (Griliches 1977). Given the pure cross-sectional nature of our analysis, we do not allow for the endogeneity from “ability bias” or dynamics. To allow for nonzero correlations between α_{it} , edu_{it} and \mathbf{z}_{it} , a panel data approach is required, which has its own challenges, as education and experience variables tend to very slow moving (if at all) for many individuals in the panel. Time delays between changes in education and experience and the wage outcomes also further complicate the interpretation of the mean estimates of β_{it} which we shall be reporting. To partially address the possible dynamic spillover effects, we provide estimates of the distribution of β_{it} using cross-sectional data from two different sample periods, and investigate the extent to which the distribution of return to education has changed over time, by gender and the level of educational achievements.¹¹

We estimate the categorical distribution of the return to education in (6.3) using the May and Outgoing Rotation Group (ORG) supplements of the Current Population Survey (CPS) data, as in Lemieux (2006b, c).¹² We pool observations from 1973 to 1975 for the first sample period, $t = \{1973-1975\}$ and observations from 2001 to 2003 for the second sample period, $t = \{2001-2003\}$. Following Lemieux (2006b), we consider subsamples of those with less than 12 years of education, “high school or less,” and those with more than 12 years of education, “postsecondary education,” as well as the combined sample. We also present results by gender. The summary statistics are reported in Table 4. As to be expected, the mean log wages are higher for those with postsecondary education (for male and female), with the number of years of schooling and experience rising by about one year across the two sub-period samples. There are also important differences across male and female, and the two educational groupings, which we hope to capture in our estimation.

We treat the cross-section observations in the two sample periods, $t = \{1973-1975\}$ and $\{2001-2003\}$, as *repeated* cross sections, rather than a panel data since the data in these two periods do not cover the same individuals, and represent random samples from the population of wage earners in two periods. It should also be noted that sample sizes (n_t), although quite large, are much larger during $\{2001-2003\}$, which could be a factor when we come to compare estimates from the two sample periods. For example, for both male and female $n_{73-75} = 111,632$ as compared to $n_{01-03} = 511,819$, a difference which becomes more pronounced when we consider the number observations in postsecondary/female category—which rises from 12,882 for the first period to 100,007 in the second period.

¹¹ Time variations in return to education have also been investigated in the literature as a possible explanation of increasing wage inequality in the USA. See, for example, the papers by Lemieux (2006b, c).

¹² The data are retrieved from <https://www.openicpsr.org/openicpsr/project/116216/version/V1/view>.

Table 4 Summary statistics of the May and outgoing rotation group (ORG) supplements of the current population survey (CPS) data across two periods, 1973–1975 and 2001–2003, by years of education and gender

	1973–1975		2001–2003		All	
	High school or less	Postsecondary education	High school or less	Postsecondary education		
<i>Both male and female</i>						
log wage	1.59 (0.50)	1.94 (0.53)	1.47 (0.47)	1.88 (0.57)	1.69 (0.53)	1.71 (0.57)
edu.	10.64 (2.11)	15.21 (1.65)	11.29 (1.68)	14.96 (1.82)	12.02 (2.89)	13.41 (2.53)
age	36.74 (13.85)	34.90 (11.58)	37.96 (12.93)	39.87 (11.33)	36.18 (13.23)	39.06 (12.07)
expr.	20.10 (14.44)	13.69 (11.41)	20.67 (12.95)	18.91 (11.17)	18.17 (13.91)	19.65 (11.98)
marriage	0.67 (0.47)	0.70 (0.46)	0.52 (0.50)	0.60 (0.49)	0.68 (0.47)	0.57 (0.50)
nonwhite	0.11 (0.32)	0.08 (0.27)	0.15 (0.36)	0.14 (0.35)	0.10 (0.30)	0.15 (0.35)
<i>n</i>	77,899	33,733	216,136	295,683	111,632	511,819

Table 4 continued

	1973–1975			2001–2003		
	High school or less	Postsecondary education	All	High school or less	Postsecondary education	All
<i>Male</i>						
log wage	1.76 (0.48)	2.07 (0.53)	1.86 (0.52)	1.57 (0.48)	2.00 (0.58)	1.81 (0.58)
edu.	10.44 (2.26)	15.29 (1.69)	12.00 (3.08)	11.19 (1.82)	15.02 (1.84)	13.31 (2.64)
age	36.79 (13.82)	35.29 (11.24)	36.31 (13.07)	37.21 (12.70)	40.24 (11.30)	38.89 (12.04)
expr.	20.35 (14.49)	14.00 (11.06)	18.32 (13.81)	20.02 (12.75)	19.22 (11.08)	19.58 (11.86)
marriage	0.73 (0.44)	0.76 (0.43)	0.74 (0.44)	0.53 (0.50)	0.64 (0.48)	0.59 (0.49)
nonwhite	0.10 (0.30)	0.06 (0.24)	0.09 (0.29)	0.14 (0.34)	0.13 (0.33)	0.13 (0.34)
<i>n</i>	44,299	20,851	65,150	116,129	144,138	260,267
<i>Female</i>						
log wage	1.35 (0.41)	1.71 (0.47)	1.45 (0.46)	1.77 (0.54)	1.36 (0.43)	1.61 (0.54)
edu.	10.89 (1.87)	15.08 (1.59)	12.05 (2.60)	14.90 (1.79)	11.42 (1.49)	13.52 (2.40)

