

# Identification and estimation of dynamic random coefficient models

Wooyong Lee\*

September 9, 2022

([click here to view the latest version](#))

## Abstract

I study panel data linear models with predetermined regressors (e.g. lagged dependent variable) that allow the coefficients as well as the intercept to be individual-specific, allowing for unobserved heterogeneity in the effects of regressors on the dependent variable. I show that the model is not point-identified in a short panel context but rather partially identified, and I characterize sharp identified sets of the mean, variance and CDF of the coefficient distributions. The characterization is general, allowing discrete, continuous and unbounded data. A computationally efficient estimation and inference procedure is proposed, based on a fast and exact global polynomial optimization algorithm. The method is applied to study lifecycle earnings dynamics in U.S. households in the Panel Study of Income Dynamics (PSID) dataset. The results suggest substantial unobserved heterogeneity in earnings persistence, which implies that the households face different levels of earnings risk that lead to heterogeneity in their consumption and savings behaviors.

---

\*Economics Discipline Group, University of Technology Sydney. Address: 14-28 Ultimo Road, Ultimo, NSW 2007, Australia. Tel: (02) 9514 3074. Email: [wooyong.lee@uts.edu.au](mailto:wooyong.lee@uts.edu.au)

# 1 Introduction

Panel data linear models with predetermined regressors (e.g. lagged dependent variable) are extremely popular in empirical research (Arellano and Bond, 1991; Blundell and Bond, 1998). Most of these models allow for fixed effects, which are individual-specific intercepts that allow for unobserved heterogeneity in the levels of dependent variables. Fixed effects offer a flexible method of controlling for unobserved heterogeneity in the levels, helping researchers to explore various research questions such as the effectiveness of a policy. Fixed effects models are well-understood for short panel data (i.e., panel data with a small number of waves).

In addition to the heterogeneity in the levels, there is ample evidence that individuals have unobserved heterogeneity in the effects of regressors on the dependent variable. For example, firms have different degrees of efficiency of labor for production, individuals have different rates of return to education, and households have different degrees of persistence in their earnings with respect to their past earnings. Such heterogeneous effects are crucial mechanisms for heterogeneous responses to exogenous shocks and policies, such as employment subsidies, tuition subsidies, and income tax reform. Heterogeneous effects also have a first-order influence on outcomes of various economic models. For example, heterogeneity in earnings persistence governs heterogeneity in earnings risk that households experience, which leads to heterogeneous motive for precautionary savings in the lifecycle model of consumption.

This paper studies a panel data linear model with predetermined regressors that allows for unobserved heterogeneity in both the levels and the effects of regressors (i.e., a dynamic random coefficient model) in a short panel context. Consider a stylized example:

$$Y_{it} = \beta_{i0} + \beta_{i1}Y_{i,t-1} + \varepsilon_{it},$$

where all variables are scalars and  $\varepsilon_{it}$  is uncorrelated with the current history of  $Y_{it}$  (up to  $t - 1$ ) but potentially correlated with its future values. In this model, both the intercept ( $\beta_{i0}$ ) and the coefficient ( $\beta_{i1}$ ) are individual-specific, reflecting the heterogeneity in the levels and the effects of regressors, respectively. The model also allows lagged dependent variable  $Y_{i,t-1}$  to be a regressor, reflecting dynamics.

Analysis of this model is challenging in short panels since it is impossible to learn about individual values of the  $\beta_i$ s with a small number of waves. An important work by Chamberlain (1993), recently published as Chamberlain (2022), showed that the mean of  $\beta_i$ s in dynamic random coefficient models are not point-identified, which implies that it is not consistently estimable. Since this negative result in the 1990s, there has not been

much progress in the literature. Arellano and Bonhomme (2012) showed that when the regressors are binary, the mean of  $\beta_i$ s for some subpopulations are identifiable and hence consistently estimable, but they did not provide a general identification result that allows for non-binary regressors. Most research on random coefficient models in short panels focus on non-dynamic contexts (Chamberlain, 1992; Wooldridge, 2005; Arellano and Bonhomme, 2012; Graham and Powell, 2012), but such contexts exclude important dynamic mechanisms such as feedback from the current dependent variable to the future regressors. For example, a firm’s labor purchase decision next year might depend on this year’s output because the firm might learn about its own efficiency of labor from the output. A researcher might also be interested in the dynamic mechanisms themselves. For example, earnings persistence of a household with respect to its past earnings is an important parameter because high earnings persistence makes earnings shocks last, reducing a household’s consumption smoothing ability and hence household welfare.

This paper is the first to present a general identification result for dynamic random coefficient models. Identification results for various features of  $\beta_i$ s are presented, including the mean, variance and CDF of  $\beta_i$ s. This paper proposes a computationally feasible estimation and inference procedure for these features, an essential step of which is to use a fast and exact algorithm for solving global polynomial optimization problems. The estimation and inference method is then applied to learn about unobserved heterogeneity in lifecycle earnings dynamics across U.S. households in the Panel Study of Income Dynamics (PSID) dataset. The results of this paper are presented in three steps.

First, I show that dynamic random coefficient models are partially identified, which means that there exist finite lower and upper bounds for the parameters of interest such as the mean, variance and CDF of  $\beta_i$ s that are estimable with data. The result is general, allowing the data and coefficients to be discrete, continuous or unbounded. I give a simple expression for the bounds of the mean of  $\beta_i$ s, which explicitly shows that the bounds are finite even if the data and coefficients are unbounded, as long as certain moments of data are finite. A key idea for the results is to recast the identification problem into a linear programming problem (Honoré and Tamer, 2006; Mogstad, Santos, and Torgovitsky, 2018; Torgovitsky, 2019), which becomes an infinite-dimensional problem when the data or coefficients are continuous. I then use the dual representation of infinite-dimensional linear programming (Galichon and Henry, 2009; Schennach, 2014) to obtain sharp bounds for the parameters of interest.

Second, I show that the sharp bounds can be computed efficiently by exploiting the linear structure of the model. Computing sharp bounds obtained from dual representation involves solving a nested optimization problem in which a researcher maximizes an

objective function that contains another minimization problem. An important computational issue is that the inner minimization problem is a global minimization problem of a possibly non-convex function. I show that for random coefficient models, the inner objective function is a polynomial. I then use the semidefinite relaxation algorithm (Lasserre, 2010, 2015), a fast and exact algorithm for solving the global polynomial optimization problem, to efficiently compute the sharp bounds. The algorithm also allows computationally tractable inference about the parameters based on testing moment inequalities (Chernozhukov, Lee, and Rosen, 2013; Romano, Shaikh, and Wolf, 2014; Chernozhukov, Chetverikov, and Kato, 2019; Bai, Santos, and Shaikh, 2019). For researchers interested in using the semidefinite relaxation approach to global polynomial optimization, I offer a general-purpose R package `optpoly` that implements the approach<sup>1</sup>.

Third, I estimate a reduced-form lifecycle model of earnings dynamics. Lifecycle earnings processes are key inputs in various economic models, including models of lifecycle consumption dynamics (Hall and Mishkin, 1982; Blundell, Pistaferri, and Preston, 2008; Blundell, Pistaferri, and Saporta-Eksten, 2016; Arellano, Blundell, and Bonhomme, 2017). Specifying an earnings process that highlights features of real data is important for calibrating and drawing conclusions from these models. This paper investigates unobserved heterogeneity in earnings of U.S. households in Panel Study of Income Dynamics (PSID) dataset. Guvenen (2007, 2009) pointed out that, when allowing for unobserved heterogeneity in time trend of earnings (called heterogeneous income profiles, HIP), earnings persistence of the permanent income process is estimated to be significantly smaller than 1, which is the estimate from the model without the heterogeneity in time trend (called restricted income profiles, RIP). I estimate a more general model that also allows for unobserved heterogeneity in the earnings persistence itself, and I find that both HIP and RIP have similar estimates of the average earnings persistence that is significantly smaller than 1. This suggests that misspecifying HIP as RIP or vice versa may not lead to serious misspecification when the earnings persistence is allowed to be heterogeneous. I also find evidence of substantial unobserved heterogeneity in earnings persistence itself, which implies that households face different levels of earnings risk that lead to heterogeneity in their consumption and savings behaviors.

Identification results from this paper can be extended generally to other structural models to allow for heterogeneous effects. For example, the results can be applied to allow for individual-specific coefficients and intercept in probit and logit regressions. They can also be applied to vector-valued regressions, such as panel data Vector Autoregressive (VAR) models and the systems of panel data regressions.

---

<sup>1</sup>Available at <https://github.com/wooyong/optpoly>.

The remainder of this paper is structured as follows. Section 2 introduces a dynamic random coefficient model. Sections 3 and 4 present identification results about the model. Sections 5 and 6 introduce estimation, inference and computation methods about the parameters of interest. Section 7 checks performance of the inference method by simulation. Section 8 applies the methods to lifecycle earnings dynamics. Section 9 concludes. All proofs and tables appear in the Appendix.

## 2 Model and motivating examples

The dynamic random coefficient model is specified as:

$$Y_{it} = Z'_{it}\gamma_i + X'_{it}\beta_i + \varepsilon_{it}, \quad t = 1, \dots, T, \quad (1)$$

where  $i$  is an index of individuals,  $T$  is the length of panel data,  $(Y_{it}, Z_{it}, X_{it}) \in \mathbb{R} \times \mathbb{R}^q \times \mathbb{R}^p$  are observed real vectors at time  $t = 1, \dots, T$ , and  $\varepsilon_{it} \in \mathbb{R}$  is an idiosyncratic error term at time  $t$ . Let  $Y_i = (Y_{i1}, \dots, Y_{iT})$  be the full history of  $\{Y_{it}\}$  and  $Y_i^t = (Y_{i1}, \dots, Y_{it})$  be the history of  $\{Y_{it}\}$  up to time  $t$ . Define  $X_i, X_i^t, Z_i, Z_i^t$  similarly. Assume:

$$\mathbb{E}(\varepsilon_{it} | \gamma_i, \beta_i, Z_i, X_i^t) = 0 \quad (2)$$

so that the error term is mean independent of the full history of  $\{Z_{it}\}$  (strict exogeneity) but of current history of  $\{X_{it}\}$  (sequential exogeneity). The presence of a sequentially exogenous regressor  $\{X_{it}\}$  makes (1) a dynamic model. For example, lagged dependent variable  $Y_{i,t-1}$  can be included in  $X_{it}$ .

The model is studied in a short panel context, which corresponds to the asymptotics that the number of individuals  $N \rightarrow \infty$ , but  $T$  is fixed. The random coefficients  $(\gamma_i, \beta_i)$  are unobserved random variables that follow a nonparametric distribution which can freely correlate among themselves and to  $(Z_i, X_{i1})$ . This is the sense that a random coefficient model extends a fixed effects model.

The following notation is used throughout the paper. Let  $W_i = (Y_i', Z_i', X_i')' \in \mathcal{W}$  be the vector of data and  $V_i = (\gamma_i', \beta_i')' \in \mathcal{V}$  be the vector of coefficients. Then,  $\varepsilon_{it}$  is understood as a deterministic function of  $(W_i, V_i)$  by the relationship  $\varepsilon_{it} = Y_{it} - Z'_{it}\gamma_i - X'_{it}\beta_i$ .

I consider a parameter  $\theta$  that has the form:

$$\theta = \mathbb{E}(m(Y_i, Z_i, X_i, \gamma_i, \beta_i)) = \mathbb{E}(m(W_i, V_i))$$

for some known function  $m$ . Identification results are presented for a generic function  $m$ ,

but I focus on the case in which  $m$  is either a polynomial or an indicator function with respect to  $V_i$  which allows for computationally feasible estimation and inference. This choice of  $m$  includes many important parameters of interest. For example,  $\theta$  can be an element of the mean of random coefficients  $\mathbb{E}(\beta_i)$  or an element of the second moments  $\mathbb{E}(\beta_i\beta_i')$ .  $\theta$  can also be the error variance  $\mathbb{E}(\varepsilon_{it}^2)$  because  $\varepsilon_{it}^2 = (Y_{it} - Z'_{it}\gamma_i - X'_{it}\beta_i)^2$  is a quadratic polynomial in  $(\gamma_i, \beta_i)$ . Another example of  $\theta$  is the CDF of  $\beta_i$  evaluated at  $b$ , in which case  $m$  is set to be  $m = \mathbf{1}(\beta_i \leq b)$  which yields  $\theta = \mathbb{E}(\mathbf{1}(\beta_i \leq b)) = \mathbb{P}(\beta_i \leq b)$ .

**Example 1** (Household earnings). One of the simplest examples of (1) is the AR(1) model with heterogeneous coefficients:

$$Y_{it} = \gamma_i + \beta_i Y_{i,t-1} + \varepsilon_{it}, \quad (3)$$

where all variables are scalars. This is a special case of (1), with  $Z_{it} = 1$  and  $X_{it} = Y_{i,t-1}$ .

The AR(1) model is a popular choice for empirical specification of the lifecycle earnings process, where  $Y_{it}$  is the log-earnings net of demographic variables, which is an important input in the lifecycle model of consumption and savings behavior<sup>2</sup>. Persistence of earnings ( $\beta_i$ ) governs earnings risk experienced by households, which is a fundamental motive of precautionary savings. Specifying an earnings process that highlights features of real data is important for drawing conclusions from the model. The literature often models it as an AR(1) process with no coefficient heterogeneity (Lillard and Weiss, 1979; Blundell, Low, and Preston, 2013; Gu and Koenker, 2017) or more simply as a unit root process, which is an AR(1) process with  $\gamma_i = 0$  and  $\beta_i = 1$  (Hall and Mishkin, 1982; Meghir and Pistaferri, 2004; Kaplan and Violante, 2014).

Guvenen (2007, 2009) estimated a variation of (3) where  $\beta_i = \beta$  is homogeneous and the time trend is heterogeneous. He pointed out that  $\beta$  is estimated to be significantly less than 1 when the time trend is allowed to be heterogeneous, in contrast to earlier findings that  $\beta$  is estimated to be close to 1 (e.g. Abowd and Card, 1989; Topel and Ward, 1992). I find in the application that, when  $\beta_i$  is allowed to be heterogeneous,  $\mathbb{E}(\beta_i)$  is estimated to be significantly less than 1 regardless of whether the time trend is heterogeneous or not. Other studies that allow for coefficient heterogeneity in earnings include Browning, Ejrnaes, and Alvarez (2010) and Alan, Browning, and Ejrnæs (2018), with factor structure on the coefficients.  $\square$

**Example 2** (Household consumption behavior). Consider a model of lifecycle consumption behavior:

$$C_{it} = \gamma_{i0} + \gamma_{i1} Y_{it} + \beta_i A_{it} + v_{it}, \quad (4)$$

---

<sup>2</sup>In the literature, it is standard to add a transitory shock to (3).

where all variables are scalars,  $C_{it}$  is non-durable consumption,  $Y_{it}$  is earnings, and  $A_{it}$  is asset holdings at time  $t$ , all measured in logs and net of demographic variables. In this model,  $Y_{it}$  may be taken as strictly exogenous, meaning that the future earnings stream is unaffected by the current consumption choice. However,  $A_{it}$  must be taken as sequentially exogenous since assets and consumptions interrelate through the intertemporal budget constraint.

(4) can be considered an approximation of the consumption rule derived from a structural model (Blundell, Pistaferri, and Saporta-Eksten, 2016). One parameter of interest is  $\gamma_{i1}$ , the elasticity of consumption to earnings. This quantity measures a household's ability to smooth consumption against exogenous changes in earnings, such as exogenous earnings shocks, which is a determinant of a household's consumption smoothing ability and hence household welfare. Similar to the case of Example 1, the literature focuses on models with no coefficient heterogeneity<sup>3</sup>.

Another parameter of interest is  $\beta_i$ , the elasticity of consumption to asset holdings, which measures a household's ability to smooth consumption against exogenous changes to assets. (4) allows a researcher to estimate this quantity while being agnostic about the evolution of assets over time (i.e., under nonparametric evolution of the assets).  $\square$

Results from this paper also extend to a multivariate version of (1), a system of random coefficient models. For example, a researcher can combine (3) and (4) in Examples 1 and 2 and consider a joint lifecycle model of earnings and consumption behavior. The resulting model allows the coefficients from the two processes to freely correlate among themselves and to  $(Y_{i0}, A_{i1})$ , allowing for correlation between earnings and consumption processes. Full description of the multivariate model can be found in the Online Appendix.

### 3 Identification of the mean

This and the next sections present identification results for the dynamic random coefficient model defined in (1). This section focuses on identification of the mean of random coefficients, which gives intuition for the general result in the next section.

Consider identifying a parameter that has the form:

$$\mu_e = \mathbb{E}(e'_\gamma \gamma_i + e'_\beta \beta_i) = \mathbb{E}(e' V_i)$$

where  $e_\gamma$  and  $e_\beta$  are real-valued vectors that the researcher chooses and  $e \equiv (e'_\gamma, e'_\beta)'$ . For example, if  $e_\gamma = 0$  and  $e_\beta = (1, 0, \dots, 0)'$ , then  $\mu_e$  is the mean of the first entry of  $\beta_i$ .

---

<sup>3</sup>See Jappelli and Pistaferri (2010) for a survey.