Table 4 continued

	1973–1975		2001–2003	
	High school or less	All	High school or less	All
age	36.67 (13.88)	34.27 (12.09)	38.83 (13.14)	39.52 (11.35)
expr.	19.78 (14.36)	13.19 (11.94)	18.61 (11.24)	21.41 (13.13)
marriage	0.60 (0.49)	0.60 (0.49)	0.56 (0.50)	0.51 (0.50)
nonwhite	0.13 (0.33)	0.10 (0.30)	0.15 (0.36)	0.17 (0.38)
<i>n</i>	33,600	12,882	151,545	100,007

Postsecondary Education” stands for the subsample with years of education higher than 12 and “High School or Less” stands for subsample with years of education less than or equal to 12). edu., and expr. are in years. marriage and nonwhite are dummy variables. *n* is the sample size. We report mean and standard deviation (in parentheses) of each variable. The data are from the May and Outgoing Rotation Group (ORG) supplements of the Current Population Survey (CPS) data retrieved from <https://www.openicpsr.org/openicpsr/project/116216/version/V1/view>

Table 5 Estimates of the distribution of the return to education across two periods, 1973–1975 and 2001–2003, by years of education and gender

	High school or less		Postsecondary edu.		All	
	1973–1975	2001–2003	1973–1975	2001–2003	1973–1975	2001–2003
<i>Both male and female</i>						
π	0.4843 (4188.8)	0.5069 (0.0269)	0.4398 (0.0502)	0.3537 (0.0091)	0.4719 (0.0485)	0.3463 (0.0047)
β_L	0.0608 (5.0939)	0.0382 (0.0014)	0.0624 (0.0035)	0.0866 (0.0009)	0.0558 (0.0020)	0.0645 (0.0004)
β_H	0.0619 (4.8132)	0.0920 (0.0019)	0.1103 (0.0032)	0.1401 (0.0007)	0.0941 (0.0022)	0.1263 (0.0004)
β_H / β_L	1.0194 (6.2938)	2.4102 (0.0428)	1.7680 (0.0618)	1.6178 (0.0111)	1.6879 (0.0295)	1.9567 (0.0080)
$E(\beta_i)$	0.0614	0.0647	0.0893	0.1212	0.0760	0.1049
s.d. (β_i)	0.0006	0.0269	0.0238	0.0256	0.0191	0.0294
n	77,899	216,136	33,733	295,683	111,632	511,819
<i>Male</i>						
π	N/a	0.4939 (0.0399)	0.4706 (0.0707)	0.3201 (0.0104)	0.4802 (0.0815)	0.3290 (0.0053)
β_L	0.0637	0.0404 (0.0019)	0.0534 (0.0046)	0.0743 (0.0012)	0.0536 (0.0030)	0.0548 (0.0005)
β_H	0.0637	0.0911 (0.0026)	0.0995 (0.0042)	0.1308 (0.0009)	0.0875 (0.0031)	0.1192 (0.0005)
β_H / β_L	1.0000	2.2526 (0.0534)	1.8641 (0.1038)	1.7603 (0.0209)	1.6312 (0.0459)	2.1772 (0.0144)
$E(\beta_i)$	0.0637	0.0661	0.0778	0.1128	0.0712	0.0980
s.d. (β_i)	0.0000	0.0253	0.0230	0.0264	0.0169	0.0303
n	44,299	116,129	20,851	144,138	65,150	260,267

Table 5 continued

	High school or less		Postsecondary edu.		All	
	1973–1975	2001–2003	1973–1975	2001–2003	1973–1975	2001–2003
<i>Female</i>						
π	0.4999 (0.5047)	0.5166 (0.0283)	0.4526 (0.0829)	0.3906 (0.0167)	0.4566 (0.0810)	0.3608 (0.0086)
β_L	0.0441 (0.0133)	0.0348 (0.0016)	0.0823 (0.0053)	0.0979 (0.0013)	0.0628 (0.0033)	0.0751 (0.0007)
β_H	0.0723 (0.0159)	0.0972 (0.0025)	0.1310 (0.0055)	0.1473 (0.0011)	0.1028 (0.0038)	0.1333 (0.0007)
β_H / β_L	1.6392 (0.1565)	2.7934 (0.0700)	1.5913 (0.0539)	1.5048 (0.0121)	1.6357 (0.0353)	1.7756 (0.0090)
$E(\beta_i)$	0.0582	0.0650	0.1090	0.1280	0.0845	0.1123
s.d. (β_i)	0.0141	0.0312	0.0242	0.0241	0.0199	0.0280
n	33,600	100,007	12,882	151,545	46,482	251,552

This table reports the estimates of the distribution of β_i with the quadratic in experience specification (6.2), using $S = 4$ order moments of edu., "Postsecondary Edu." stands for the subsample with years of education higher than 12 and "High School or Less" stands for those with years of education less than or equal to 12. s.d. (β_i) corresponds to the square root of estimated var (β_i). n is the sample size. "n/a" is inserted when the estimates show homogeneity of β_i and π is not identified and cannot be estimated. Estimated standard errors are reported in parentheses.

Table 6 Estimates of γ associated with control variables z_i with specification (6.2) across two periods, 1973–1975 and 2001–2003, by years of education and gender, which complements Table 5

	High school or less		Postsecondary Edu.		All	
	1973–1975	2001–2003	1973–1975	2001–2003	1973–1975	2001–2003
<i>Both male and female</i>						
exper.	0.0305 (0.0004)	0.0319 (0.0002)	0.0415 (0.0008)	0.0354 (0.0003)	0.0310 (0.0003)	0.0321 (0.0002)
exper. ² ($\times 10^2$)	-0.0490 (0.0009)	-0.0505 (0.0005)	-0.0826 (0.0022)	-0.0652 (0.0007)	-0.0499 (0.0008)	-0.0537 (0.0005)
marriage	0.1120 (0.0036)	0.0751 (0.0020)	0.0886 (0.0059)	0.0770 (0.0020)	0.1085 (0.0031)	0.0818 (0.0014)
nonwhite	-0.0922 (0.0047)	-0.0775 (0.0024)	-0.0424 (0.0088)	-0.0571 (0.0025)	-0.0715 (0.0042)	-0.0667 (0.0018)
gender	0.4157 (0.0029)	0.2298 (0.0017)	0.2962 (0.0050)	0.2023 (0.0018)	0.3892 (0.0025)	0.2167 (0.0013)
<i>n</i>	77,899	216,136	33,733	295,683	111,632	511,819
<i>Male</i>						
exper.	0.0369 (0.0005)	0.0366 (0.0003)	0.0516 (0.0011)	0.0405 (0.0005)	0.0389 (0.0005)	0.0371 (0.0003)
exper. ² ($\times 10^2$)	-0.0589 (0.0012)	-0.0589 (0.0008)	-0.1016 (0.0029)	-0.0752 (0.0011)	-0.0635 (0.0010)	-0.0629 (0.0007)
marriage	0.1940 (0.0053)	0.1123 (0.0028)	0.1497 (0.0085)	0.1344 (0.0031)	0.1828 (0.0045)	0.1316 (0.0021)
nonwhite	-0.1241 (0.0065)	-0.1165 (0.0035)	-0.1172 (0.0127)	-0.1010 (0.0039)	-0.1178 (0.0058)	-0.1093 (0.0027)
<i>n</i>	44,299	116,129	20,851	144,138	65,150	260,267

Table 6 continued

	High school or less		Postsecondary Edu.		All	
	1973–1975	2001–2003	1973–1975	2001–2003	1973–1975	2001–2003
<i>Female</i>						
exper.	0.0223 (0.0006)	0.0265 (0.0003)	0.0271 (0.0011)	0.0313 (0.0004)	0.0208 (0.0005)	0.0272 (0.0003)
exper. ² ($\times 10^2$)	-0.0376 (0.0013)	-0.0411 (0.0008)	-0.0564 (0.0030)	-0.0576 (0.0010)	-0.0338 (0.0012)	-0.0450 (0.0006)
marriage	0.0115 (0.0048)	0.0317 (0.0028)	-0.0005 (0.0079)	0.0262 (0.0026)	0.0118 (0.0041)	0.0322 (0.0019)
nonwhite	-0.0581 (0.0065)	-0.0441 (0.0033)	0.0395 (0.0117)	-0.0236 (0.0033)	-0.0202 (0.0058)	-0.0315 (0.0024)
<i>n</i>	33,600	100,007	12,882	151,545	46,482	251,552

This table reports the estimates of γ in (6.2). “Postsecondary Edu.” stands for the subsample with years of education higher than 12 and “High School or Less” stands for those with years of education less than or equal to 12. Standard errors of the estimates of the coefficients associated with control variables are estimated based on Theorem 3 and reported in parentheses. n is the sample size

We report estimates of π_t , $\beta_{L,t}$ and $\beta_{H,t}$, as well as corresponding mean and standard deviations (denoted by s.d. ($\hat{\beta}_{it}$)) of the return to education (β_{it}) for $t = \{1973-1975\}$ and $\{2001-2003\}$. For a given π_t , the ratio $\beta_{H,t}/\beta_{L,t}$ provides a measure of within-group heterogeneity and allows us to augment information on changes in mean with changes in the distribution of return of education. The estimates for the distribution of the return to education (β_{it}) are summarized in Table 5, with the estimation results for control variables (such as experience, experienced squared, and other individual specific characteristic) reported in Table 6.

As can be seen from Table 5, estimates of s.d. (β_{it}) are strictly positive for all subgroups, except for the “high school or less” group during the first sample period. For this group during the first period the estimate of s.d. (β_{it}) for the male subsample is zero, π is not identified, and we have identical estimates for β_L and β_H . For this subsample, the associated estimates and their standard errors are shown as unavailable (n/a). In case of the female subsample as well as both male and female subsamples where the estimates of s.d. ($\hat{\beta}_{it}$) are close to zero and π is poorly estimated, only the mean of the return to education is informative. In the case of the samples where the estimates of s.d. (β_{it}) are strictly positive, the estimate of the ratio $\beta_{H,t}/\beta_{L,t}$ provides a good measure of within-group heterogeneity of return to education. The estimates of $\beta_{H,t}/\beta_{L,t}$ lie between 1.50 and 2.79, with the high estimate obtained for the females with high school or less education during $\{2001-2003\}$, and the low estimate is obtained for females with postsecondary education during the same period.

As our theory suggests the mean estimates of return to education, $E(\beta_{it})$ are very precisely estimated and inferences involving them tend to be robust to conditional error heteroskedasticity. The results in Table 5 show that estimates of $E(\beta_{it})$ have increased over the two sample periods $t = \{1973-1975\}$ to $t = \{2001-2003\}$, regardless of gender or educational grouping. The postsecondary educational group show larger increases in the estimates of $E(\beta_{it})$ as compared to those with high school or less. Estimates of $E(\beta_{it})$ increase by 36 percent for the postsecondary group, while the estimates of mean return to education rise only by around 5 percent in the case of those with high school or less. This result holds for both genders. Comparing the mean returns across the two educational groups, we find that mean return to education of individuals with postsecondary education is 45 percent higher than those with high school or less in the $\{1973-1975\}$ period, but this gap increases to 87 percent in the second period, $\{2001-2003\}$. Similar patterns are observed in the subsamples by gender. The estimates suggest rising between group heterogeneity, which is mainly due to the increasing returns to education for the postsecondary group.