I present results regarding identification of  $\mu_e$  in two subsections. The first subsection shows that  $\mu_e$  is generally not point-identified. The following subsection shows that  $\mu_e$  is non-trivially partially identified.

### 3.1 Failure of point-identification

This subsection shows that  $\mu_e$  is generally not point-identified, by considering a specific example of (1) and showing that  $\mu_e$  is not point-identified in that example.

The example considered is the AR(1) model with heterogeneous coefficients in which two waves are observed:

$$Y_{it} = \gamma_i + \beta_i Y_{i,t-1} + \varepsilon_{it}, \quad \mathbb{E}(\varepsilon_{it} | \gamma_i, \beta_i, Y_i^{t-1}) = 0, \quad t = 1, 2. \quad (5)$$

The following proposition states that  $\mathbb{E}(\beta_i)$  is not point-identified in this model, which implies that there is no consistent estimator for  $\mathbb{E}(\beta_i)$ .

**Proposition 1.** *Consider the model defined in (5). Assume that  $(Y_{i0}, Y_{i1}, Y_{i2}, \gamma_i, \beta_i) \in \mathcal{C}$ , where  $\mathcal{C}$  is a compact subset of  $\mathbb{R}^5$ . Assume also that  $(Y_{i0}, Y_{i1}, Y_{i2}, \gamma_i, \beta_i)$  are absolutely continuous with respect to the Lebesgue measure and that their joint density is strictly positive on  $\mathcal{C}$ . Then,  $\mathbb{E}(\beta_i)$  is not point-identified.*

The same result holds when  $(Y_{i0}, Y_{i1}, Y_{i2}, \gamma_i, \beta_i)$  is discrete and the number of support points of  $(\gamma_i, \beta_i)$  is not too small relative to that of  $(Y_{i0}, Y_{i1}, Y_{i2})$ . Chamberlain (1993), recently published as Chamberlain (2022), showed that  $\mathbb{E}(\beta_i)$  is not point-identified in (5) when  $Y_{it}$ s are discrete and  $\varepsilon_{it}$  is mean independent of  $Y_i^{t-1}$ . Proposition 1 generalizes the result, showing that point-identification also fails with stronger assumptions and continuous data. Failure of point-identification in both discrete and continuous cases in (5) suggests that it is a general feature of dynamic random coefficient models.

Proof of Proposition 1 implies that  $\mathbb{E}(\beta_i)$  is point-identified if and only if there exists an unbiased estimator of  $\beta_i$  in the AR(1) process, which is worth stating separately:

**Lemma 1.** *Under assumptions of Proposition 1,  $\mathbb{E}(\beta_i)$  is point-identified if and only if there exists an unbiased estimator of  $\beta_i$  in individual time series, that is, a function  $S^*(Y_{i0}, Y_{i1}, Y_{i2})$  which is a linear functional on the space of finite and countably additive signed Borel measures that are absolutely continuous with respect to the Lebesgue measure such that*

$$\mathbb{E}(S^*(Y_{i0}, Y_{i1}, Y_{i2}) | \beta_i) = \beta_i$$

*almost surely. When such  $S^*$  exists,  $\mathbb{E}(\beta_i)$  is identified by  $\mathbb{E}(\beta_i) = \mathbb{E}(S^*(Y_{i0}, Y_{i1}, Y_{i2}))$ .*



I prove Proposition 1 by showing that there is no unbiased estimator of  $\beta_i$ . The intuition for Lemma 1 is as follows. Since the distribution of  $\beta_i$  is unrestricted, information on individual  $\beta_i$  can only be obtained from its individual time series. In long panel context, such information can be obtained by a time series estimator of  $\beta_i$  that is consistent as  $T \rightarrow \infty$ . In short panel context, however, such estimator is not reliable because  $T$  is finite. Lemma 1 shows that a time series estimator that is unbiased for finite  $T$  is the only reliable information on  $\beta_i$  for point-identification in short panels.

### 3.2 Partial identification

A natural question following the last subsection is whether the data are informative at all about  $\mu_e = \mathbb{E}(e'V_i)$ , or whether there is no information. This subsection shows that the data are informative about  $\mu_e$ . I show that there are finite bounds  $L$  and  $U$  such that

$$L \leq \mu_e \leq U$$

which are estimable with data.

I first concisely write (1) and (2) defining  $R_{it} = (Z'_{it}, X'_{it})'$  to be the vector of regressors at time  $t$ :

$$Y_{it} = R'_{it}V_i + \varepsilon_{it}, \quad t = 1, \dots, T, \quad (6)$$

and

$$\mathbb{E}(\varepsilon_{it}|V_i, Z_i, X_i^t) = 0. \quad (7)$$

In this section and throughout the paper, I use unconditional moment restrictions that are implications of (7). It is known that the set of unconditional moment restrictions of the form

$$\mathbb{E}(g(V_i, Z_i, X_i^t)\varepsilon_{it}) = 0, \quad (8)$$

indexed by a suitable class of functions  $g$ , is equivalent to the conditional moment restriction in (7) (Bierens, 1990; Andrews and Shi, 2013). I choose the class to be the set of polynomial functions and use its finite subset for estimation and inference. Such finite subset of unconditional moment restrictions contains less information than (7), but it yields a computationally feasible estimation and inference procedure. Partial identification results based on (7) are established in the Online Appendix.

Consider the following assumptions:

**Assumption 1.** Random variables  $(W_i, V_i)_{t=1}^T$  and  $(\varepsilon_{it})_{t=1}^T$  satisfy (6).

**Assumption 2.**  $\sum_{t=1}^T R_{it}R'_{it}$  is positive definite with probability 1.

**Assumption 3.** Random variables  $(W_i, V_i)_{t=1}^T$  and  $(\varepsilon_{it})_{t=1}^T$  satisfy, for all  $t = 1, \dots, T$ :

$$\begin{aligned}\mathbb{E}((R'_{it}V_i)\varepsilon_{it}) &= 0, \\ \mathbb{E}((Z'_i, X^{t'}_i)'\varepsilon_{it}) &= 0.\end{aligned}$$

Assumption 1 states that the dynamic random coefficient model is correctly specified. Assumption 2 is a no-multicollinearity assumption imposed on individual time series. This is stronger than the assumption that  $\mathbb{E}(\sum_{t=1}^T R_{it}R'_{it})$  is positive definite, a standard assumption in dynamic fixed effect models. A stronger assumption is required because the distribution of  $V_i$  is unrestricted and each  $V_i$  can only be learned from its individual data. If there are individuals whose data are not informative about  $V_i$ , then the data are not informative about  $\mathbb{E}(e'V_i)$  because the missing  $V_i$  values might be arbitrarily large or small<sup>4</sup>. Assumption 3 considers a specific choice of unconditional moment restrictions that are implications of (7). The first restriction states that the “error term”  $(\varepsilon_{it})$  is orthogonal to the “explained term”  $(R'_{it}V_i)$ . The second restriction states that  $\varepsilon_{it}$  is orthogonal to the full history of  $Z_{it}$  and the current history of  $X_{it}$ .

The following theorem shows that  $\mu_e$  is partially identified under Assumptions 1 to 3 and additional regularity conditions. This theorem is a special case of Theorem 2 presented in the next section.

**Theorem 1.** *Suppose that Assumptions 1 to 3 hold, and assume additional regularity conditions which will be stated as Assumption 5 in the next section. In addition, assume that  $W_i$  are absolutely continuous with respect to the Lebesgue measure. For brevity of notation, define*

$$\mathcal{R}_i = \sum_{t=1}^T R_{it}R'_{it} \quad \text{and} \quad \mathcal{Y}_i = \sum_{t=1}^T R_{it}Y_{it}.$$

Then  $L \leq \mu_e \leq U$  where

$$[L, U] = \left[ \tilde{V} - \frac{1}{2}\sqrt{\mathcal{E}\mathcal{D}}, \quad \tilde{V} + \frac{1}{2}\sqrt{\mathcal{E}\mathcal{D}} \right]$$

and

$$\begin{aligned}\tilde{V} &= \frac{1}{2}\mathbb{E}(\mathcal{R}_i^{-1}\mathcal{Y}_i) + \frac{1}{2}\mathbb{E}(\mathcal{R}_i)^{-1}\mathbb{E}(\mathcal{Y}_i), \\ \mathcal{E} &= e'\mathbb{E}(\mathcal{R}_i^{-1})e - e'\mathbb{E}(\mathcal{R}_i)^{-1}e, \\ \mathcal{D} &= \mathbb{E}(\mathcal{Y}'_i\mathcal{R}_i^{-1}\mathcal{Y}_i) - \mathbb{E}(\mathcal{Y}_i)'\mathbb{E}(\mathcal{R}_i)^{-1}\mathbb{E}(\mathcal{Y}_i),\end{aligned}$$

where  $\mathcal{E} \geq 0$  and  $\mathcal{D} \geq 0$  and they are zero if and only if  $\mathcal{R}_i^{-1}e$  and  $\mathcal{R}_i^{-1}\mathcal{Y}_i$  are degenerate across

---

<sup>4</sup>Graham and Powell (2012) studied violation of Assumption 2 in a non-dynamic context.

individuals, respectively. In addition,  $[L, U]$  are the sharp bounds of  $\mu_e$  if Assumption 3 is replaced with the following implication of Assumption 3:

$$\begin{aligned} \sum_{t=1}^T \mathbb{E}((R'_{it} V_i) \varepsilon_{it}) &= 0, \\ \sum_{t=1}^T \mathbb{E}(R_{it} \varepsilon_{it}) &= 0. \end{aligned} \tag{9}$$

The closed-form expressions in Theorem 1 give intuition for when  $L$  and  $U$  are finite. In particular,  $L$  and  $U$  are finite even if  $(Y_i, R_i, V_i)$  are unbounded, as long as  $\mathbb{E}(\mathcal{R}_i)$ ,  $\mathbb{E}(\mathcal{Y}_i)$ ,  $\mathbb{E}(\mathcal{R}_i^{-1} \mathcal{Y}_i)$  and  $\mathbb{E}(\mathcal{Y}'_i \mathcal{R}_i^{-1} \mathcal{Y}_i)$  are finite. Note that  $\mathcal{R}_i$  is the squared design matrix of individual time series and that  $\mathcal{R}_i^{-1} \mathcal{Y}_i$  is the OLS estimator of  $V_i$  from individual time series.

I now explain the intuition behind Theorem 1, focusing on the upper bound  $U$ . Consider a Lagrangian where the objective function is the parameter of interest  $e' V_i$  and the constraints are the moment functions in (9):

$$Q(\lambda, \mu, W_i, V_i) = e' V_i + \lambda \sum_{t=1}^T (R'_{it} V_i) \varepsilon_{it} + \mu' \sum_{t=1}^T R_{it} \varepsilon_{it},$$

where  $\lambda \in \mathbb{R}$  and  $\mu$  has the same dimension as  $R_{it}$ . Note that  $\mathbb{E}(Q) = \mathbb{E}(e' V_i) = \mu_e$  because the constraints have zero expectations.

If I substitute  $\varepsilon_{it} = Y_{it} - R_{it} V_i$  into  $Q$  and use the notation of  $\mathcal{R}_i$  and  $\mathcal{Y}_i$  in Theorem 1, I obtain the expression:

$$Q(\lambda, \mu, W_i, V_i) = e' V_i + \lambda \mathcal{Y}'_i V_i - \lambda V'_i \mathcal{R}_i V_i + \mu' \mathcal{Y}_i - \mu' \mathcal{R}_i V_i.$$

This is a quadratic polynomial in  $V_i$  whose second order derivative is

$$\frac{d^2 Q}{dV_i dV'_i} = -2\lambda \mathcal{R}_i.$$

If  $\lambda > 0$ , then the second order derivative is a negative definite matrix, in which case  $Q$  attains a global maximum at the solution to the first-order condition  $dQ/dV_i = 0$ . Let  $P = \max_{v \in \mathcal{V}} Q(\lambda, \mu, W_i, v)$  be the resulting maximum, which is only a function of  $(\lambda, \mu, W_i)$  since  $V_i$  is “maximized out.” Then:

$$P(\lambda, \mu, W_i) \geq Q(\lambda, \mu, W_i, V_i),$$

which implies:

$$\mathbb{E}(P(\lambda, \mu, W_i)) \geq \mathbb{E}(Q(\lambda, \mu, W_i, V_i)) = \mu_e.$$

This shows that  $\mathbb{E}(P)$  is an upper bound of  $\mu_e$  for any  $\lambda > 0$  and  $\mu$ . I then obtain the smallest upper bound by minimizing  $\mathbb{E}(P)$  with respect to  $\lambda > 0$  and  $\mu$ :

$$\min_{\lambda > 0, \mu} \mathbb{E}(P(\lambda, \mu, W_i)) \geq \mu_e.$$

This coincides with  $U$  in Theorem 1, which is the sharp upper bound of  $\mu_e$  under (9). The lower bound can be obtained by repeating the same argument with  $\lambda < 0$ .

## 4 Identification of the general parameters

This section presents a general partial identification result for dynamic random coefficient models. I consider a parameter of interest of the form

$$\theta = \mathbb{E}(m(W_i, V_i))$$

for some known function  $m : \mathcal{W} \times \mathcal{V} \mapsto \mathbb{R}$ . I consider a generic set of unconditional moment restrictions:

**Assumption 4.** Random vectors  $(W_i, V_i)$  satisfy:

$$\mathbb{E}(\phi_k(W_i, V_i)) = 0, \quad k = 1, \dots, K,$$

where  $\phi_k : \mathcal{W} \times \mathcal{V} \mapsto \mathbb{R}$  are moment functions and  $K \in \mathbb{N}$  is the number of moments.

$\varepsilon_{it}$  does not appear in Assumption 4 because  $\varepsilon_{it}$  is understood as a deterministic function of  $(W_i, V_i)$  by the relationship  $\varepsilon_{it} = Y_{it} - R'_{it}V_i$ . Assumption 4 can also be considered as generic moment equalities without connection to random coefficient models. A more general case that involves conditional moment restrictions is studied in the Online Appendix.

The following explains that the moment restrictions considered in the previous section is a special case of Assumption 4.

**Example 3.** Consider identification of  $\mathbb{E}(e'V_i)$  discussed in the previous section. Assumption 3 implies the following moment functions. The  $\phi_k$ s for  $k = 1, \dots, T$  are

$$\phi_k(W_i, V_i) = (R'_{ik}V_i)(Y_{ik} - R'_{ik}V_i).$$

The  $\phi_k$ s for  $k > T$  are entries of the vectors

$$(Z'_i, X_i^{t'})'(Y_i - R'_{it}V_i), \quad t = 1, \dots, T$$

which is a  $(qT + pt)$ -dimensional vector for each  $t$ .

I characterize the identified set of  $\theta$  under Assumption 4 and additional regularity conditions. The idea is to recognize that the identified set can be characterized using linear programs. I then show that their dual programs yield a tractable characterization of the identified set.

Let  $P_{W,V} \in \mathcal{M}_{W \times V}$ , where  $\mathcal{M}_{W \times V}$  is the linear space of finite and countably additive signed Borel measures on  $\mathcal{W} \times \mathcal{V}$  equipped with the total variation norm. Let  $P_W \in \mathcal{M}_W$  be the marginal distribution of  $W_i$  that the econometrician observes. The sharp identified set  $I$  of  $\theta$  is *defined* by:

$$I \equiv \left\{ \int m(w, v) dP \mid \begin{aligned} &P \in \mathcal{M}_{W \times V}, \quad P \geq 0, \quad \int dP = 1, \\ &\int \phi_k(w, v) dP = 0, \quad k = 1, \dots, K, \\ &\int P(w, dv) = P_W(w) \text{ for all } w \in \mathcal{W} \end{aligned} \right\}.$$

$I$  is the collection of all  $\int m(W_i, V_i) dP$  values implied from  $P$  such that (i)  $P$  is a probability distribution of  $(W_i, V_i)$ , (ii)  $P$  satisfies moment restrictions, and (iii) the marginal distribution of  $W_i$  implied from  $P$  equals the observed distribution  $P_W$ . Dependence of  $I$  on  $m$ ,  $P_W$ , and the  $\phi_k$ s are suppressed in the notation.

All defining properties of  $I$  are linear in  $P$ , which means that  $I$  is a convex set in  $\mathbb{R}$  (i.e., an interval). Therefore,  $I$  can be characterized by its lower and upper bounds. The sharp lower bound  $L$  of  $I$  is *defined* by:

$$\min_{P \in \mathcal{M}_{W \times V}, P \geq 0} \int m(w, v) dP \quad \text{subject to} \quad \begin{aligned} &\int \phi_k(w, v) dP = 0, \quad k = 1, \dots, K, \\ &\int P(w, dv) = P_W(w) \text{ for all } w \in \mathcal{W}. \end{aligned} \quad (10)$$

The constraint  $\int dP = 1$  is omitted because it is implied by the last line of (10). Note that  $P_W$  is a probability distribution.