Turning to within-group heterogeneity, we focus on the estimates of $\beta_{H,t}/\beta_{L,t}$ and first note that over the two periods, within-group heterogeneity has been rising mainly in the case of those with high school or less, for both male and female. For the combined male and female samples and the male subsample, there is little evidence of within-group heterogeneity for the first period $\{1973-1975\}$. However, for the second period $\{2001-2003\}$ we find a sizeable degree of within-group heterogeneity where $\beta_{H,t}/\beta_{L,t}$ is estimated to be around 2.41, with s.d. (β_{it}) ≈ 0.03 . For the female subsample with high school or less, little evidence of heterogeneity was found for the first period, estimates of $\beta_{H,t}/\beta_{L,t}$ increase to 2.79 for the second sample period,

that corresponds to a commensurate rise in s.d. (β_i) to 0.032. The pattern of within-group heterogeneity is very different for those with postsecondary educational. For this group, we in fact observe a slight decline in the estimates of $\beta_{H,t}/\beta_{L,t}$ by gender and over two sample periods.

Overall, our estimates of return to education and the within and between group comparisons are in line with the evidence of rising wage inequality documented in the literature (Corak 2013).

7 Conclusion

In this paper, we consider random coefficient models for repeated cross sections in which the random coefficients follow categorical distributions. Identification is established using moments of the random coefficients in terms of the moments of the underlying observations. We propose two-step generalized method of moments to estimate the parameters of the categorical distributions. The consistency and asymptotic normality of the GMM estimators are established without the IID assumption typically assumed in the literature. Small sample properties of the proposed estimator are investigated by means of Monte Carlo experiments and shown to be robust to heterogeneously generated regressors and errors, although relatively large samples are required to estimate the parameters of the underlying categorical distributions. This is largely due to the highly nonlinear mapping between the parameters of the categorical distribution and the higher-order moments of the coefficients. This problem is likely to become more pronounced with a larger number of categories and coefficients.

In the empirical application, we apply the model to study the evolution of returns to education over two sub-periods, also considered in the literature by Lemieux (2006b). Our estimates show that mean (ex post) returns to education have risen over the periods from 1973–1975 to 2001–2003 mainly in the case of individuals with postsecondary education, and this result is robust by gender. We find evidence of within-group heterogeneity in the case of high school or less educational group as compared to those with postsecondary education.

In our model specification, the number of categories, K , is treated as a tuning parameter and assumed to be known. An information criterion, as in Bonhomme and Manresa (2015) and Su et al. (2016), to determine K could be considered. Further investigation of models with multiple regressors subject to parameter heterogeneity is also required. These and other related issues are topics for future research.

Supplementary Information The online version contains supplementary material available at <https://doi.org/10.1007/s00181-023-02402-0>.

Acknowledgements We would like to thank Timothy Armstrong, Hidehiko Ichimura, Esfandiar Maa-soumi, Geert Ridder, Ron Smith and Hayun Song for helpful comments, and two anonymous referees for constructive comments and suggestions.

Funding Open access funding provided by SCEL, Statewide California Electronic Library Consortium

Declarations

Conflict of interest The authors declare that they have no conflict of interest.

Ethical approval The authors did not receive support from any organization for the submitted work. This article does not contain any studies with human participants or animals performed by any of the authors.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article’s Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article’s Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

A.1 Proofs

We include proofs and technical details in this section.

Proof of Theorem 1 Sum (2.6) over i and rearrange terms,

$$\begin{aligned} &\left(\frac{1}{n} \sum_{i=1}^n E(x_i^r)\right) E(\beta_i^r) + \frac{1}{n} \sum_{i=1}^n E(u_i^r) = \frac{1}{n} \sum_{i=1}^n E(\tilde{y}_i^r) \\ &\quad - \sum_{q=2}^{r-1} \binom{r}{q} \left(\frac{1}{n} \sum_{i=1}^n E(x_i^{r-q}) E(u_i^q)\right) E(\beta_i^{r-q}). \end{aligned} \tag{A.1.1}$$

Note that

$$\frac{1}{n} \sum_{i=1}^n E(x_i^{r-q}) E(u_i^q) = \left(\frac{1}{n} \sum_{i=1}^n E(x_i^{r-q})\right) \sigma_q + \frac{1}{n} \sum_{i=1}^n E(x_i^{r-q}) (E(u_i^q) - \sigma_q),$$

and

$$\left| \frac{1}{n} \sum_{i=1}^n E(x_i^{r-q}) (E(u_i^q) - \sigma_q) \right| \leq \sup_i |E(x_i^{r-q})| \left| \frac{1}{n} \sum_{i=1}^n (E(u_i^q) - \sigma_q) \right| = O(n^{-1/2}),$$

by Assumption 1(b) and 2(b), then by taking $n \rightarrow \infty$ on both sides of (A.1.1), we have (2.8). Similar steps for (2.7) give (2.9). □

Proof of Theorem 2 Let $m_r = E(\beta_i^r)$, $r = 1, 2, \dots, 2K - 1$, which are taken as known. We show that

$$m_r = \sum_{k=1}^K \pi_k b_k^r, \tag{A.1.2}$$

$r = 0, 1, 2, \dots, 2K - 1$, has a unique solution $\theta = (\pi', \mathbf{b}')'$, with $b_1 < b_2 < \dots < b_K$ and $\pi_k \in (0, 1)$ imposed.

Let

$$q(\lambda) = \prod_{k=1}^K (\lambda - b_k) = \lambda^K + (-1)^1 b_1^* \lambda^{K-1} + \dots + (-1)^K b_K^*, \quad (\text{A.1.3})$$

be the polynomial with K distinct roots b_1, b_2, \dots, b_K . Note that for each k , $(b_k^r)_{r=0}^{2K-1}$ satisfies the linear homogeneous recurrence relation,

$$b_k^{K+r} = b_1^* b_k^{K+r-1} + (-1)^1 b_2^* b_k^{K+r-2} + \dots + (-1)^{K-1} b_K^* b_k^r, \quad (\text{A.1.4})$$

for $r = 0, 1, \dots, K - 1$, since q is the characteristic polynomial of the linear recurrence relation (A.1.4) and b_k is a root of q (Rosen 2006, Chapter 5.2). $(m_r)_{r=0}^{2K-1}$ is a linear combination of $(b_1^r)_{r=0}^{2K-1}, (b_2^r)_{r=0}^{2K-1}, \dots, (b_K^r)_{r=0}^{2K-1}$ by (A.1.2), then $(m_r)_{r=0}^{2K-1}$ also satisfies the linear recurrence relation (A.1.4), i.e.,

$$m_{K+r} = b_1^* m_{K+r-1} + (-1)^1 b_2^* m_{K+r-2} + \dots + (-1)^{K-1} b_K^* m_r, \quad (\text{A.1.5})$$

for $r = 0, 1, \dots, K - 1$. (A.1.5) is a linear system of K equations in terms of $(b_k^*)_{k=1}^K$. In matrix form,

$$\mathbf{M} \mathbf{b}^* = \mathbf{m}, \quad (\text{A.1.6})$$

where

$$\mathbf{M} = \begin{pmatrix} 1 & m_1 & \dots & m_{K-1} \\ m_1 & m_2 & \dots & m_K \\ \vdots & \vdots & \ddots & \vdots \\ m_{K-1} & m_K & \dots & m_{2K-2} \end{pmatrix},$$

$\mathbf{D} = \text{diag}((-1)^{K-1}, (-1)^{K-2}, \dots, 1)$, $\mathbf{b}^* = (b_K^*, b_{K-1}^*, \dots, b_1^*)'$, and $\mathbf{m} = (m_K, m_{K+1}, \dots, m_{2K-1})'$.

Denote $\psi_k = (1, b_k, b_k^2, \dots, b_k^{K-1})'$ and $\Psi = (\psi_1, \psi_2, \dots, \psi_K)$. Then

$$\mathbf{M}_k = \begin{pmatrix} 1 & b_k & \dots & b_k^{K-1} \\ b_k & b_k^2 & \dots & b_k^K \\ \vdots & \vdots & \ddots & \vdots \\ b_k^{K-1} & b_k^K & \dots & b_k^{2K-2} \end{pmatrix} = \psi_k \psi_k',$$

and $\mathbf{M} = \sum_{k=1}^K \pi_k \mathbf{M}_k = \Psi \text{diag}(\boldsymbol{\pi}) \Psi'$. Note that Ψ' is a Vandermonde matrix then $\det(\Psi) = \prod_{1 \leq k < k' \leq K} (b_{k'} - b_k) > 0$ since $b_1 < b_2 < \dots < b_K$.

$$\begin{aligned} \det(\mathbf{MD}) &= \det(\Psi \text{diag}(\boldsymbol{\pi}) \Psi') \det(\mathbf{D}) \\ &= \left(\prod_{1 \leq k < k' \leq K} (b_{k'} - b_k) \right)^2 \left(\prod_{k=1}^K \pi_k \right) \left((-1)^{\frac{1}{2}K(K-1)} \right) \neq 0, \end{aligned}$$

since $\pi_k \in (0, 1)$ for any k . Then, we can identify $(b_k^*)_{k=1}^K$ by $(m_r)_{r=0}^{2K-1}$ in (A.1.6), and hence the characteristic polynomial is determined, and we can identify $(b_k)_{k=1}^K$ by (A.1.3).

Since both $(b_k)_{k=1}^K$ and $(m_r)_{r=1}^{2K-1}$ are identified, the first K equations of (A.1.2) is

$$\Psi' \boldsymbol{\pi} = (1, m_1, m_2, \dots, m_{K-1})',$$

and $\boldsymbol{\pi}$ is identified by inverting the Vandermonde matrix Ψ' , which completes the proof. □