(10) is a linear program in  $P$ , with the caveat that  $P$  is an infinite-dimensional object. (10) is not a tractable characterization of  $L$  in the sense that the estimation methods that it imply are computationally infeasible for random coefficient models. For example, (10)

can be solved by discretizing the space of  $(W_i, V_i)$  and solving the discretized problem (Honoré and Tamer, 2006; Günsilius, 2019), which is computationally infeasible for random coefficient models because the dimension of  $(W_i, V_i)$  is large.  $W_i$  contains the full history of regressors and dependent variables and  $V_i$  contains all random coefficients. For the random coefficient model with  $R$  regressors and  $T$  waves,  $P$  is a distribution on a  $(RT + R + T)$ -dimensional space.

My approach is to use the dual representation of (10) obtained by the duality theorem for infinite-dimensional linear programming (Galichon and Henry, 2009; Schennach, 2014). I assume the following regularity conditions:

**Assumption 5.** The following conditions hold.

- (i)  $\mathcal{W} \times \mathcal{V}$  is a compact set in a Euclidean space.
- (ii)  $(m, \phi_1, \dots, \phi_K)$  are bounded Borel measurable functions on  $\mathcal{W} \times \mathcal{V}$ .
- (iii) The following set is closed:

$$\left\{ \left( \int \phi_1 dP, \dots, \int \phi_K dP, \int P(\cdot, dv), \int m dP \right) \mid P \in \mathcal{M}_{\mathcal{W} \times \mathcal{V}}, P \geq 0 \right\} \subseteq \mathbb{R}^K \times \mathcal{M}_{\mathcal{W}} \times \mathbb{R}.$$

A sufficient condition for Assumption 5 (iii) is that the joint distribution of  $(W_i, V_i)$  in the data generating process, or its observationally equivalent one, is strictly positive on  $\mathcal{W} \times \mathcal{V}$  (Anderson, 1983, Theorem 9).

The following theorem characterizes  $I$  using the dual representation of (10) and the corresponding problem for the sharp upper bound.

**Theorem 2.** Suppose Assumptions 4 and 5 hold. Let  $\lambda_k \in \mathbb{R}$  for  $k = 1, \dots, K$ . Then  $I = [L, U]$  where:

$$L = \max_{\lambda_1, \dots, \lambda_K} \mathbb{E} \left[ \min_{v \in \mathcal{V}} \left\{ m(W_i, v) + \sum_{k=1}^K \lambda_k \phi_k(W_i, v) \right\} \right], \quad (11)$$

and

$$U = \min_{\lambda_1, \dots, \lambda_K} \mathbb{E} \left[ \max_{v \in \mathcal{V}} \left\{ m(W_i, v) + \sum_{k=1}^K \lambda_k \phi_k(W_i, v) \right\} \right]. \quad (12)$$

$\theta$  is point-identified if and only if  $L = U$ . Proof of Theorem 2 then implies a necessary and sufficient condition for point-identification of  $\theta$ , which is worth stating separately:

**Lemma 2.** Suppose that assumptions of Theorem 2 hold. Suppose also that  $(W_i, V_i)$  are absolutely continuous with respect to the Lebesgue measure and that their joint density is strictly positive

on  $\mathcal{W} \times \mathcal{V}$ . Then  $\theta$  is point-identified if and only if there exists a function  $S^*$  which is a linear functional on  $\mathcal{M}_{\mathcal{W}}$  and real numbers  $\lambda_1^*, \dots, \lambda_K^* \in \mathbb{R}$  such that:

$$m(W_i, V_i) + \sum_{k=1}^K \lambda_k^* \phi_k(W_i, V_i) = S^*(W_i)$$

almost surely on  $\mathcal{W} \times \mathcal{V}$ . When such  $S^*$  exists,  $\theta$  is identified by  $\theta = \mathbb{E}(S^*(W_i))$ .

Lemma 2 states that  $\theta$  is point-identified if and only if the Lagrangian reduces to a function of data only.  $S^*$  can be thought of as an unbiased estimator because the term  $\sum_{k=1}^K \lambda_k^* \phi_k(W_i, V_i)$  has zero expectation.

Theorem 2 and Lemma 2 do not explicitly involve dynamic random coefficient models. Theorem 2 is a general duality result for models of moment equalities, where the moment functions contain both observables and unobservables (Schennach, 2014; Li, 2018). In general, it is not obvious that Theorem 2 leads to a computationally feasible estimation and inference procedure. I show in the next sections that, for dynamic random coefficient models, I can obtain a computationally tractable estimation and inference procedure by exploiting that it is a linear model.

## 5 Estimation and inference

This section explains estimation and inference procedure for the identified sets discussed in the previous sections, focusing on describing the procedure. The next section discusses computation of the objects involved in the procedure.

### 5.1 Estimation

Theorem 2 characterizes the lower and upper bounds in the population. In practice, a researcher does not observe the population distribution  $P_{\mathcal{W}}$  but instead observes a finite sample  $(W_1, \dots, W_N)$  of size  $N$  which are i.i.d.  $P_{\mathcal{W}}$ . A natural approach for estimating  $L$  and  $U$  is to replace expectations in (11) and (12) with sample means (the plug-in principle), which is equivalent to considering the empirical version of (10) where  $P_{\mathcal{W}}$  is replaced with the empirical distribution  $\hat{P}_{\mathcal{W}}$ . I define  $\hat{L}$  as an estimator for  $L$ :

$$\hat{L} = \max_{\lambda_1, \dots, \lambda_K} \frac{1}{N} \sum_{i=1}^N \min_{v \in \mathcal{V}} \left\{ m(W_i, v) + \sum_{k=1}^K \lambda_k \phi_k(W_i, v) \right\} \equiv \max_{\lambda \in \mathbb{R}^K} \frac{1}{N} \sum_{i=1}^N G_L(\lambda, W_i), \quad (13)$$

and  $\hat{U}$  as an estimator for  $U$ :

$$\hat{U} = \min_{\lambda_1, \dots, \lambda_K} \frac{1}{N} \sum_{i=1}^N \max_{v \in \mathcal{V}} \left\{ m(W_i, v) + \sum_{k=1}^K \lambda_k \phi_k(W_i, v) \right\} \equiv \min_{\lambda \in \mathbb{R}^K} \frac{1}{N} \sum_{i=1}^N G_U(\lambda, W_i), \quad (14)$$

where  $\lambda \in \mathbb{R}^K$ . Let  $[\hat{L}, \hat{U}]$  be the plug-in bound.

The plug-in bound is used as a key object for estimation and inference, but it is not obvious to compute. Computation of the bound requires solving two types of optimization problems — the inner optimization problem over  $\mathcal{V}$  and the outer optimization problem with respect to  $\lambda_1, \dots, \lambda_K$ . Each problem has its own difficulties. The inner problem must be solved globally, but its objective function is not necessarily convex. It must also be solved fast because it needs to be solved for each  $i$  and for each step of the outer problem. The outer problem must be solved globally as well, and it might be an optimization over a large dimensional space. The next section discusses how to deal with these computational issues. In this section, I discuss estimation and inference assuming that we can solve the two optimization problems numerically.

In what follows, I show consistency of lower plug-in bound in (13) to population lower bound in (11). The consistency of upper plug-in bound follows by the same argument.

In (13), the solution function of the inner optimization problem

$$G_L(\lambda, w) = \min_{v \in \mathcal{V}} \left\{ m(w, v) + \sum_{k=1}^K \lambda_k \phi_k(w, v) \right\}$$

is a deterministic function given the model (i.e., given  $m$  and the  $\phi_k$ s) and  $(\lambda, w)$ . Therefore, what is studied here is consistency of the statistical object

$$\hat{L} = \max_{\lambda} \hat{L}(\lambda) = \max_{\lambda} \frac{1}{N} \sum_{i=1}^N G_L(\lambda, W_i) \quad (15)$$

as an estimator for

$$L = \max_{\lambda} L(\lambda) = \max_{\lambda} \mathbb{E} (G_L(\lambda, W_i)). \quad (16)$$

$\hat{L}(\lambda)$  is the objective function of an M-estimation problem in which  $L(\lambda)$  is the population objective and  $\lambda$  is the parameter that is M-estimated. Consistency then follows by replicating the analysis of M-estimation. The regularity conditions of M-estimation are satisfied by the fact that  $G_L$  is concave in  $\lambda$ .

**Proposition 2.**  $G_L(\lambda, W_i)$  is globally concave in  $\lambda$ , which implies global concavity of  $\hat{L}(\lambda)$ .



**Proposition 3.** *Suppose that  $L$  exists and is finite, and that  $\operatorname{argmax}_{\lambda} L(\lambda)$  is attained in  $\mathbb{R}^K$ .  $\hat{L}$  then converges to  $L$  in probability.*

## 5.2 Inference

This subsection discusses construction of a confidence interval for the parameter  $\theta \in \mathbb{R}$  whose sharp identified set is  $[L, U]$  in Theorem 2. Any value  $\theta \in [L, U]$  must satisfy:

$$\begin{aligned}\theta &\geq L = \max_{\lambda \in \mathbb{R}^K} \mathbb{E}(G_L(\lambda, W_i)), \\ \theta &\leq U = \min_{\lambda \in \mathbb{R}^K} \mathbb{E}(G_U(\lambda, W_i)),\end{aligned}$$

which implies

$$\begin{aligned}\theta &\geq \mathbb{E}(G_L(\lambda, W_i)) \quad \text{for all } \lambda \in \mathbb{R}^K, \\ \theta &\leq \mathbb{E}(G_U(\lambda, W_i)) \quad \text{for all } \lambda \in \mathbb{R}^K.\end{aligned}$$

These imply the following moment inequality conditions:

$$\begin{aligned}\mathbb{E}(G_L(\lambda, W_i) - \theta) &\leq 0 \quad \text{for all } \lambda \in \mathbb{R}^K, \\ \mathbb{E}(\theta - G_U(\lambda, W_i)) &\leq 0 \quad \text{for all } \lambda \in \mathbb{R}^K.\end{aligned}\tag{17}$$

(17) is a moment inequalities model with infinite number of moment restrictions (indexed by  $\lambda \in \mathbb{R}^K$ ). For computational tractability, I choose finite number of moment inequalities from (17). Let  $\Lambda_F$  be a finite subset of  $\mathbb{R}^K$ . Consider a moment inequalities model:

$$\begin{aligned}\mathbb{E}(G_L(\lambda, W_i) - \theta) &\leq 0 \quad \text{for all } \lambda \in \Lambda_F, \\ \mathbb{E}(\theta - G_U(\lambda, W_i)) &\leq 0 \quad \text{for all } \lambda \in \Lambda_F.\end{aligned}\tag{18}$$

Since  $\Lambda_F$  is a subset of  $\mathbb{R}^K$ , I can use (18) to make a conservative inference about  $\theta$  in (17).

How conservative (18) is relative to (17) depends on how much information is contained in (18) relative to (17). Formal analysis of comparison between (17) and (18) is beyond the scope of this paper<sup>5</sup>, but two observations provide guidance on how to choose an informative  $\Lambda_F$  in practice. First, the inequalities in (17) bind at two  $\lambda$  values, namely  $\lambda_L^* = \operatorname{argmax}_{\lambda} \mathbb{E}(G_L(\lambda, W_i))$  and  $\lambda_U^* = \operatorname{argmin}_{\lambda} \mathbb{E}(G_U(\lambda, W_i))$ , and the inequalities are loose for  $\lambda$ s that are distant from them (because  $G_L$  is concave and  $G_U$  is convex). This

---

<sup>5</sup>Galichon and Henry (2011) studied reduction of the number of model restrictions without losing information. Their approach applies to the case in which the model outcomes, which are moment values in moment inequalities models, have discrete support.

means that most of information in (17) is contained in the neighborhood of  $\lambda_L^*$  and  $\lambda_U^*$ <sup>6</sup>. Second, concavity of  $G_L$  (and convexity of  $G_U$ ) implies that  $G_L$  and  $G_U$  are continuous, which means that considering a finite set of points in the neighborhood of  $\lambda_L^*$  and  $\lambda_U^*$  does not lead to serious loss of information compared to considering all points in the neighborhood. These observations lead to a practical strategy for choosing an informative  $\Lambda_F$ : estimate  $\lambda_L^*$  and  $\lambda_U^*$  using (13) and (14) and select a finite number of points in the neighborhoods of them. I check performance of this strategy by simulation in Section 7.

Given  $\Lambda_F$ , (18) is a standard moment inequalities model although  $\Lambda_F$  can be a large set. The literature on many moment inequalities (Romano, Shaikh, and Wolf, 2014; Bai, Santos, and Shaikh, 2019; Chernozhukov, Chetverikov, and Kato, 2019) propose procedures for computing a confidence interval  $[L_\alpha, U_\alpha]$  of  $\theta$  that has an asymptotic size of  $\alpha$  for large  $\Lambda_F$ . Among the proposed methods, a procedure based on multiplier bootstrap by Chernozhukov, Chetverikov, and Kato (2019) is particularly appealing because the bootstrap does not require re-computation of  $G_L$  and  $G_U$  which have high computational costs. Their procedure uses the following test statistic, computed for each  $\theta \in \mathbb{R}$ :

$$T_{CCK}(\theta) = \max \left\{ \max_{\lambda \in \Lambda_F} \left\{ \frac{\sqrt{N}(\mu_{G_L}(\lambda) - \theta)}{\sigma_{G_L}(\lambda)} \right\}, \max_{\lambda \in \Lambda_F} \left\{ \frac{\sqrt{N}(\theta - \mu_{G_U}(\lambda))}{\sigma_{G_U}(\lambda)} \right\} \right\},$$

where

$$\mu_{G_L}(\lambda) = \frac{1}{N} \sum_{i=1}^N G_L(\lambda, W_i) \quad \text{and} \quad \sigma_{G_L}^2(\lambda) = \frac{1}{N} \sum_{i=1}^N (G_L(\lambda, W_i) - \mu_{G_L}(\lambda))^2$$

and where  $\mu_{G_U}(\lambda)$  and  $\sigma_{G_U}^2(\lambda)$  are defined similarly with  $G_U$ .

$T_{CCK}$  is then compared to a critical value  $c_{CCK}(\alpha)$ , computed using multiplier bootstrap. Each multiplier bootstrap replication simulates independent standard normal random draws  $e_1, \dots, e_N \in \mathbb{R}$  and computes:

$$c_{CCK} = \max \left\{ \max_{\lambda \in \Lambda_F} \left\{ \frac{1}{\sqrt{N}} \sum_{i=1}^N e_i \frac{G_L(\lambda, W_i) - \mu_{G_L}(\lambda)}{\sigma_{G_L}(\lambda)} \right\}, \max_{\lambda \in \Lambda_F} \left\{ \frac{1}{\sqrt{N}} \sum_{i=1}^N e_i \frac{\mu_{G_U}(\lambda) - G_U(\lambda, W_i)}{\sigma_{G_U}(\lambda)} \right\} \right\}.$$

The critical value  $c_{CCK}(\alpha)$  is then the  $100 \times (1 - \alpha)$  percentile of the bootstrapped  $c_{CCK}$  values. The confidence interval is then the set of  $\theta$  for which  $T_{CCK}(\theta) \leq c_{CCK}(\alpha)$ . Note that  $c_{CCK}(\alpha)$  does not depend on  $\theta$ , because  $c_{CCK}$  does not. This allows computationally efficient search of  $\theta$  because  $c_{CCK}(\alpha)$  can be computed only once and be fixed.

---

<sup>6</sup>This relates to a step in the inference procedure of Chernozhukov, Lee, and Rosen (2013), in which they compute a set of moment restrictions that are likely to bind.

The inference procedure naturally extends to a vector-valued parameter  $\theta \in \mathbb{R}^d$ , by considering (17) for every entry of  $\theta$ . For example, the moment inequalities for  $\theta = (\theta_1, \theta_2) \in \mathbb{R}^2$  are:

$$\begin{aligned} \mathbb{E}(G_{L1}(\lambda, W_i) - \theta_1) &\leq 0 \quad \text{for all } \lambda \in \mathbb{R}^K, \\ \mathbb{E}(\theta_1 - G_{U1}(\lambda, W_i)) &\leq 0 \quad \text{for all } \lambda \in \mathbb{R}^K, \\ \mathbb{E}(G_{L2}(\lambda, W_i) - \theta_2) &\leq 0 \quad \text{for all } \lambda \in \mathbb{R}^K, \\ \mathbb{E}(\theta_2 - G_{U2}(\lambda, W_i)) &\leq 0 \quad \text{for all } \lambda \in \mathbb{R}^K, \end{aligned} \tag{19}$$

where  $G_{Uk}$  and  $G_{Lk}$  are  $G_L$  and  $G_U$  in (17) for  $\theta_k, k = 1, 2$ . Inference can then be performed by the same procedure, giving a confidence region in  $\mathbb{R}^2$ . This extension can be used to compute a confidence interval for the variance of random coefficients which involves the first and the second moments.

### 5.3 Estimation and inference under over-identification

In practice, the plug-in bound may yield an empty set, in which case  $\hat{L}$  diverges to  $+\infty$  and  $\hat{U}$  diverges to  $-\infty$ . This happens when the empirical distribution  $\hat{P}_W$  does not satisfy the moment restrictions, which may happen even if the population distribution  $P_W$  satisfies them. This can be compared to over-identification in the generalized method of moments (GMM) estimation, where the GMM objective may be strictly positive in the sample even if the moments are correctly specified. In this case, the empirical version of (10) (where  $P_W$  is replaced with  $\hat{P}_W$ ) does not have a feasible solution, giving an empty plug-in bound.