Proof of Theorem 4 Denote

$$\Phi_0(\boldsymbol{\theta}, \boldsymbol{\sigma}, \boldsymbol{\gamma}) = \mathbf{g}_0(\boldsymbol{\theta}, \boldsymbol{\sigma}, \boldsymbol{\gamma})' \mathbf{A} \mathbf{g}_0(\boldsymbol{\theta}, \boldsymbol{\sigma}, \boldsymbol{\gamma}),$$

where we stack the left-hand side of (3.7) and transform $\mathbf{m}_\beta = h(\boldsymbol{\theta})$ to get $\mathbf{g}_0(\boldsymbol{\theta}, \boldsymbol{\sigma}, \boldsymbol{\gamma})$. We suppress and the argument $\hat{\boldsymbol{\gamma}}$ and denote $\boldsymbol{\eta} = (\boldsymbol{\theta}', \boldsymbol{\sigma}')'$ for notation simplicity and proceed by verifying the conditions of Newey and McFadden (1994, Theorem 2.1). Theorem 2 provides the identification results which together with the positive definiteness of \mathbf{A} verifies that $\Phi_0(\boldsymbol{\eta}, \boldsymbol{\gamma})$ is uniquely minimized to 0 at $\boldsymbol{\eta}_0$. The compactness of the parameter space holds by Assumption 4(a). Note that $\mathbf{g}_0(\boldsymbol{\eta}, \boldsymbol{\gamma})$ is a polynomial in $\boldsymbol{\eta}$, which is continuous in $\boldsymbol{\eta}$. $\mathbf{g}_0(\boldsymbol{\eta}, \boldsymbol{\gamma})$ is bounded on $\Theta \times \mathcal{S}$. We proceed by verify the uniform convergence condition. The additive terms in $\hat{\mathbf{g}}_n(\boldsymbol{\eta}, \hat{\boldsymbol{\gamma}}) - \mathbf{g}_0(\boldsymbol{\eta}, \boldsymbol{\gamma})$ are of the form $H_{n,1} h^{(r,q)}(\boldsymbol{\eta})$ or $H_{n,2}$, where

$$\begin{aligned} |H_{n,1}| &= \left| \frac{1}{n} \sum_{i=1}^n x_i^{r-q+s_r} - \rho_{0,r-q+s_r} \right| \\ &\leq \left| \frac{1}{n} \sum_{i=1}^n x_i^{r-q+s_r} - \frac{1}{n} \sum_{i=1}^n E(x_i^{r-q+s_r}) \right| + \left| \frac{1}{n} \sum_{i=1}^n E(x_i^{r-q+s_r}) - \rho_{0,r-q+s_r} \right| \\ &= O_p(n^{-1/2}), \end{aligned}$$

$h^{(r,q)}(\eta)$ is a polynomial in η , and

$$\begin{aligned} |H_{n,2}| &= \left| \frac{1}{n} \sum_{i=1}^n \hat{y}_i^r x_i^{s_r} - \rho_{r,s_r} \right| \\ &\leq \left| \frac{1}{n} \sum_{i=1}^n \hat{y}_i^r x_i^{s_r} - \frac{1}{n} \sum_{i=1}^n E(\tilde{y}_i^r x_i^{s_r}) \right| + \left| \frac{1}{n} \sum_{i=1}^n E(\tilde{y}_i^r x_i^{s_r}) - \rho_{r,s_r} \right| \\ &= O_p(n^{-1/2}). \end{aligned}$$

$H_{n,1} = O_p(n^{-1/2})$ and $H_{n,2} = O_p(n^{-1/2})$ are due to Assumption 2(a) and 4(c).

By the compactness of $\Theta \times \mathcal{S}$, $\sup_{\eta \in \Theta \times \mathcal{S}} h^{(r,q)}(\eta) < C < \infty$ for some positive constant C . By triangle inequality, we have

$$\sup_{\eta \in \Theta \times \mathcal{S}} \|\hat{\mathbf{g}}_n(\eta, \hat{\boldsymbol{\gamma}}) - \mathbf{g}_0(\eta, \boldsymbol{\gamma})\| \rightarrow_p 0, \tag{A.1.7}$$

as $n \rightarrow \infty$. Following the proof of Newey and McFadden (1994, Theorem 2.1),

$$\begin{aligned} &\left| \hat{\Phi}_n(\eta, \hat{\boldsymbol{\gamma}}) - \Phi_0(\eta, \boldsymbol{\gamma}) \right| \\ &\leq \left| [\hat{\mathbf{g}}_n(\eta, \hat{\boldsymbol{\gamma}}) - \mathbf{g}_0(\eta, \boldsymbol{\gamma})]' \mathbf{A}_n [\hat{\mathbf{g}}_n(\eta, \hat{\boldsymbol{\gamma}}) - \mathbf{g}_0(\eta, \boldsymbol{\gamma})] \right| \\ &\quad + \left| \mathbf{g}_0(\eta, \boldsymbol{\gamma})' (\mathbf{A}_n + \mathbf{A}'_n) [\hat{\mathbf{g}}_n(\eta, \hat{\boldsymbol{\gamma}}) - \mathbf{g}_0(\eta, \boldsymbol{\gamma})] \right| \\ &\quad + \left| \mathbf{g}_0(\eta, \boldsymbol{\gamma})' (\mathbf{A}_n - \mathbf{A}) \mathbf{g}_0(\eta, \boldsymbol{\gamma}) \right| \\ &\leq \|\hat{\mathbf{g}}_n(\eta, \hat{\boldsymbol{\gamma}}) - \mathbf{g}_0(\eta, \boldsymbol{\gamma})\|^2 \|\mathbf{A}_n\| + 2 \|\mathbf{g}_0(\eta, \boldsymbol{\gamma})\| \|\hat{\mathbf{g}}_n(\eta, \hat{\boldsymbol{\gamma}}) - \mathbf{g}_0(\eta, \boldsymbol{\gamma})\| \|\mathbf{A}_n\| \\ &\quad + \|\mathbf{g}_0(\eta, \boldsymbol{\gamma})\|^2 \|\mathbf{A}_n - \mathbf{A}\|. \end{aligned}$$