There are two approaches for dealing with this issue. First, a researcher may obtain a point estimate that minimizes distance between the model and the data<sup>7</sup>. Second, a researcher may directly obtain a confidence interval without insisting on a point estimate, assuming that the model is correctly specified. For the first approach, consider the following relaxation of the moment restrictions:

$$|\mathbb{E}(\phi_k(W_i, V_i))| \leq \delta, \quad k = 1, \dots, K, \tag{20}$$

where  $\delta \geq 0$ , which reduces to Assumption 4 when  $\delta = 0$ . This can be thought of as an absolute-value GMM criterion. The smallest  $\delta$  that makes (20) hold with the empirical

---

<sup>7</sup>Andrews and Kwon (2019) study this approach for standard GMM estimation of moment equality models without unobservables.

distribution, denoted by  $\delta^*$ , is given by:

$$\max_{\lambda_1, \dots, \lambda_K} \frac{1}{N} \sum_{i=1}^N \min_{v \in \mathcal{V}} \left\{ \sum_{k=1}^K \lambda_k \phi_k(W_i, v) \right\} \quad \text{subject to} \quad \sum_{k=1}^K |\lambda_k| \leq 1, \quad (21)$$

which can be computed using the computation methods in the next section. Then, for  $\delta \geq \delta^*$ , the plug-in lower bound  $\tilde{L}$  under the relaxation (20) is given by the bound with a negative  $L^1$  penalty:

$$\tilde{L} = \max_{\lambda_1, \dots, \lambda_K} \left[ \frac{1}{N} \sum_{i=1}^N \min_{v \in \mathcal{V}} \left\{ m(W_i, v) + \sum_{k=1}^K \lambda_k \phi_k(W_i, v) \right\} - \delta \sum_{k=1}^K |\lambda_k| \right]. \quad (22)$$

Proofs for (21) and (22) can be found in the Online Appendix, which replicate the proof of Theorem 2. The relaxed upper bound  $\tilde{U}$  is given similarly to (22), with a positive  $L^1$  penalty. When  $\delta > \delta^*$ ,  $\tilde{L}$  (and  $\tilde{U}$ ) computes the smallest (and largest) value of  $\theta$  among the ones that attain the *near-minimum* of the absolute-value GMM criterion in (20).

Although (22) resolves the empty set problem, it has two drawbacks. First, it is an ad-hoc approach, having no formal justification for why relaxation of moment conditions is a good idea. Second, the procedure may yield a point identified set (or a small interval) even if the model is partially identified. The literature dealt with the second problem by choosing  $\delta$  that is reasonably larger than  $\delta^*$  (Mogstad, Santos, and Torgovitsky, 2018), but how much larger it should be remains a question. In the remainder of this subsection, I discuss a more principled approach, which is directly computing a confidence interval without insisting on a point estimate.

Note that the inference procedure that tests (18) does not involve the plug-in bound per se. The plug-in bound is involved only in the step of choosing  $\Lambda_F$ , which I propose to be the set of  $\lambda$ s that are close to the solutions to the plug-in bound problems. The inference procedure is valid regardless of whether the plug-in bound is empty or not; the issue is that there is no guidance for choosing  $\Lambda_F$  when the plug-in bound is empty. In what follows, I propose a strategy for choosing  $\Lambda_F$  when the plug-in bound is empty.

I propose to consider a grid of positive real numbers  $\{\delta_1, \dots, \delta_M\}$  such that  $\delta_m \geq \delta^*$  for all  $m \in \{1, \dots, M\}$ . Then, for each  $\delta_m$ , I compute the relaxed plug-in bounds:

$$\begin{aligned} \tilde{\lambda}_L(\delta_m) &= \operatorname{argmax}_{\lambda_1, \dots, \lambda_K} \left[ \frac{1}{N} \sum_{i=1}^N \min_{v \in \mathcal{V}} \left\{ m(W_i, v) + \sum_{k=1}^K \lambda_k \phi_k(W_i, v) \right\} - \delta_m \sum_{k=1}^K |\lambda_k| \right]. \\ \tilde{\lambda}_U(\delta_m) &= \operatorname{argmin}_{\lambda_1, \dots, \lambda_K} \left[ \frac{1}{N} \sum_{i=1}^N \max_{v \in \mathcal{V}} \left\{ m(W_i, v) + \sum_{k=1}^K \lambda_k \phi_k(W_i, v) \right\} + \delta_m \sum_{k=1}^K |\lambda_k| \right]. \end{aligned}$$

I then propose to choose points in the neighborhoods of *every*  $\tilde{\lambda}_L(\delta_m)$  and  $\tilde{\lambda}_U(\delta_m)$ . I check performance of this approach by simulation in Section 7.

When  $\delta^* = 0$ , i.e. when the plug-in bound is non-empty, a researcher may choose  $M = 1$  with  $\delta_1 = 0$ , in which case the procedure reduces to the previous procedure in Section 5.2. This means that the inference procedure with relaxed bounds generalizes the procedure discussed in Section 5.2.

## 6 Computation

This section discusses computation of objects involved in estimation and inference. In particular, the discussion focuses on computation of the two optimization problems in the plug-in lower bound in (13), which apply similarly to the other objects such as plug-in upper bound in (14), moment inequalities in (18) and relaxed plug-in bounds in (22).

### 6.1 The inner problem

The inner optimization problem of (13) is to evaluate the function

$$G_L(\lambda, w) = \min_{v \in \mathcal{V}} \left\{ m(w, v) + \sum_{k=1}^K \lambda_k \phi_k(w, v) \right\} \quad (23)$$

for each fixed  $w = W_i$ , where  $i = 1, \dots, N$ , given the value of  $\lambda \in \mathbb{R}^K$ .

Evaluation of  $G_L$  involves a global minimization of a possibly non-convex function, and  $G_L$  must be evaluated for each  $w = W_i$  and for each step of the outer optimization problem. In the simple case that  $\mathcal{V}$  is discrete or low-dimensional, the inner problem can be solved by enumerating all points in  $\mathcal{V}$  or the grid points of  $\mathcal{V}$ . However, for random coefficient models, these are not likely to be the case.

This subsection shows that  $G_L$  can be computed fast and exactly when  $m$  and  $\phi_k$ s are polynomials in  $v$ , in which case evaluation of  $G_L$  is equivalent to globally minimizing a polynomial, for which a fast and exact algorithm exists. The polynomial case is useful for computing bounds of many interesting parameters, such as the moments of random coefficients. The following examples describe some of them.

**Example 4.** In Section 3, I showed identification of the mean parameter  $\mathbb{E}(e'V_i)$  under Assumptions 1 to 3. In this case, the  $m$  function is given by  $m(W_i, V_i) = e'V_i$ , which is a linear function of  $V_i$  and hence a first-order polynomial. The moment functions under

Assumption 3 consist of the functions

$$(R'_{it}V_i)(Y_{it} - R'_{it}V_i), \quad t = 1, \dots, T, \quad (24)$$

and the entries of the vectors

$$(Z'_{it}, X'^t_{it})'(Y_{it} - R'_{it}V_i), \quad t = 1, \dots, T, \quad (25)$$

which are at most second-order polynomials of  $V_i$ .

**Example 5.** Consider identification of an element of  $\mathbb{E}(V_i V'_i)$ . Then  $m$  is an element of  $V_i V'_i$ , which is a second-order polynomial of  $V_i$ . Consider the moment restriction  $\mathbb{E}((R'_i V_i)^3 \varepsilon_{it}) = 0$ , in which case the  $\phi_k$ s consist of the functions

$$(R'_{it}V_i)^3(Y_{it} - R'_{it}V_i), \quad t = 1, \dots, T, \quad (26)$$

which are fourth-order polynomials of  $V_i$ . A researcher may also consider the moment functions in Assumption 3, in which case the additional  $\phi_k$ s are set to be (24) and (25).

In Examples 4 and 5, the moment functions are chosen so that they yield finite lower and upper bounds for the parameters of interest. As a practical strategy for obtaining finite bounds, I choose  $\phi_k$ s so that the inner objective function is an even order polynomial whose order is strictly larger than the order of the parameter of interest. In Examples 4 and 5, I choose (24) to obtain a second order polynomial and (26) to obtain a fourth order polynomial as inner objectives. The inner objective polynomial then has its leading coefficient positive or negative depending on the signs of  $\lambda$ , which leads to finite inner minimum or maximum that yields finite lower and upper bounds.

The polynomial case can be extended to allow  $m$  to be an indicator function of  $V_i$ . An indicator function partitions  $\mathcal{V}$  into two exclusive sets, and the indicator function is constant within each set. A researcher can then compute the global optimum in each partition, and then the optimum of the two. This extension is useful for computing bounds for CDFs of random coefficients, which is described as the following example.

**Example 6.** Let  $V_{i1}$  be the first entry of  $V_i \in \mathbb{R}^{q+p}$ , and let  $v^0 \in \mathbb{R}$ . Consider identification of the CDF of  $V_{i1}$  evaluated at  $v^0$ . I set

$$m(W_i, V_i) = \mathbf{1}(V_{i1} \leq v^0),$$

which is an indicator function of  $V_i$ . Consider the same set of moment restrictions as Example 4, in which case the  $\phi_k$ s are at most second-order polynomials in  $V_i$ . The  $m$

function then partitions  $\mathcal{V}$  into two exclusive sets  $\mathcal{V}_1 = \{(v_1, \dots, v_{q+p}) \mid v_1 \leq v\}$  and  $\mathcal{V}_2 = \{(v_1, \dots, v_{q+p}) \mid v_1 > v\}$ , and  $m = 1$  on  $\mathcal{V}_1$  and  $m = 0$  on  $\mathcal{V}_2$ . The inner objective function is then a second-order polynomial within each of  $\mathcal{V}_1$  and  $\mathcal{V}_2$ , for which I can compute the minimum. I can then evaluate the inner objective by taking the smaller optimum between those in  $\mathcal{V}_1$  and  $\mathcal{V}_2$ .

The next two subsections discuss a fast and exact computation method for global optimization of polynomials. The first considers a simple case of quadratic polynomials for which the global solution can be obtained in a closed-form. The second considers generic polynomials for which the global optimization problem is solved numerically.

### 6.1.1 Global optimization of quadratic polynomials

I first consider a simple case of quadratic polynomials. I express a quadratic polynomial in standard form:

$$Q(v) = v'Av + b'v + c$$

where  $A$  is a  $\dim(v) \times \dim(v)$  symmetric matrix,  $b$  is a  $\dim(v)$ -dimensional vector, and  $c \in \mathbb{R}$ . If the inner objective of (23) is expressed in this standard form,  $(A, b, c)$  are functions of the data  $w$ .

Quadratic polynomials can be solved efficiently using quadratic optimization softwares. In practice, a researcher can use a heuristic but faster (that is, closed-form) method to improve the speed. If  $A$  is positive definite,  $Q$  attains the global minimum at the solution to the first order condition, which is given by:

$$\min_{v \in \mathcal{V}} Q(v) = c - \frac{1}{4}b'A^{-1}b. \quad (27)$$

If  $A$  is not positive definite, the minimum of  $Q$  is negative infinity unless  $A$  has a zero eigenvalue. If  $A$  has a zero eigenvalue,  $Q$  has a finite minimum if the first-order condition,  $2Av + b = 0$ , has an infinite number of solutions for which the value of  $Q$  is the same. Otherwise, the minimum of  $Q$  is negative infinity.

In the context of (23), if the data  $w$  is from a continuous distribution,  $A$  has a zero eigenvalue with probability zero. Therefore, for continuous data, a researcher may rule out the possibility of zero eigenvalue in practice, simply using (27) to express (23) in a closed-form if and only if  $A$  is positive definite; otherwise the solution is negative infinity.

The heuristic method discussed in the above applies when  $\mathcal{V}$  in (23) is unbounded. In some cases, a researcher may consider restricting  $\mathcal{V}$  to be a bounded set, such as restricting the autoregressive parameter of the AR(1) model to be in  $[0, 1]$ . In that case, the heuris-

tic method can be modified to incorporate the constraint using the Lagrange multiplier method. Alternatively, quadratic optimization softwares can be used.

### 6.1.2 Global optimization of generic polynomials

When  $m$  and  $\phi_k$ s are polynomials of generic order, a closed-form solution is not available, but it can be solved numerically. The idea is to transform the problem into a convex optimization problem (Lasserre, 2010, 2015). The resulting algorithm is fast and it computes an *exact* solution. This subsection summarizes the main idea of the algorithm. A formal discussion can be found in Lasserre (2010, 2015).

Consider computing the global minimum of a fourth-order polynomial in two variables  $(v_1, v_2)$ . Let  $u(v) = (1, v_1, v_2, v_1^2, v_1v_2, v_2^2)'$  be the vector of monomials up to the second order and  $u_j(v)$  be the  $j$ -th entry of  $u(v)$ . Let  $\{p_j(v)\}$  be the collection of all monomials up to the fourth order, which are unique entries of  $u(v)u(v)'$ . Let  $J$  be the cardinality of  $\{p_j(v)\}$ . I express a fourth-order polynomial in standard form:

$$\pi(v) = \sum_{j=1}^J a_j p_j(v),$$

where  $a_j$  is the coefficient on the monomial  $p_j(v)$ .

Consider minimization of  $\pi(v)$  with respect to  $v \in \mathcal{V}$ . The minimum of  $\pi(v)$  over  $\mathcal{V}$  equals the solution of the minimization problem:

$$\min_{P_V \in \mathcal{M}_V, \int dP_V = 1} \int \pi(v) dP_V \quad (28)$$

where  $P_V$  is a probability distribution on  $\mathcal{V}$ . (28) is minimized at the point-mass distribution concentrated at the minimizer of  $\pi(v)$ . Since  $\pi(v)$  is a linear combination of the  $p_j$ s, I can rewrite (28) as:

$$\min_{P_V \in \mathcal{M}_V, \int dP_V = 1} \sum_{j=1}^J a_j \int p_j(v) dP_V,$$

which can be rewritten further as:

$$\min_{M_1, \dots, M_J \in \mathbb{R}, M_1=1} \sum_{j=1}^J a_j M_j \quad \text{subject to} \quad M_j = \int p_j(v) dP_V \text{ for some } P_V \in \mathcal{M}_V. \quad (29)$$

Except for the fact that the constraint is complicated, (29) is a minimization over  $\mathbb{R}^J$  and the objective is linear (thus convex) in the choice variables.



The idea then is to replace the constraint in (29) with a convex constraint that only involves  $(M_1, \dots, M_J)$ . The constraint in (29) tells that  $(M_1, \dots, M_J)$  must be moments of some underlying distribution, and its necessary condition can be characterized using a matrix<sup>8</sup>. For example, a random variable  $X$  must satisfy  $\mathbb{E}(X^2) - \mathbb{E}(X)^2 \geq 0$ , because  $\text{Var}(X)$  must be nonnegative. This is equivalent to the condition:

$$\begin{pmatrix} 1 & \mathbb{E}(X) \\ \mathbb{E}(X) & \mathbb{E}(X^2) \end{pmatrix} \text{ is positive semidefinite.}$$

This example can be generalized. Define linear operator  $\mathcal{L}$  that maps a polynomial to  $\mathbb{R}$  by the relationship:

$$\mathcal{L}\left(\sum_j a_j p_j(v)\right) = \sum_j a_j M_j.$$

If  $(M_1, \dots, M_J)$  are moments, then it must be that:

$$\mathcal{L}(u(v)u(v)') \text{ is positive semidefinite} \quad (30)$$

where the operator  $\mathcal{L}$  is applied to each element of  $u(v)u(v)'$ .  $\mathcal{L}(u(v)u(v)')$  is a matrix that only involves  $(M_1, \dots, M_J)$ .

(30) is a convex constraint because the set of positive semidefinite matrices is a convex set. Replacing the constraint in (29) with (30) yields a convex optimization problem:

$$\min_{M_1, \dots, M_J \in \mathbb{R}} \sum_{j=1}^J a_j M_j \quad \text{subject to} \quad \mathcal{L}(u(v)u(v)') \text{ is positive semidefinite.} \quad (31)$$

The constraint can be handled more efficiently than a generic convex constraint so that the optimization problem has its own name—semidefinite program (SDP)—an optimization problem in which a matrix that involves the choice variables is constrained to be positive semidefinite.

The SDP approach to polynomial optimization solves (31), the *semidefinite relaxation*, which can be solved fast and reliably using SDP solvers. The algorithm offers *certificate of optimality*, a condition for the optimal value of  $(M_1, \dots, M_J)$ , satisfying which means that the solution to (31) equals the global optimum. For researchers interested in using the semidefinite relaxation approach to global polynomial optimization, I offer a general-purpose R package `optpoly` that implements the approach<sup>9</sup>.

<sup>8</sup>See Lasserre (2010, 2015) for the necessary and sufficient condition that involves a sequence of matrices.

<sup>9</sup>Available at <https://github.com/wooyong/optpoly>.