By (A.1.7) and the boundedness of \mathbf{g}_0 , $\sup_{\eta \in \eta} \left| \hat{\Phi}_n(\eta, \hat{\boldsymbol{\gamma}}) - \Phi_n(\eta, \boldsymbol{\gamma}) \right| \rightarrow_p 0$, which completes the proof. □

Proof of Theorem 5 We denote $\eta = (\boldsymbol{\theta}', \boldsymbol{\sigma}')'$ for notation simplicity. The first-order condition, $\nabla_{\eta} \hat{\mathbf{g}}_n(\hat{\eta}, \hat{\boldsymbol{\gamma}}) \mathbf{A}_n \hat{\mathbf{g}}_n(\hat{\eta}, \hat{\boldsymbol{\gamma}}) = \mathbf{0}$, holds with probability 1. Denote $\hat{\mathbf{G}}(\eta, \boldsymbol{\gamma}) = \nabla_{\eta} \hat{\mathbf{g}}_n(\eta, \boldsymbol{\gamma})$ and expand $\hat{\mathbf{g}}_n(\hat{\eta}, \hat{\boldsymbol{\gamma}})$ in the first-order condition around η_0 , we have

$$\begin{aligned} \sqrt{n}(\hat{\eta} - \eta_0) &= - \left[\hat{\mathbf{G}}(\hat{\eta}, \hat{\boldsymbol{\gamma}})' \mathbf{A}_n \hat{\mathbf{G}}(\bar{\eta}, \bar{\boldsymbol{\gamma}}) \right]^{-1} \hat{\mathbf{G}}(\hat{\eta}, \hat{\boldsymbol{\gamma}})' \mathbf{A}_n (\sqrt{n} \hat{\mathbf{g}}_n(\eta_0, \hat{\boldsymbol{\gamma}})) \\ &= - \left[\hat{\mathbf{G}}(\hat{\eta}, \hat{\boldsymbol{\gamma}})' \mathbf{A}_n \hat{\mathbf{G}}(\bar{\eta}, \bar{\boldsymbol{\gamma}}) \right]^{-1} \hat{\mathbf{G}}(\hat{\eta}, \hat{\boldsymbol{\gamma}})' \mathbf{A}_n [\sqrt{n} \hat{\mathbf{g}}_n(\eta_0, \boldsymbol{\gamma}_0) \\ &\quad + \nabla_{\boldsymbol{\gamma}} \hat{\mathbf{g}}_n(\eta_0, \bar{\boldsymbol{\gamma}}) \sqrt{n}(\hat{\boldsymbol{\gamma}} - \boldsymbol{\gamma}_0)], \end{aligned}$$

where $\bar{\eta}$ and $\bar{\boldsymbol{\gamma}}$ are between $\hat{\eta}$ and η_0 ; and $\hat{\boldsymbol{\gamma}}$ and $\boldsymbol{\gamma}_0$, respectively. Note that by term-by-term convergence, we have $\hat{\mathbf{G}}(\hat{\eta}, \hat{\boldsymbol{\gamma}}), \hat{\mathbf{G}}(\bar{\eta}, \bar{\boldsymbol{\gamma}}) \rightarrow_p \mathbf{G}_0$ and $\nabla_{\boldsymbol{\gamma}} \hat{\mathbf{g}}_n(\eta_0, \bar{\boldsymbol{\gamma}}) \rightarrow_p$

$\nabla_{\boldsymbol{\gamma}} \mathbf{g}_0(\boldsymbol{\eta}_0, \boldsymbol{\gamma}_0) = \mathbf{G}_{0,\boldsymbol{\gamma}}$. By Assumption 4(b), $\mathbf{A}_n \rightarrow_p \mathbf{A}$. By Assumption 5(a) and (b) and Slutsky theorem,

$$\sqrt{n}(\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}_0) \rightarrow_d (\mathbf{G}'_0 \mathbf{A} \mathbf{G}_0)^{-1} \mathbf{G}'_0 \mathbf{A} (\boldsymbol{\zeta} + \mathbf{G}_{0,\boldsymbol{\gamma}} \boldsymbol{\zeta}_{\boldsymbol{\gamma}}),$$

which completes the proof. □

Further details for Example 4 We need to verify the invertibility of the matrix

$$\mathbf{B} = \begin{pmatrix} 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 1 & 0 & 1 & 0 \\ b_{1L}b_{2L} & b_{1L}b_{2H} & b_{1H}b_{2L} & b_{1H}b_{2H} \end{pmatrix}.$$

The span of the first three rows of \mathbf{B} is

$$\mathcal{S} = \{(\alpha_1 + \alpha_3, \alpha_1, \alpha_2 + \alpha_3, \alpha_3)' : \alpha_1, \alpha_2, \alpha_3 \in \mathbb{R}\}.$$