The solution to (31) (i.e., the SDP solution) is less than or equal to the solution to (29) since a necessary condition is weaker than the original condition. The semidefinite relaxation approach solves a sequence of the SDP programs, or a *hierarchy* of the SDP programs, until the certificate of optimality is obtained, which is known to be obtained in a finite number of steps under suitable conditions. Even if a researcher does not solve the hierarchy of the SDPs, he/she can take an SDP solution as a lower bound for (29), in which case the resulting plug-in bound is a non-sharp but a valid bound for  $\theta$ .

## 6.2 The outer problem

I turn to the outer optimization problem of (13). A researcher needs to solve the optimization problem:

$$\max_{\lambda \in \mathbb{R}^K} \frac{1}{N} \sum_{i=1}^N G_L(\lambda, W_i).$$

Assume that the researcher can evaluate  $G_L$  exactly using the algorithms in the previous subsection. The remaining difficulty then is how to solve the optimization problem where  $K$  can be potentially large.

Recall that  $G_L$  is globally concave, as shown in Proposition 2. This implies that there is only one local maximum in the outer optimization problem, which is also the global maximum. Under suitable conditions,  $G_L$  is differentiable when  $K = 1$  (Milgrom and Segal, 2002, Theorem 3), which can be extended to show that  $G_L$  is directionally differentiable for  $K > 1$ . This implies that the researcher can solve the outer problem using fast convex optimization algorithms such as gradient descent methods.

Concavity of the outer problem comes from the concavity of  $G_L$ , and solving the inner problem exactly by the polynomial optimization algorithm is crucial for computational tractability of the outer problem. This is an important distinction from Schennach (2014) and Li (2018) who studied generic moment equality models. I focus on random coefficient models and exploit the linear structure of the model to achieve computational tractability for the models with large dimensions. If a researcher uses general-purpose global optimization methods that involve random approximation errors (e.g. simulated annealing) to solve the inner problem, then  $G_L$  is no longer concave and the researcher cannot use fast convex optimization algorithms for the outer problem. This is problematic when  $K$  is large, which is often the case in applications of random coefficient models.

## 7 Simulation

This section examines the performance of the inference procedure discussed in Section 5. The simulation uses the AR(1) model given in (3) as the data generating process (DGP):

$$Y_{it} = \gamma_i + \beta_i Y_{i,t-1} + \varepsilon_{it}, \quad t = 1, \dots, T.$$

In DGP,  $\gamma_i \in \mathbb{R}$  and  $\beta_i \in [0, 1]$  follow Normal and Beta distributions respectively, and their joint distribution is given by Gaussian copula.  $\varepsilon_{it}$  follows an independent Normal distribution with mean zero and variance varying over  $t$ . Conditional on  $(\gamma_i, \beta_i)$  where  $\beta_i \leq 0.9$ ,  $Y_{i0}$  is generated from the stationary distribution implied by  $(\gamma_i, \beta_i)$  and the variance of  $\varepsilon_{i1}$ :

$$Y_{i0} \sim N \left( \frac{\gamma_i}{1 - \beta_i}, \frac{\text{Var}(\varepsilon_{i1})}{1 - \beta_i^2} \right).$$

In contrast, conditional on  $(\gamma_i, \beta_i)$  where  $\beta_i > 0.9$ ,  $Y_{i0}$  is generated from an independent Normal distribution because the stationary distribution implied by  $(\gamma_i, \beta_i)$  produces extreme values when  $\beta_i$  is close to 1. Parameter values of the distribution of  $(\gamma_i, \beta_i)$  and  $\varepsilon_{it}$  are set based on the estimates of the income process in the application.

Simulation data are generated in two steps. The first step simulates a dataset of 100,000 observations from the parametric model described above. The second step then creates Monte Carlo samples from the 100,000 observations in the first step, by sampling observations with replacement. I consider the 100,000 observations as a *finite population* from which Monte Carlo samples are generated. Using a finite population is necessary because the identified set of the parametric model itself is infeasible to compute, while the identified set for the finite population can be computed exactly using (13) and (14).

Table 1 lists population bounds for  $\mathbb{E}(\beta_i)$  computed from the finite population of size 100,000 with  $T = 5, 10$  or  $15$ . For each  $T$ , I compute the bounds using moment conditions:

$$\begin{aligned} \mathbb{E}((\gamma_i + \beta_i Y_{i,t-1})\varepsilon_{it}) &= 0, \quad t = 1, \dots, T, \\ \mathbb{E}(\varepsilon_{it}) &= 0, \quad t = 1, \dots, T, \\ \mathbb{E}(Y_{i,t-1-s}\varepsilon_{it}) &= 0, \quad s = 0, \dots, \min\{L, T\}, \quad t = 1, \dots, T, \end{aligned}$$

where  $L = 3, 5$  or  $7$ . I also restrict  $(\gamma_i, \beta_i) \in \mathcal{V} = [-3, 3] \times [0, 1]$  during computation, which is true for the finite populations used in the simulation. I use the same restrictions in the application with  $T = 15$  and  $L = 5$ .

For each  $(T, L)$ , I create Monte Carlo replications by sampling  $N = 750$  or  $1000$  observations from the finite population with replacement. I then compute confidence interval

for  $\mathbb{E}(\beta_i)$  in each Monte Carlo replication, using the procedure with relaxed bounds in Section 5.3. The grid of  $\{\delta_m\}$  is set to be  $\delta_m \in \{1.25\delta^*, 1.50\delta^*, 1.75\delta^*, \dots, 2.75\delta^*, 3\delta^*\}$ , which is an equispaced grid of 8 values.

For each  $\tilde{\lambda}_L(\delta_m)$  and each  $\tilde{\lambda}_U(\delta_m)$ , I sample  $P = 25, 50$  or  $75$  points from its neighborhood by adding Gaussian noise whose standard deviation is inversely proportional to the gradient of the bounds at the  $\tilde{\lambda}$ s. This means that the size of  $\Lambda_F$  — the number of moment inequalities — is  $8P$ . The critical value is computed with 2000 multiplier bootstrap replications.

Tables 2 and 3 present coverage probabilities for combinations of  $N, T, L, P$ . Each coverage probability is computed with 1000 Monte Carlo replications. Simulation results suggest that the proposed inference procedure produces conservative but reasonable coverage probabilities.

## 8 Application to lifecycle earnings dynamics

### 8.1 Overview

Lifecycle earnings dynamics is a key input in various macroeconomic models. For example, in models of consumption and savings dynamics (Hall and Mishkin, 1982; Blundell, Pistaferri, and Preston, 2008; Blundell, Pistaferri, and Saporta-Eksten, 2016; Arellano, Blundell, and Bonhomme, 2017), households facing higher risk in earnings dynamics accumulate more precautionary savings to smooth consumption. Households save more when they experience a positive earnings shock, which is used to maintain consumption during a negative earnings shock. Specifying earnings process that highlights features of real data is important for calibrating and drawing conclusions from these models.

When used as an input, it is common to specify earnings dynamics using a parsimonious linear model. It consists of permanent and transitory income processes<sup>10</sup>:

$$Y_{it} = z_{it} + \varepsilon_{it}, \quad z_{it} = \rho z_{i,t-1} + \eta_{it},$$

where  $Y_{it}$  is log-earnings net of common trends on observables such as demographics and years of experience,  $\{z_{it}\}$  is a permanent income process and  $\{\varepsilon_{it}\}$  is a transitory income process.  $\eta_{it}$  and  $\varepsilon_{it}$  are i.i.d. mean zero shocks.

Guvenen (2007, 2009) studied two leading views on unobserved heterogeneity in earn-

---

<sup>10</sup>As Guvenen (2007) points out, these are stylized versions of what is used in the literature, but they still capture features important for the discussion.

ings dynamics. Consider two earnings processes:

$$\begin{aligned} Y_{it} &= \alpha_i + z_{it} + \varepsilon_{it}, & z_{it} &= \rho z_{i,t-1} + \eta_{it}, & (\text{RIP}) \\ Y_{it} &= \alpha_i + \beta_i h_{it} + z_{it} + \varepsilon_{it}, & z_{it} &= \rho z_{i,t-1} + \eta_{it}, & (\text{HIP}) \end{aligned} \quad (32)$$

where  $h_{it}$  is potential years of experience and  $(\alpha_i, \beta_i)$  are heterogeneous deviations from common trends. These two models are called the Restricted Income Profiles (RIP) process and the Heterogeneous Income Profiles (HIP) process, respectively. In both models,  $\rho$  a key parameter because it represents earnings risk that households face. The literature reports  $0.5 < \rho < 0.7$  and  $\text{Var}(\beta_i) > 0$  for HIP process (e.g. Lillard and Weiss, 1979; Baker, 1997), which means households experience modest persistence and heterogeneous trends. In contrast, MaCurdy (1982) tested hypothesis that  $\text{Var}(\beta_i) = 0$  and did not reject the hypothesis. The literature reports  $\rho \approx 1$  for RIP process (e.g. Abowd and Card, 1989; Topel and Ward, 1992), meaning households experience extreme persistence and homogeneous trends. Guvenen (2007) studied implications of the two models on consumption data and found that HIP is more consistent with features of consumption data. Guvenen (2009) pointed out that misspecifying HIP process as a RIP process leads to an upward biased estimator of  $\rho$ , obtaining  $\rho \approx 1$ .

While there is vast literature on heterogeneity in  $\beta_i$  and its influence on  $\rho$ , there is relatively little work on investigating heterogeneity in  $\rho$  itself. Recent studies on it include Browning, Ejrnaes, and Alvarez (2010) and Alan, Browning, and Ejrnæs (2018) in which heterogeneity in  $\rho$  is given by a factor structure. In this section, I investigate heterogeneity in  $\rho$  by estimating generalization of (32) where  $\rho = \rho_i$ . I treat the modified model as a random coefficient model, meaning that distribution of  $\rho_i$  and its dependence to  $(\alpha_i, \beta_i, Y_{i0})$  are unrestricted. Distributions of  $\eta_{it}$ s are also not restricted and may depend on  $\rho_i$ , allowing for heteroskedasticity.

In the remainder of this section, I find that, when  $\rho = \rho_i$  is allowed to be heterogeneous, both RIP and HIP have similar estimates of  $\mathbb{E}(\rho_i)$  that is significantly less than 1. Confidence intervals for  $\mathbb{E}(\rho_i)$  in the two processes have substantial overlap, both having upper confidence limits around 0.6 at 90% confidence level. Confidence intervals for the CDF of  $\rho_i$  are also similar between the two processes. These results suggest that choosing RIP over HIP or vice versa may not lead to serious misspecification when  $\rho$  is allowed to be heterogeneous. I also find evidence of heterogeneity in  $\rho_i$ , obtaining a confidence interval for  $\text{Var}(\rho_i)$  in RIP process that has a lower confidence limit of 0.067 at 90% confidence level. This implies a lower confidence limit of 0.258 for standard deviation of  $\rho_i$ .

## 8.2 Data and model

I use data on U.S. households from the Panel Study of Income Dynamics (PSID) dataset. The data is based on the dataset of Guvenen (2009), who estimated RIP and HIP processes using PSID earnings data. I constructed data of  $N = 800$  and  $T = 15$  from his dataset, the detail of which can be found in the Online Appendix.

I follow Guvenen (2009) and use the potential experience as a measure of experience,  $h = \text{age} - \max\{\text{years of schooling}, 12\} - 6$ , which is a strictly exogenous variable since it is a deterministic function of age. In addition, because my method requires that there is no multicollinearity in each individual time series (Assumption 2), I removed 40 individuals (that is, five percent of data) who have smallest variations in their reported incomes, giving a dataset of  $N = 760$  and  $T = 15$ . Estimation results with this removal does not qualitatively differ from the results with full data, which can be found in the Online Appendix.

To apply my method to income processes, I transform two models in (32) to random coefficient models. I first separate  $\varepsilon_{it}$  from  $Y_{it}$  using a simulation based de-noising method inspired by Arellano and Bonhomme (2018), described in the Online Appendix, obtaining simulated observations of the permanent processes  $\tilde{Y}_{it}$  for the RIP and the HIP processes:

$$\begin{aligned} \tilde{Y}_{it} &= \alpha_i + z_{it}, & z_{it} &= \rho z_{i,t-1} + \eta_{it}, & \text{(RIP)} \\ \tilde{Y}_{it} &= \alpha_i + \beta_i h_{it} + z_{it}, & z_{it} &= \rho z_{i,t-1} + \eta_{it}. & \text{(HIP)} \end{aligned} \tag{33}$$

Estimation results with de-noising does not qualitatively differ from the results without de-noising, which can also be found in the Online Appendix. I then use quasi-differencing to (33) and transform them to random coefficient models. Quasi-differencing each of (33) yields, respectively:

$$\begin{aligned} \tilde{Y}_{it} &= \alpha_i(1 - \rho_i) + \rho_i \tilde{Y}_{i,t-1} + \eta_{it} && \equiv \tilde{\alpha}_i + \rho_i \tilde{Y}_{i,t-1} + \eta_{it}, & \text{(RIP)} \\ \tilde{Y}_{it} &= \alpha_i(1 - \rho_i) + \beta_i \rho_i + \beta_i(1 - \rho_i) h_{it} + \rho_i \tilde{Y}_{i,t-1} + \eta_{it} && \equiv \tilde{\alpha}_i + \tilde{\beta}_i h_{it} + \rho_i \tilde{Y}_{i,t-1} + \eta_{it}. & \text{(HIP)} \end{aligned}$$

These are standard random coefficient models. Note that  $h_{it}$  is strictly exogenous since it is a deterministic function of age.

## 8.3 Strategy for estimation and inference

For each model, I compute confidence intervals for  $\mathbb{E}(\rho_i)$ ,  $\text{Var}(\rho_i)$  and  $\mathbb{P}(\rho_i \leq r)$  for a grid of  $r = 0, 0.1, \dots, 0.9, 1$ . For  $\mathbb{E}(\rho_i)$  and  $\mathbb{P}(\rho_i \leq r)$ , I use moment restrictions stated in

Example 4. In particular, I use for the RIP process:

$$\mathbb{E}((\tilde{\alpha}_i + \rho_i \tilde{Y}_{i,t-1})\eta_{it}) = 0, \quad \mathbb{E}(\eta_{it}) = 0, \quad \mathbb{E}(\tilde{Y}_{i,t-1-s}\eta_{it}) = 0,$$

for  $s = 0, \dots, 5$ . I use for the HIP process:

$$\mathbb{E}((\tilde{\alpha}_i + \tilde{\beta}_i h_{it} + \rho_i \tilde{Y}_{i,t-1})\eta_{it}) = 0, \quad \mathbb{E}(\eta_{it}) = 0, \quad \mathbb{E}(\tilde{Y}_{i,t-1-r}\eta_{it}) = 0, \quad \mathbb{E}(h_{i,t-s}\eta_{it}) = 0,$$

for  $r = 0, \dots, 5$  and  $s = -5, \dots, -1, 0, 1, \dots, 5$ . These make the inner objective as a second order polynomial. I then solve the inner optimization problem as a closed-form.

I use additional moment restrictions to compute confidence intervals for  $\text{Var}(\rho_i)$ . Additional moment conditions for the RIP process are:

$$\mathbb{E}((\tilde{\alpha}_i + \rho_i \tilde{Y}_{i,t-1})^3 \eta_{it}) = 0, \quad \mathbb{E}(\tilde{\alpha}_i^k \eta_{it}) = 0, \quad \mathbb{E}(\rho_i^k \eta_{it}) = 0,$$

for  $k = 1, 2$ . Additional moments for the HIP process are:

$$\mathbb{E}((\tilde{\alpha}_i + \tilde{\beta}_i h_{it} + \rho_i \tilde{Y}_{i,t-1})^3 \eta_{it}) = 0, \quad \mathbb{E}(\tilde{\alpha}_i^k \eta_{it}) = 0, \quad \mathbb{E}(\tilde{\beta}_i^k \eta_{it}) = 0, \quad \mathbb{E}(\rho_i^k \eta_{it}) = 0,$$

for  $k = 1, 2$ . The first additional moment restriction in both models was stated in Example 5, which makes the inner objective as a fourth order polynomial. The rest of additional moment restrictions then add low order terms to the inner objective. These additional restrictions yield finite lower and upper bounds on the second moments of  $(\tilde{\alpha}_i, \tilde{\beta}_i, \rho_i)$ . The inner problem is then solved using the SDP method with hierarchy of length two.

With these moment restrictions, I compute confidence intervals using the procedure in Section 5.3. Tuning parameters for the inference procedure are the same as those in simulations, sampling  $P = 50$  points in the neighborhood of each  $\tilde{\lambda}_L(\delta_m)$  and  $\tilde{\lambda}_U(\delta_m)$ .

## 8.4 Estimation results

Confidence intervals for  $\mathbb{E}(\rho_i)$  and  $\text{Var}(\rho_i)$  are given in Table 4. Both models estimate  $\mathbb{E}(\rho_i)$  to be significantly less than 1, in contrast to the models with homogeneous  $\rho = \rho_i$  where the literature estimates  $\rho \approx 1$  for the RIP process. Moreover, their confidence intervals show substantial overlap, having similar upper confidence limits, which suggests that specifying homogeneous or heterogeneous  $\beta$  does not lead to serious misspecification when  $\rho$  is allowed to be heterogeneous. Confidence interval for  $\text{Var}(\rho_i)$  suggest heterogeneity in  $\rho_i$ , having a lower confidence limit of 0.067 for RIP process implying standard deviation of 0.258. This implies that households face different levels of earnings

risk that lead to heterogeneity in their consumption and savings behaviors.

Confidence intervals for the CDF of  $\rho_i$  for a grid of values are given in Table 5. They suggest substantial heterogeneity in  $\rho_i$ . For example, confidence intervals for RIP process suggest that at least 15% of households are estimated to have  $\rho_i \leq 0.5$ , while another at least 15% of households are estimated to have  $\rho_i > 0.5$ . Confidence intervals for the two CDFs show substantial overlap as well.

## 9 Conclusion

This paper studies identification and estimation of dynamic random coefficient models in a short panel context. The model extends the widely-used panel data linear model with fixed effects (Arellano and Bond, 1991; Blundell and Bond, 1998), by allowing for individual-specific coefficients and intercept. I show that the model is not point-identified but rather partially identified, and I characterize sharp identified sets of the parameters of interest using the dual representation of infinite-dimensional linear programming. I propose a computationally feasible estimation and inference procedure for the parameters of interest, which is based on testing many moment inequalities. Computational feasibility is achieved using a fast and exact algorithm for global polynomial optimization.

Using my method, I estimate unobserved heterogeneity in earnings persistence across U.S. households using the PSID dataset. I find that the average earnings persistence is significantly less than 1 when it is allowed to be heterogeneous. I also find evidence that, when earnings persistence is allowed to be heterogeneous, choosing RIP over HIP or vice versa may not lead to serious misspecification of the persistence. Estimates for variance and CDF of earnings persistence suggest that there is substantial degree of unobserved heterogeneity in it.

## 10 Acknowledgements

This is based on my PhD dissertation at the University of Chicago. I am indebted to Stéphane Bonhomme, Alexander Torgovitsky, and Guillaume Pouliot for their guidance and support. I thank Manuel Arellano, Timothy Armstrong, Antonio Galvao, Greg Kaplan, Roger Koenker, Zhipeng Liao, Jack Light, Azeem Shaikh, Shuyang Sheng, Panagiotis Toulis, and Ying Zhu for helpful discussions and comments. I also thank seminar participants at the Econometrics Workshop and Econometrics Working Group at the University of Chicago.



## References

- Abowd, John M and David Card. 1989. "On the covariance structure of earnings and hours changes." *Econometrica* 57 (2):411–445.
- Alan, Sule, Martin Browning, and Mette Ejrnæs. 2018. "Income and consumption: A micro semistructural analysis with pervasive heterogeneity." *Journal of Political Economy* 126 (5):1827–1864.
- Anderson, Edward J. 1983. "A review of duality theory for linear programming over topological vector spaces." *Journal of Mathematical Analysis and Applications* 97 (2):380–392.
- Andrews, Donald WK and Soonwoo Kwon. 2019. "Inference in moment inequality models that is robust to spurious precision under model misspecification." *Working paper* .
- Andrews, Donald WK and Xiaoxia Shi. 2013. "Inference based on conditional moment inequalities." *Econometrica* 81 (2):609–666.
- Arellano, Manuel, Richard Blundell, and Stéphane Bonhomme. 2017. "Earnings and consumption dynamics: a nonlinear panel data framework." *Econometrica* 85 (3):693–734.
- Arellano, Manuel and Stephen Bond. 1991. "Some tests of specification for panel data: Monte Carlo evidence and an application to employment equations." *Review of Economic Studies* 58 (2):277–297.
- Arellano, Manuel and Stéphane Bonhomme. 2012. "Identifying distributional characteristics in random coefficients panel data models." *Review of Economic Studies* 79 (3):987–1020.
- . 2018. "Recovering latent variables by matching." *Working paper* .
- Bai, Yuehao, Andres Santos, and Azeem M Shaikh. 2019. "A practical method for testing many moment inequalities." *Working Paper* .
- Baker, Michael. 1997. "Growth-rate heterogeneity and the covariance structure of life-cycle earnings." *Journal of Labor Economics* 15 (2):338–375.
- Bierens, Herman J. 1990. "A consistent conditional moment test of functional form." *Econometrica: Journal of the Econometric Society* :1443–1458.
- Blundell, Richard and Stephen Bond. 1998. "Initial conditions and moment restrictions in dynamic panel data models." *Journal of Econometrics* 87 (1):115–143.
- Blundell, Richard, Hamish Low, and Ian Preston. 2013. "Decomposing changes in income risk using consumption data." *Quantitative Economics* 4 (1):1–37.
- Blundell, Richard, Luigi Pistaferri, and Ian Preston. 2008. "Consumption inequality and partial insurance." *American Economic Review* 98 (5):1887–1921.

- Blundell, Richard, Luigi Pistaferri, and Itay Saporta-Eksten. 2016. "Consumption inequality and family labor supply." *American Economic Review* 106 (2):387–435.
- Browning, Martin, Mette Ejrnaes, and Javier Alvarez. 2010. "Modelling income processes with lots of heterogeneity." *Review of Economic Studies* 77 (4):1353–1381.
- Chamberlain, Gary. 1992. "Efficiency bounds for semiparametric regression." *Econometrica* 60 (3):567–596.
- . 1993. "Feedback in panel data models." *Working paper* .
- . 2022. "Feedback in panel data models." *Journal of Econometrics* 226 (1):4–20.
- Chernozhukov, Victor, Denis Chetverikov, and Kengo Kato. 2019. "Inference on causal and structural parameters using many moment inequalities." *Review of Economic Studies* 86 (5):1867–1900.
- Chernozhukov, Victor, Sokbae Lee, and Adam M Rosen. 2013. "Intersection bounds: Estimation and inference." *Econometrica* 81 (2):667–737.
- Galichon, Alfred and Marc Henry. 2009. "A test of non-identifying restrictions and confidence regions for partially identified parameters." *Journal of Econometrics* 152 (2):186–196.
- . 2011. "Set identification in models with multiple equilibria." *Review of Economic Studies* 78 (4):1264–1298.
- Graham, Bryan S and James L Powell. 2012. "Identification and estimation of average partial effects in "irregular" correlated random coefficient panel data models." *Econometrica* 80 (5):2105–2152.
- Gu, Jiaying and Roger Koenker. 2017. "Unobserved heterogeneity in income dynamics: An empirical Bayes perspective." *Journal of Business & Economic Statistics* 35 (1):1–16.
- Gunsilius, Florian. 2019. "Bounds in continuous instrumental variable models." *Working paper, arXiv preprint arXiv:1910.09502* .
- Guvenen, Fatih. 2007. "Learning your earning: Are labor income shocks really very persistent?" *American Economic Review* 97 (3):687–712.
- . 2009. "An empirical investigation of labor income processes." *Review of Economic dynamics* 12 (1):58–79.
- Hall, Robert E and Frederic S Mishkin. 1982. "The sensitivity of consumption to transitory income: Estimates from panel data on households." *Econometrica* 50 (2):461–481.
- Honoré, Bo E and Elie Tamer. 2006. "Bounds on parameters in panel dynamic discrete choice models." *Econometrica* 74 (3):611–629.

- Jappelli, Tullio and Luigi Pistaferri. 2010. "The consumption response to income changes." *Annual Review of Economics* 2:479–506.
- Kaplan, Greg and Giovanni L Violante. 2014. "A model of the consumption response to fiscal stimulus payments." *Econometrica* 82 (4):1199–1239.
- Kiefer, Jack. 1959. "Optimum experimental designs." *Journal of the Royal Statistical Society: Series B* 21 (2):272–304.
- Lasserre, Jean-Bernard. 2010. *Moments, positive polynomials and their applications*. World Scientific.
- . 2015. *An introduction to polynomial and semi-algebraic optimization*. Cambridge University Press.
- Li, Lixiong. 2018. "Identification of structural and counterfactual parameters in a large class of structural econometric models." *Working paper*.
- Lillard, Lee A and Yoram Weiss. 1979. "Components of variation in panel earnings data: American scientists 1960-70." *Econometrica* 47 (2):437–454.
- MaCurdy, Thomas E. 1982. "The use of time series processes to model the error structure of earnings in a longitudinal data analysis." *Journal of Econometrics* 18 (1):83–114.
- Meghir, Costas and Luigi Pistaferri. 2004. "Income variance dynamics and heterogeneity." *Econometrica* 72 (1):1–32.
- Milgrom, Paul and Ilya Segal. 2002. "Envelope theorems for arbitrary choice sets." *Econometrica* 70 (2):583–601.
- Mogstad, Magne, Andres Santos, and Alexander Torgovitsky. 2018. "Using instrumental variables for inference about policy relevant treatment parameters." *Econometrica* 86 (5):1589–1619.
- Newey, Whitney K and Daniel McFadden. 1994. "Large sample estimation and hypothesis testing." *Handbook of Econometrics* 4:2111–2245.
- Romano, Joseph P, Azeem M Shaikh, and Michael Wolf. 2014. "A practical two-step method for testing moment inequalities." *Econometrica* 82 (5):1979–2002.
- Schennach, Susanne M. 2014. "Entropic latent variable integration via simulation." *Econometrica* 82 (1):345–385.
- Topel, Robert H and Michael P Ward. 1992. "Job mobility and the careers of young men." *Quarterly Journal of Economics* 107 (2):439–479.
- Torgovitsky, Alexander. 2019. "Nonparametric inference on state dependence in unemployment." *Working paper, Available at SSRN:2564305*.
- Van der Vaart, Aad W. 2000. *Asymptotic statistics*, vol. 3. Cambridge university press.

Wooldridge, Jeffrey M. 2005. "Fixed-effects and related estimators for correlated random-coefficient and treatment-effect panel data models." *Review of Economics and Statistics* 87 (2):385–390.

# Appendices

## A Proofs

**Proof of Proposition 1.** This proof uses Lemma 3 in the Online Appendix, which is a straightforward extension of Lemma 2 in Section 4 that also incorporates conditional moment restrictions. A key assumption for the lemma, Assumption 7 (iv) in the Online Appendix, is satisfied by the assumption that the joint density of  $(Y_{i0}, Y_{i1}, Y_{i2}, \gamma_i, \beta_i)$  is strictly positive (Anderson, 1983, Theorem 9).

For notational simplicity, assume  $\mathcal{C} = \mathcal{C}_0^5$  where  $\mathcal{C}_0$  is a compact subset of  $\mathbb{R}$ . The proof can be easily modified for a generic compact set  $\mathcal{C}$ .

Suppose that  $\mathbb{E}(\beta_i)$  is point-identified, from which I draw contradiction. Lemma 3 in the Online Appendix tells that, if  $\mathbb{E}(\beta_i)$  is point-identified, there exists  $f^* : \mathcal{C}_0^3 \mapsto \mathbb{R}$ ,  $g_1^* : \mathcal{C}_0^3 \mapsto \mathbb{R}$  and  $g_2^* : \mathcal{C}_0^4 \mapsto \mathbb{R}$  which are linear functionals on the spaces of finite and countably additive signed Borel measures that are absolutely continuous with respect to the Lebesgue measure such that

$$f^*(Y_{i0}, Y_{i1}, Y_{i2}) + g_1^*(\gamma_i, \beta_i, Y_{i0})\varepsilon_{i1} + g_2^*(\gamma_i, \beta_i, Y_{i0}, Y_{i1})\varepsilon_{i2} = \beta_i \quad (34)$$

almost surely on  $\mathcal{C}_0^5$ . Substituting  $\varepsilon_{it} = Y_{it} - \gamma_i - \beta_i Y_{i,t-1}$  in (34) yields:

$$f^*(Y_{i0}, Y_{i1}, Y_{i2}) + g_1^*(\gamma_i, \beta_i, Y_{i0})(Y_{i1} - \gamma_i - \beta_i Y_{i0}) + g_2^*(\gamma_i, \beta_i, Y_{i0}, Y_{i1})(Y_{i2} - \gamma_i - \beta_i Y_{i1}) = \beta_i. \quad (35)$$

Consider any  $\gamma, \tilde{\gamma}, \beta, y_0, y_1, y_2 \in \mathcal{C}_0$  such that  $\gamma \neq \tilde{\gamma}$ . I evaluate (35) at  $(\gamma, \beta, y_0, y_1, y_2)$  and at  $(\tilde{\gamma}, \beta, y_0, y_1, y_2)$ , and I take difference:

$$\begin{aligned} & (y_1 - \tilde{\gamma} - \beta y_0)\Delta_{\tilde{\gamma}, \gamma} g_1^* - (\tilde{\gamma} - \gamma)g_1^*(\gamma, \beta, y_0) \\ & + (y_2 - \tilde{\gamma} - \beta y_1)\Delta_{\tilde{\gamma}, \gamma} g_2^* - (\tilde{\gamma} - \gamma)g_2^*(\gamma, \beta, y_0, y_1) = 0 \end{aligned} \quad (36)$$

where  $\Delta_{\tilde{\gamma}, \gamma} g_1^* = g_1^*(\tilde{\gamma}, \beta, y_0) - g_1^*(\gamma, \beta, y_0)$  and  $\Delta_{\tilde{\gamma}, \gamma} g_2^*$  is defined similarly.

In (36),  $y_2$  only appears in the third term. Also, (36) must hold almost surely for all

$\gamma, \tilde{\gamma}, \beta, y_0, y_1, y_2 \in \mathcal{C}_0$  such that  $\gamma \neq \tilde{\gamma}$ , and in particular for any  $y_2 \in \mathcal{C}_0$ . These imply that, almost surely:

$$\Delta_{\tilde{\gamma}, \gamma} g_2^* = 0,$$

which means that  $g_2^*$  is almost surely a constant function over  $\gamma$ :

$$g_2^*(\gamma, \beta, y_0, y_1) = g_2^*(\beta, y_0, y_1). \quad (37)$$

If not, i.e. if  $\Delta_{\tilde{\gamma}, \gamma} g_2^* \neq 0$  on a subset of  $\mathcal{C}_0^6$  with positive Lebesgue measure, one can change the value of  $y_2$  without changing  $(\gamma, \tilde{\gamma}, \beta, y_0, y_1)$  within this subset to violate (36) with positive measure. Next, consider any  $\gamma, \beta, \tilde{\beta}, y_0, y_1, y_2 \in \mathcal{C}_0$  such that  $\beta \neq \tilde{\beta}$ . I evaluate (35) at  $(\gamma, \beta, y_0, y_1, y_2)$  and  $(\gamma, \tilde{\beta}, y_0, y_1, y_2)$ , and I take difference:

$$\begin{aligned} & (y_1 - \gamma - \tilde{\beta}y_0)\Delta_{\tilde{\beta}, \beta} g_1^* - (\tilde{\beta} - \beta)y_0 g_1^*(\gamma, \beta, y_0) \\ & + (y_2 - \gamma - \tilde{\beta}y_1)\Delta_{\tilde{\beta}, \beta} g_2^* - (\tilde{\beta} - \beta)y_1 g_2^*(\beta, y_0, y_1) = \tilde{\beta} - \beta. \end{aligned} \quad (38)$$

where  $\Delta_{\tilde{\beta}, \beta} g_1^* = g_1^*(\gamma, \tilde{\beta}, y_0) - g_1^*(\gamma, \beta, y_0)$  and  $\Delta_{\tilde{\beta}, \beta} g_2^*$  is defined similarly. In (38),  $y_2$  only appears in the third term. This implies  $g_2^*(\beta, y_0, y_1) = g_2^*(y_0, y_1)$  almost surely, similarly to the argument for (37). Then (36) simplifies to:

$$(y_1 - \tilde{\gamma} - \beta y_0)\Delta_{\tilde{\gamma}, \gamma} g_1^* - (\tilde{\gamma} - \gamma)g_1^*(\gamma, \beta, y_0) - (\tilde{\gamma} - \gamma)g_2^*(y_0, y_1) = 0. \quad (39)$$

Let  $\hat{\gamma} \in \mathcal{C}$  be such that  $\hat{\gamma} - \tilde{\gamma} = \tilde{\gamma} - \gamma$ . I evaluate (39) at  $(\gamma, \tilde{\gamma}, \beta, y_0, y_1)$  and  $(\tilde{\gamma}, \hat{\gamma}, \beta, y_0, y_1)$ , and I take difference:

$$(y_1 - \hat{\gamma} - \beta y_0)(\Delta_{\hat{\gamma}, \tilde{\gamma}} g_1^* - \Delta_{\tilde{\gamma}, \gamma} g_1^*) - (\hat{\gamma} - \tilde{\gamma})\Delta_{\tilde{\gamma}, \gamma} g_1^* - (\tilde{\gamma} - \gamma)\Delta_{\tilde{\gamma}, \gamma} g_1^* = 0. \quad (40)$$