$(b_{1L}b_{2L}, b_{1L}b_{2H}, b_{1H}b_{2L}, b_{1H}b_{2H})' \notin \mathcal{S}$ is equivalent to $b_{1H}b_{2H} - b_{1H}b_{2L} \neq b_{1L}b_{2H} - b_{1L}b_{2L}$. This can be verified by

$$(b_{1H}b_{2H} - b_{1H}b_{2L}) - (b_{1L}b_{2H} - b_{1L}b_{2L}) = (b_{1H} - b_{1L})(b_{2H} - b_{2L}) > 0,$$

given that $b_{1L} < b_{1H}$ and $b_{2L} < b_{2H}$ hold. □

References

- Ahn SC, Lee YH, Schmidt P (2001) GMM estimation of linear panel data models with time-varying individual effects. *J Econom* 101(2):219–255
- Ahn SC, Lee YH, Schmidt P (2013) Panel data models with multiple time-varying individual effects. *J Econom* 174(1):1–14
- Andrews DWK (2001) Testing when a parameter is on the boundary of the maintained hypothesis. *Econometrica* 69(3):683–734
- Arellano M, Bonhomme S (2012) Identifying distributional characteristics in random coefficients panel data models. *Rev Econ Stud* 79(3):987–1020
- Becker GS (1962) Investment in human capital: a theoretical analysis. *J Polit Econ* 70(5, Part 2):9–49
- Becker GS (1964) Human capital: a theoretical and empirical analysis, with special reference to education. The University of Chicago Press, Chicago
- Beran R (1993) Semiparametric random coefficient regression models. *Ann Inst Stat Math* 45(4):639–654
- Beran R, Hall P (1992) Estimating coefficient distributions in random coefficient regressions. *Ann Stat* 20(4):1970–1984
- Beran R, Millar PW (1994) Minimum distance estimation in random coefficient regression models. *Ann Stat* 22(4):1976–1992
- Beran R, Feuerverger A, Hall P (1996) On nonparametric estimation of intercept and slope distributions in random coefficient regression. *Ann Stat* 24(6):2569–2592
- Bick A, Blandin A, Rogerson R (2022) Hours and wages. *Q J Econ* 137:1901–1962
- Bonhomme S, Manresa E (2015) Grouped patterns of heterogeneity in panel data. *Econometrica* 83(3):1147–1184
- Breunig C, Hoderlein S (2018) Specification testing in random coefficient models. *Quant Econ* 9(3):1371–1417

- Corak M (2013) Income inequality, equality of opportunity, and intergenerational mobility. *J Econ Perspect* 27(3):79–102
- Foster A, Hahn J (2000) A consistent semiparametric estimation of the consumer surplus distribution. *Econ Lett* 69(3):245–251
- Gautier E, Hoderlein S (2015). A triangular treatment effect model with random coefficients in the selection equation. Working Paper. [arXiv:1109.0362](https://arxiv.org/abs/1109.0362)
- Gautier E, Kitamura Y (2013) Nonparametric estimation in random coefficients binary choice models. *Econometrica* 81(2):581–607
- Griliches Z (1977) Estimating the returns to schooling: some econometric problems. *Econometrica* 45(1):1–22
- Hausman JA (1981) Exact consumer's surplus and deadweight loss. *Am Econ Rev* 71(4):662–676
- Hausman JA, Newey WK (1995) Nonparametric estimation of exact consumers surplus and deadweight loss. *Econometrica* 63(6):1445–1476
- Heckman JJ (2001) Micro data, heterogeneity, and the evaluation of public policy: nobel lecture. *J Polit Econ* 109(4):673–748
- Heckman JJ, Humphries JE, Veramendi G (2018) Returns to education: the causal effects of education on earnings, health, and smoking. *J Polit Econ* 126(S1):S197–S246
- Hoderlein S, Klemelä J, Mammen E (2010) Analyzing the random coefficient model nonparametrically. *Econom Theor* 26(3):804–837
- Hoderlein S, Holzmann H, Meister A (2017) The triangular model with random coefficients. *J Econom* 201(1):144–169
- Hsiao C, Pesaran MH (2008) Random coefficient models. In: Mátyás L, Sevestre P (eds) *The econometrics of panel data*, chapter 6. Springer, Berlin, pp 185–213
- Ichimura H, Thompson TS (1998) Maximum likelihood estimation of a binary choice model with random coefficients of unknown distribution. *J Econ* 86(2):269–295
- Lemieux T (2006a) The “mincer equation” thirty years after schooling, experience, and earnings. In: Grossbard S (ed) *Jacob Mincer a pioneer of modern labor economics*, chapter 11. Springer, New York, pp 127–145
- Lemieux T (2006b). Post-secondary education and increasing wage inequality. Working Paper No. 12077, National Bureau of Economic Research
- Lemieux T (2006c) Postsecondary education and increasing wage inequality. *Am Econ Rev* 96(2):195–199
- Masten MA (2018) Random coefficients on endogenous variables in simultaneous equations models. *Rev Econ Stud* 85(2):1193–1250
- Matzkin RL (2012) Identification in nonparametric limited dependent variable models with simultaneity and unobserved heterogeneity. *J Econom* 166(1):106–115
- Mincer J (1974) *Schooling, Experience and Earnings*, National Bureau of Economic Research, New York. 0-87014-265-8
- Newey K, McFadden D (1994) Large sample estimation and hypothesis. In: Engle RF, McFadden DL (eds) *Handbook of econometrics*, volume 4, chapter 36. Elsevier, Amsterdam, pp 2112–2245
- Nicholls D, Pagan A (1985) Varying coefficient regression. In: Hannan EJ, Krishnaiah PR, Rao MM (eds) *Handbook of statistics*, volume 5, chapter 16. Elsevier, Amsterdam, pp 413–449
- Pesaran MH, Zhou Q (2018) To pool or not to pool: revisited. *Oxford Bull Econ Stat* 80(2):185–217
- Rosen K (2006) *Discrete mathematics and its applications*, 6th edn. McGraw-Hill Education, New York
- Su L, Shi Z, Phillips PC (2016) Identifying latent structures in panel data. *Econometrica* 84(6):2215–2264

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.