In (40),  $y_1$  only appears in the first term, which implies  $\Delta_{\hat{\gamma}, \tilde{\gamma}} g_1^* - \Delta_{\tilde{\gamma}, \gamma} g_1^* = 0$  almost surely. (40) then simplifies to:

$$(\hat{\gamma} - \tilde{\gamma})\Delta_{\tilde{\gamma}, \gamma} g_1^* + (\tilde{\gamma} - \gamma)\Delta_{\tilde{\gamma}, \gamma} g_1^* = 0, \quad (41)$$

which implies  $\Delta_{\tilde{\gamma}, \gamma} g_1^* = 0$  since  $\hat{\gamma} - \tilde{\gamma} = \tilde{\gamma} - \gamma \neq 0$ . This means that  $g_1^*$  is almost surely a constant function over  $\gamma$ , i.e.  $g_1^*(\gamma, \beta, y_0) = g_1^*(\beta, y_0)$ . Then (38) simplifies to:

$$\begin{aligned} & (y_1 - \gamma - \tilde{\beta}y_0)\Delta_{\tilde{\beta}, \beta} g_1^* - (\tilde{\beta} - \beta)y_0 g_1^*(\beta, y_0) \\ & + (y_2 - \gamma - \tilde{\beta}y_1)\Delta_{\tilde{\beta}, \beta} g_2^* - (\tilde{\beta} - \beta)y_1 g_2^*(y_0, y_1) = \tilde{\beta} - \beta. \end{aligned} \quad (42)$$

Let  $\hat{\beta} \in \mathcal{C}$  be such that  $\hat{\beta} - \tilde{\beta} = \tilde{\beta} - \beta$ . Evaluating (42) with  $(\gamma, \hat{\beta}, \tilde{\beta}, y_0, y_1, y_2)$  and  $(\gamma, \tilde{\beta}, \beta, y_0, y_1, y_2)$  and taking difference yields  $g_1^*(\beta, y_0) = g_1^*(y_0)$ , similarly to the argument for  $g_1^*(\gamma, \beta, y_0) = g_1^*(\beta, y_0)$  from (41). Then (35) simplifies to:

$$f^*(y_0, y_1, y_2) + g_1^*(y_0)(y_1 - \gamma - \beta y_0) + g_2^*(y_0, y_1)(y_2 - \gamma - \beta y_1) = \beta$$

almost surely for all  $(\gamma, \beta, y_0, y_1, y_2)$ . This is a linear identity in  $(\gamma, \beta)$ , so their coefficients must coincide on both sides. I then obtain  $g_1^* + g_2^* = 0$  from the coefficients on  $\gamma$  and  $-y_0 g_1^* - y_1 g_2^* = 1$  from the coefficients on  $\beta$ . Solving these for  $(g_1^*, g_2^*)$  yields, almost surely:

$$g_1^* = \frac{1}{y_1 - y_0}, \quad g_2^* = \frac{-1}{y_1 - y_0}.$$

However,  $g_1^*$  cannot be a function of  $y_1$ , which is contradiction.  $\square$

**Proof of Lemma 1.** As discussed in the proof of Proposition 1, Lemma 3 in the Online Appendix tells that, if  $\mathbb{E}(\beta_i)$  is point-identified, there exists  $(f^*, g_1^*, g_2^*)$  such that (34) holds almost surely on  $\mathcal{C}_0^5$ . Then (34) implies  $S^*(Y_{i0}, Y_{i1}, Y_{i2}) = f^*(Y_{i0}, Y_{i1}, Y_{i2})$  because:

$$\mathbb{E}(f^*(Y_{i0}, Y_{i1}, Y_{i2})|\beta_i) = \mathbb{E}(\beta_i - g_1^*(\gamma_i, \beta_i, Y_{i0})\varepsilon_{i1} - g_2^*(\gamma_i, \beta_i, Y_{i0}, Y_{i1})\varepsilon_{i2}|\beta_i) = \beta_i.$$

Conversely, if there exists  $S^*(Y_{i0}, Y_{i1}, Y_{i2})$  such that  $\mathbb{E}(S^*(Y_{i0}, Y_{i1}, Y_{i2})|\beta_i) = \beta_i$ , then  $\mathbb{E}(S^*(Y_{i0}, Y_{i1}, Y_{i2})) = \mathbb{E}(\mathbb{E}(S^*(Y_{i0}, Y_{i1}, Y_{i2})|\beta_i)) = \mathbb{E}(\beta_i)$ , which completes the proof.  $\square$

**Proof of Theorem 1.** For the proof, it suffices to show the sharpness of  $[L, U]$  when Assumption 3 is replaced with (9). In what follows, I show that  $U$  is the sharp upper bound under (9). The same argument applies to  $L$ . This proof relies on Theorem 2 in Section 4.

By Theorem 2, the sharp upper bound is:

$$\min_{\lambda, \mu} \mathbb{E} \left( \max_v \left[ e'v + \mu' \sum_{t=1}^T R_{it}(Y_{it} - R'_{it}v) + \lambda \sum_{t=1}^T (R'_{it}v)(Y_{it} - R'_{it}v) \right] \right)$$

where  $\mu$  has the same dimension as  $R_{it}$ , and  $\lambda$  is scalar. With the notation of  $\mathcal{R}_i$  and  $\mathcal{Y}_i$  in the statement of Theorem 1, I can write the above concisely as

$$\min_{\mu, \lambda} \mathbb{E} \left( \max_v \left[ e'v + \mu' \mathcal{Y}_i - \mu' \mathcal{R}_i v + \lambda \mathcal{Y}_i' v - v'(\lambda \mathcal{R}_i) v \right] \right).$$

The objective function of the inner maximization problem is a quadratic polynomial in  $v$ . As discussed in Section 6.1.1, it suffices to consider  $\lambda > 0$  which yields a negative definite

second derivative matrix. The inner maximization problem then has a closed-form:

$$\min_{\lambda > 0, \mu} \mathbb{E} \left( \mu' \mathcal{Y}_i + \frac{1}{4\lambda} [e + \lambda \mathcal{Y}_i - \mathcal{R}_i \mu]' \mathcal{R}_i^{-1} [e + \lambda \mathcal{Y}_i - \mathcal{R}_i \mu] \right). \quad (43)$$

I solve this problem with respect to  $\mu$  for a fixed  $\lambda$ . The first order condition with respect to  $\mu$  given  $\lambda$  is:

$$\mathbb{E}(\mathcal{Y}_i) + \frac{1}{2\lambda} \mathbb{E}(\mathcal{R}_i) \mu - \frac{1}{2\lambda} e - \frac{1}{2} \mathbb{E}(\mathcal{Y}_i) = 0.$$

The optimal  $\mu$  that solves the first order solution is  $\mu^* = \mathbb{E}(\mathcal{R}_i)^{-1} [e - \lambda \mathbb{E}(\mathcal{Y}_i)]$ . I substitute this into (43) and solve the problem with respect to  $\lambda$ . The first order condition is:

$$\frac{1}{\lambda^2} \left[ e' \mathbb{E}(\mathcal{R}_i)^{-1} e - e' \mathbb{E}(\mathcal{R}_i^{-1}) e \right] = \mathbb{E}(\mathcal{Y}_i)' \mathbb{E}(\mathcal{R}_i)^{-1} \mathbb{E}(\mathcal{Y}_i) - \mathbb{E}(\mathcal{Y}_i' \mathcal{R}_i^{-1} \mathcal{Y}_i).$$

Since  $\lambda > 0$ , the optimal  $\lambda$  that solves the first order condition is:

$$\lambda^* = \sqrt{\frac{e' \mathbb{E}(\mathcal{R}_i^{-1}) e - e' \mathbb{E}(\mathcal{R}_i)^{-1} e}{\mathbb{E}(\mathcal{Y}_i' \mathcal{R}_i^{-1} \mathcal{Y}_i) - \mathbb{E}(\mathcal{Y}_i)' \mathbb{E}(\mathcal{R}_i)^{-1} \mathbb{E}(\mathcal{Y}_i)}}. \quad (44)$$

Substituting (44) into (43) yields the expression of  $\tilde{U}$  in Theorem 1.

The numerator and the denominator in (44) are both weakly positive, and they are zero if and only if  $\mathcal{R}_i^{-1} e$  and  $\mathcal{R}_i^{-1} \mathcal{Y}_i$  are degenerate across individuals, respectively. To show this, define the functions  $\mathcal{E}(\mathcal{R}_i) = e' \mathcal{R}_i^{-1} e$  and  $\mathcal{D}(\mathcal{Y}_i, \mathcal{R}_i) = \mathcal{Y}_i' \mathcal{R}_i^{-1} \mathcal{Y}_i$ , and apply the following proposition to  $\mathcal{E}$  and  $\mathcal{D}$ .  $\square$

**Proposition 4** (Kiefer, 1959, Lemma 3.2). *For an integer  $l > 0$ , let  $A_1, \dots, A_l$  be  $n \times m$  matrices and  $B_1, \dots, B_l$  be nonsingular positive definite and symmetric  $n \times n$  matrices. Let  $a_1, \dots, a_l$  be positive real numbers such that  $\sum_k a_k = 1$ . Then*

$$\sum_{k=1}^l a_k A_k' B_k^{-1} A_k - \left[ \sum_{k=1}^l a_k A_k \right]' \left[ \sum_{k=1}^l a_k B_k \right]^{-1} \left[ \sum_{k=1}^l a_k A_k \right] \geq 0$$

where ' $\geq$ ' is the partial ordering defined in terms of positive semidefinite matrices. In addition, the equality holds if and only if  $B_1^{-1} A_1 = \dots = B_l^{-1} A_l$ .

**Proof of Theorem 2.** In what follows, I show that (11) is the dual problem of (10), by applying the duality theorem of linear programming for topological vector spaces (Anderson, 1983). The same argument applies to (12).

To apply the theorem, I first rewrite (11) as a standard form of linear programming,

for which I introduce additional notation. Recall that  $\mathcal{M}_{W \times V}$  is a linear space of finite and countably additive signed Borel measures on  $\mathcal{W} \times \mathcal{V}$ . Let  $\overline{\mathcal{F}}_{W \times V}$  be the dual space of  $\mathcal{M}_{W \times V}$ , and let  $\mathcal{F}_{W \times V}$  be the space of all bounded Borel measurable functions on  $\mathcal{W} \times \mathcal{V}$ . Note that  $\mathcal{F}_{W \times V}$  is a linear subspace of  $\overline{\mathcal{F}}_{W \times V}$ .

For  $P \in \mathcal{M}_{W \times V}$  and  $f \in \overline{\mathcal{F}}_{W \times V}$ , define the *dual pairing*

$$\langle P, f \rangle = \int f dP.$$

Let  $\mathcal{M}_W$  be the linear space of finite and countably additive signed Borel measures on  $\mathcal{W}$ . Let  $\overline{\mathcal{F}}_W$  be the dual space of  $\mathcal{M}_W$ , and let  $\mathcal{F}_W$  be the space of all bounded Borel measurable functions on  $\mathcal{W}$ . Note that  $\mathcal{F}_W$  is a linear subspace of  $\overline{\mathcal{F}}_W$ . In addition, define  $\mathcal{G} = \mathbb{R}^K \times \mathcal{M}_W$  and  $\mathcal{H} = \mathbb{R}^K \times \overline{\mathcal{F}}_W$ , and let  $g = (g_1, \dots, g_K, P_g)$  and  $h = (\lambda_1, \dots, \lambda_K, f_h)$  be their generic elements. Then  $\mathcal{H}$  is the dual space of  $\mathcal{G}$ . Define the dual pairing

$$\langle g, h \rangle = \sum_{k=1}^K g_k \lambda_k + \int f_h dP_g.$$

Next, define a linear map  $A : \mathcal{M}_{W \times V} \mapsto \mathcal{G}$  by

$$A(P) = \left( \int \phi_1 dP, \dots, \int \phi_K dP, P(\cdot, \mathcal{V}) \right).$$

$A$  is a bounded (thus continuous) linear operator because  $\phi_k$ s are assumed to be bounded. Note that

$$\langle A(P), h \rangle = \sum_{k=1}^K \lambda_k \int \phi_k dP + \int_{\mathcal{W}} f_h(w) P(dw, \mathcal{V}).$$

It is straightforward to show that:

$$\int_{\mathcal{W}} f_h(w) P(dw, \mathcal{V}) = \int_{\mathcal{W} \times \mathcal{V}} f_h(w) dP(w, v).$$

Then:

$$\langle A(P), h \rangle = \sum_{k=1}^K \lambda_k \int \phi_k dP + \int f_h dP = \int \left[ \sum_{k=1}^K \lambda_k \phi_k + f_h \right] dP \equiv \langle P, A^*(h) \rangle, \quad (45)$$

where  $A^*(h) : \mathcal{H} \mapsto \overline{\mathcal{F}}_{W \times V}$  is defined as

$$A^*(h) = \sum_{k=1}^K \lambda_k \phi_k + f_h.$$



(45) shows that  $A^*$  is the adjoint of  $A$ . With these notation, I rewrite (10) as a standard form of linear programming:

$$\min_{P \in \mathcal{M}_{W \times V}} \langle P, m \rangle \quad \text{subject to} \quad A(P) = c, \quad P \geq 0, \quad (46)$$

where  $c = (0, \dots, 0, P_W)$ . I then apply the strong duality theorem under Assumption 5 and the continuity of  $A$  (Anderson, 1983, Theorem 6), which implies that the optimal solution to (46) equals to the solution to:

$$\max_{h \in \mathcal{H}} \langle c, h \rangle \quad \text{subject to} \quad m - A^*(h) \geq 0, \quad P \geq 0,$$

which I can write more concretely as:

$$\max_{\lambda_1, \dots, \lambda_K \in \mathbb{R}, f_h \in \overline{\mathcal{F}}_W} \int f_h dP_W \quad \text{subject to} \quad \sum_{k=1}^K \lambda_k \phi_k + f_h \leq m. \quad (47)$$

Now I show that (47) simplifies to (11). I rearrange the constraint of (47):

$$f_h(w) \leq m(w, v) - \sum_{k=1}^K \lambda_k \phi_k(w, v).$$

The left-hand side does not involve  $v$ . Therefore:

$$f_h(w) \leq \min_{v \in \mathcal{V}} \left[ m(w, v) - \sum_{k=1}^K \lambda_k \phi_k(w, v) \right] \quad \text{for all } w \in \mathcal{W}.$$

Since (47) maximizes the expectation of  $f_h(w)$ , the optimal  $f_h$  for a fixed  $(\lambda_1, \dots, \lambda_K)$  is given by:

$$f_h^*(w) = \min_{v \in \mathcal{V}} \left[ m(w, v) - \sum_{k=1}^K \lambda_k \phi_k(w, v) \right] \quad (48)$$

almost surely on  $\mathcal{W}$ . If not, i.e. if  $f_h^*(w)$  is strictly less than the right-hand side of (48) with positive probability, one can increase the value of the objective by increasing  $f_h^*$  on a set of positive measure. I then substitute (48) into (47), which yields:

$$\max_{\lambda_1, \dots, \lambda_K \in \mathbb{R}} \int \min_{v \in \mathcal{V}} \left[ m(w, v) - \sum_{k=1}^K \lambda_k \phi_k(w, v) \right] dP_W(w).$$

The above display remains equivalent even if the signs of  $(\lambda_1, \dots, \lambda_K)$  are switched because the  $\lambda$ s are choice variables supported on  $\mathbb{R}^K$ . Switching the signs of  $\lambda$ s in the above

gives:

$$\max_{\lambda_1, \dots, \lambda_K \in \mathbb{R}} \int \min_{v \in \mathcal{V}} \left[ m(w, v) + \sum_{k=1}^K \lambda_k \phi_k(w, v) \right] dP_W(w)$$

which is the expression in (11).  $\square$

**Proof of Lemma 2** As in (47) in the proof of Theorem 2, the sharp lower bound of  $\theta$  is given by

$$\max_{\lambda_1, \dots, \lambda_K \in \mathbb{R}, f_h \in \overline{\mathcal{F}}_W} \int f_h dP_W \quad \text{subject to} \quad \sum_{k=1}^K \lambda_k \phi_k + f_h \leq m \quad (49)$$

where all notation follow the proof of Theorem 2. Similarly, the sharp upper bound of  $\theta$  is given by

$$\min_{\lambda_1, \dots, \lambda_K \in \mathbb{R}, f_h \in \overline{\mathcal{F}}_W} \int f_h dP_W \quad \text{subject to} \quad \sum_{k=1}^K \lambda_k \phi_k + f_h \geq m. \quad (50)$$

Suppose that  $\theta$  is point-identified but there is no such  $S^* \in \overline{\mathcal{F}}_W$  and  $\lambda_1^*, \dots, \lambda_K^* \in \mathbb{R}$  such that, almost surely:

$$\sum_{k=1}^K \lambda_k^* \phi_k + S^* = m.$$

Then the solution  $(\lambda_1^l, \dots, \lambda_K^l, S^l)$  to (49) satisfies its constraint  $\sum_{k=1}^K \lambda_k^l \phi_k + S^l \leq m$  with strict inequality on a set with positive Lebesgue measure on  $\mathcal{W} \times \mathcal{V}$ . Similarly, the solution  $(\lambda_1^u, \dots, \lambda_K^u, S^u)$  to (50) satisfies its constraint  $\sum_{k=1}^K \lambda_k^u \phi_k + S^u \geq m$  with strict inequality on a set with positive Lebesgue measure on  $\mathcal{W} \times \mathcal{V}$ . Then:

$$\mathbb{E}(S^l) = \mathbb{E} \left( \sum_{k=1}^K \lambda_k^l \phi_k + S^l \right) < \mathbb{E}(m) < \mathbb{E} \left( \sum_{k=1}^K \lambda_k^u \phi_k + S^u \right) = \mathbb{E}(S^u)$$

where strict inequalities follow because the density of  $(W_i, V_i)$  is strictly positive. The above implies that the sharp lower bound  $\mathbb{E}(S^l)$  is strictly less than the sharp upper bound  $\mathbb{E}(S^u)$ , which is contradiction since  $\theta$  is assumed to be point-identified.

Conversely, suppose there exists  $(S^*, \lambda_1^*, \dots, \lambda_K^*)$  such that  $\sum_{k=1}^K \lambda_k^* \phi_k + S^* = m$ . Then:

$$\mathbb{E}(S^*) = \mathbb{E} \left( \sum_{k=1}^K \lambda_k^* \phi_k + S^* \right) = \mathbb{E}(m) = \theta,$$

which proves that  $\theta$  is point-identified by  $\mathbb{E}(S^*)$ .  $\square$

**Proof of Proposition 2.** It suffices to show that  $G_L$  is concave in  $\lambda$ . Let  $\lambda_1 = (\lambda_{11}, \dots, \lambda_{1K})$  and  $\lambda_2 = (\lambda_{21}, \dots, \lambda_{2K})$  be two distinct points in  $\mathbb{R}^K$ . Then, for any  $t \in [0, 1]$  and  $w \in \mathcal{W}$ :

$$\begin{aligned}
& G_L(t\lambda_1 + (1-t)\lambda_2, w) \\
&= \min_{v \in \mathcal{V}} \left\{ t \left[ m(w, v) + \sum_{k=1}^K \lambda_{1k} \phi_k(w, v) \right] + (1-t) \left[ m(w, v) + \sum_{k=1}^K \lambda_{2k} \phi_k(w, v) \right] \right\} \\
&\geq t \min_{v \in \mathcal{V}} \left\{ m(w, v) + \sum_{k=1}^K \lambda_{1k} \phi_k(w, v) \right\} + (1-t) \min_{v \in \mathcal{V}} \left\{ m(w, v) + \sum_{k=1}^K \lambda_{2k} \phi_k(w, v) \right\} \\
&= tG_L(\lambda_1, w) + (1-t)G_L(\lambda_2, w).
\end{aligned}$$

This is the definition of concavity.  $\square$

**Proof of Proposition 3.** When  $L$  exists and is finite, Proposition 2 implies that  $L(\lambda)$  is concave. Then  $\hat{L}$  uniformly converges to  $L$  on any compact set  $K \subseteq \mathbb{R}^K$ , as in the proof of Theorem 2.7 in Newey and McFadden (1994):

$$\sup_{\lambda \in K} |\hat{L}(\lambda) - L(\lambda)| \xrightarrow{p} 0. \quad (51)$$

Let  $\hat{\lambda} = \operatorname{argmax}_{\lambda} \hat{L}(\lambda)$  and  $\lambda_0 = \operatorname{argmax}_{\lambda} L(\lambda)$ . If there are multiple  $\operatorname{argmax}$ 's, choose any of them. Then for  $\hat{\lambda}$  that is on a compact set  $K \subseteq \mathbb{R}^K$ :

$$\begin{aligned}
|L(\lambda_0) - \hat{L}(\hat{\lambda})| &\leq L(\lambda_0) - L(\hat{\lambda}) + |L(\hat{\lambda}) - \hat{L}(\hat{\lambda})| && \text{(triangle inequality)} \\
&= \hat{L}(\lambda_0) - L(\hat{\lambda}) + |L(\hat{\lambda}) - \hat{L}(\hat{\lambda})| + o_p(1) && \text{(by (51))} \\
&\leq \hat{L}(\hat{\lambda}) - L(\hat{\lambda}) + |L(\hat{\lambda}) - \hat{L}(\hat{\lambda})| + o_p(1) && (\hat{\lambda} \text{ is argmax}) \\
&\leq 2|L(\hat{\lambda}) - \hat{L}(\hat{\lambda})| + o_p(1) = o_p(1). && \text{(by (51))}
\end{aligned}$$

Let  $\Lambda_0$  be the set of all  $\operatorname{argmax}_{\lambda} L(\lambda)$ . Let  $K_0$  be a compact set containing an open neighborhood of  $\Lambda_0$  with radius  $\varepsilon > 0$ . If such  $\varepsilon$  does not exist,  $L(\lambda)$  is a constant function, in which case consistency is immediate. If such  $\varepsilon$  exists, by Theorem 5.14 of Van der Vaart (2000):

$$\mathbb{P}(\tilde{d}(\hat{\lambda}, \Lambda_0) \geq \varepsilon \wedge \hat{\lambda} \in K_0) \longrightarrow 0$$

where  $\tilde{d}(\hat{\lambda}, \Lambda_0) = \inf\{d(\hat{\lambda}, \lambda) \mid \lambda \in \Lambda_0\}$  and  $d$  is Euclidean distance. This implies  $\hat{\lambda} \in K_0$  with probability approaching to one.  $\square$

	$T = 5$	$T = 10$	$T = 15$
$L = 3$	[0.195, 0.834]	[0.317, 0.738]	[0.371, 0.685]
$L = 5$	[0.198, 0.825]	[0.319, 0.720]	[0.372, 0.675]
$L = 7$	[0.196, 0.823]	[0.319, 0.728]	[0.368, 0.669]

Table 1: Bounds for  $\mathbb{E}(\beta_i)$  computed from the 100,000 observations from parametric model, for each  $T$  and  $L$ . I use the 100,000 observations as finite populations for which the above are population bounds for  $\mathbb{E}(\beta_i)$ .

$(T = 5)$	$L = 3$	$L = 5$	$L = 7$
$P = 25$	0.956	0.951	0.946
$P = 50$	0.971	0.963	0.959
$P = 75$	0.974	0.966	0.971
$(T = 10)$	$L = 3$	$L = 5$	$L = 7$
$P = 25$	0.933	0.929	0.894
$P = 50$	0.949	0.955	0.919
$P = 75$	0.960	0.962	0.937
$(T = 15)$	$L = 3$	$L = 5$	$L = 7$
$P = 25$	0.986	0.929	0.874
$P = 50$	0.990	0.947	0.915
$P = 75$	0.991	0.961	0.939

Table 2: Coverage probabilities of the inference procedure with sample size of  $N = 750$ . Nominal coverage probability is 0.9.

$(T = 5)$	$L = 3$	$L = 5$	$L = 7$
$P = 25$	0.959	0.943	0.949
$P = 50$	0.968	0.957	0.960
$P = 75$	0.972	0.964	0.966
$(T = 10)$	$L = 3$	$L = 5$	$L = 7$
$P = 25$	0.869	0.890	0.885
$P = 50$	0.923	0.916	0.925
$P = 75$	0.934	0.930	0.938
$(T = 15)$	$L = 3$	$L = 5$	$L = 7$
$P = 25$	0.974	0.880	0.821
$P = 50$	0.988	0.910	0.869
$P = 75$	0.990	0.933	0.904

Table 3: Coverage probabilities of the inference procedures with sample size of  $N = 1000$ . Nominal coverage probability is 0.9.

	Confidence interval of $\mathbb{E}(\rho_i)$	Confidence interval of $\text{Var}(\rho_i)$
RIP process	[0.456, 0.615]	[0.067, 0.292]
HIP process	[0.264, 0.583]	[0.000, 0.701]

Table 4: Confidence interval for  $\mathbb{E}(\rho_i)$  and  $\text{Var}(\rho_i)$ . Nominal coverage probability is 0.9.

$\mathbb{P}(\rho_i \leq r)$	RIP process	HIP process
$r = 0.0$	[0.000, 0.362]	[0.000, 0.752]
$r = 0.1$	[0.005, 0.428]	[0.004, 0.800]
$r = 0.2$	[0.025, 0.548]	[0.099, 0.818]
$r = 0.3$	[0.066, 0.627]	[0.085, 0.834]
$r = 0.4$	[0.104, 0.713]	[0.135, 0.879]
$r = 0.5$	[0.153, 0.848]	[0.205, 0.914]
$r = 0.6$	[0.209, 0.895]	[0.200, 0.983]
$r = 0.7$	[0.290, 0.944]	[0.286, 0.986]
$r = 0.8$	[0.372, 0.975]	[0.319, 1.000]
$r = 0.9$	[0.471, 0.994]	[0.353, 1.000]
$r = 1.0$	[0.550, 1.000]	[0.427, 1.000]

Table 5: Confidence intervals for  $\mathbb{P}(\rho_i \leq r)$ . Nominal coverage probability is 0.9.

THIS PAGE IS INTENTIONALLY LEFT BLANK

## B Online Appendix

### B.1 Extension to multivariate random coefficient models

Results from this paper also extend to a multivariate version of (1), a system of random coefficient models:

$$\mathbf{Y}_{it} = \mathbf{Z}_{it}'\boldsymbol{\gamma}_i + \mathbf{X}_{it}'\boldsymbol{\beta}_i + \mathbf{e}_{it},$$

where  $\mathbf{Y}_{it}$  is a  $D \times 1$  vector of dependent variables,  $\mathbf{Z}_{it}$  is a  $D \times q$  matrix of strictly exogenous regressors,  $\mathbf{X}_{it}$  is a  $D \times p$  matrix of sequentially exogenous regressors, and  $\mathbf{e}_{it}$  is a  $D \times 1$  vector of idiosyncratic error terms. Assume:

$$\mathbb{E}(\mathbf{e}_{it} | \boldsymbol{\gamma}_i, \boldsymbol{\beta}_i, \mathbf{Z}_i, \mathbf{X}_i^t) = 0,$$

which is a multivariate extension of (2). The following is an example of a multivariate random coefficient models.

**Example 7** (Joint model of household earnings and consumption behavior). A researcher can combine (3) and (4) in Examples 1 and 2 and consider a joint lifecycle model of earnings and consumption behavior. If I combine the time  $t$  consumption equation and the time  $t + 1$  earnings equation, I obtain a system of random coefficient models:

$$\begin{aligned} C_{it} &= \gamma_{i1} + \gamma_{i2}Y_{it} + \beta_{i1}A_{it} + v_{it}, \\ Y_{i,t+1} &= \gamma_{i3} + \beta_{i2}Y_{it} + \varepsilon_{it}, \end{aligned}$$

which can be written in matrix form:

$$\begin{pmatrix} C_{it} \\ Y_{i,t+1} \end{pmatrix} = \begin{pmatrix} 1 & Y_{it} & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} \gamma_{i1} \\ \gamma_{i2} \\ \gamma_{i3} \end{pmatrix} + \begin{pmatrix} A_{it} & 0 \\ 0 & Y_{it} \end{pmatrix} \begin{pmatrix} \beta_{i1} \\ \beta_{i2} \end{pmatrix} + \begin{pmatrix} v_{it} \\ \varepsilon_{it} \end{pmatrix}.$$

In this model,  $\gamma$ s and  $\beta$ s can freely correlate among themselves and to  $(Y_{i0}, A_{i1})$ , allowing for correlation between earnings and consumption processes.  $\square$

### B.2 Identification under conditional moment restrictions

This section studies moment equality models that involve both conditional and unconditional moment restrictions. Consider the following extension of Assumption 4:

**Assumption 6.** The random vectors  $(W_i, V_i)$  satisfy:

$$\begin{aligned}\mathbb{E}(\phi_k(W_i, V_i)) &= 0, \quad k = 1, \dots, K_U, \\ \mathbb{E}(\psi_k(W_i, V_i)|A_{ik}) &= 0, \quad k = 1, \dots, K_C,\end{aligned}$$

where  $\phi_k, \psi_k : \mathcal{W} \times \mathcal{V} \mapsto \mathbb{R}$  are moment functions,  $A_{i1}, \dots, A_{iK_C}$  are subvectors of  $(W_i, V_i)$  and  $K_U, K_C \in \mathbb{N}$  are the number of unconditional and conditional moments, respectively.

Under Assumption 6, I characterize the identified set of

$$\theta = \mathbb{E}(m(W_i, V_i))$$

for some known function  $m : \mathcal{W} \times \mathcal{V} \mapsto \mathbb{R}$ . To simplify notation, let  $A'_{ik}$  be the subvector of  $(W_i, V_i)$  that collects the variables not included in  $A_{ik}$ , making  $(A_{ik}, A'_{ik})$  a rearrangement of  $(W_i, V_i)$ . I then write any function  $f(w, v)$  on  $\mathcal{W} \times \mathcal{V}$  equivalently as  $f(a_k, a'_k)$  on  $\mathcal{A}_k \times \mathcal{A}'_k$ , where  $\mathcal{A}_k \times \mathcal{A}'_k$  is the rearrangement of  $\mathcal{W} \times \mathcal{V}$  according to  $(A_{ik}, A'_{ik})$ .

I assume the following regularity conditions:

**Assumption 7.** The following conditions hold.

- (i)  $\mathcal{W} \times \mathcal{V}$  is a compact set in a Euclidean space.
- (ii) The distribution of  $(W_i, V_i)$  is absolutely continuous with respect to Lebesgue measure. In addition, its density  $p$  is  $L^\infty$  with respect to the Lebesgue measure.
- (iii)  $(m, \phi_1, \dots, \phi_{K_U}, \psi_1, \dots, \psi_{K_C})$  are  $L^\infty$  with respect to the Lebesgue measure.
- (iv) The following set is closed:

$$\left\{ \left( \int \phi_1 p \, d(w, v), \dots, \int \phi_{K_U} p \, d(w, v), \int \psi_1 p \, da'_k, \dots, \int \psi_{K_C} p \, da'_k, \int m p \, d(w, v) \right) \mid p \in \mathcal{M}_{\mathcal{W} \times \mathcal{V}}, p \geq 0 \right\}.$$

Assumption 7 (ii) is restrictive, but it is useful enough for showing Proposition 1 based on it. The rest are similar to Assumption 5. A sufficient condition for Assumption 7 (iv) is that the joint density of  $(W_i, V_i)$  in the data generating process, or its observationally equivalent one, is strictly positive on  $\mathcal{W} \times \mathcal{V}$  (Anderson, 1983, Theorem 9).

Under these assumptions, the following results characterize the identified set  $I$  of  $\theta$  and provide an equivalent condition for point-identification of  $\theta$ .



**Theorem 3.** Suppose Assumptions 6 and 7 hold. Let  $\lambda_k \in \mathbb{R}$  for  $k = 1, \dots, K_U$ , and let  $\mu_k : \mathcal{A}_k \mapsto \mathbb{R}$  for  $k = 1, \dots, K_C$ . Then  $I = [L, U]$  where

$$L = \max_{\{\lambda_k\}_{k=1}^{K_U}, \{\mu_k\}_{k=1}^{K_C}} \mathbb{E} \left[ \min_{v \in \mathcal{V}} \left\{ m(W_i, v) + \sum_{k=1}^{K_U} \lambda_k \phi_k(W_i, v) + \sum_{k=1}^{K_C} \mu_k(A_k(W_i, v)) \psi_k(W_i, v) \right\} \right] \quad (52)$$

and

$$U = \min_{\{\lambda_k\}_{k=1}^{K_U}, \{\mu_k\}_{k=1}^{K_C}} \mathbb{E} \left[ \max_{v \in \mathcal{V}} \left\{ m(W_i, v) + \sum_{k=1}^{K_U} \lambda_k \phi_k(W_i, v) + \sum_{k=1}^{K_C} \mu_k(A_k(W_i, v)) \psi_k(W_i, v) \right\} \right] \quad (53)$$

where  $A_k(w, v)$  is the value of  $A_{ik}$  given  $W_i = w$  and  $V_i = v$ .

*Proof.* The proof focuses on showing (52). The same argument applies to (53).

Let  $\mathcal{M}_{W \times V}$  to be the space of finite and countably additive signed Borel measures that are absolutely continuous with respect to the Lebesgue measure. Using absolute continuity, I identify an element of  $\mathcal{M}_{W \times V}$  by its density  $p : \mathcal{W} \times \mathcal{V} \mapsto \mathbb{R}$ . Let  $p_W$  be the density of  $P_W$ . The identified set  $I$  is then defined by

$$I \equiv \left\{ \int m(w, v) p(w, v) d(w, v) \mid p \in \mathcal{M}_{W \times V}, \quad p \geq 0, \right. \\ \int \phi_k(w, v) p(w, v) d(w, v) = 0, \quad k = 1, \dots, K_U, \\ \int \psi_k(a_k, a'_k) p(a_k, a'_k) da'_k = 0 \text{ for all } a_k \in \mathcal{A}_k, \quad k = 1, \dots, K_C, \\ \left. \int p(w, v) dv = p_W(w) \text{ for all } w \in \mathcal{W} \right\},$$

where  $a_k$  is an element of  $\mathcal{A}_k$  and  $a'_k$  is an element of  $\mathcal{A}'_k$ . The second line represents unconditional moment restrictions and the third line represents conditional moment restrictions.

The lower bound of  $I$  is then given by the infinite-dimensional linear program

$$\min_{p \in \mathcal{M}_{W \times V}, p \geq 0} \int m(w, v) p(w, v) d(w, v) \quad \text{subject to} \\ \int \phi_k(w, v) p(w, v) d(w, v) = 0, \quad k = 1, \dots, K_U, \\ \int \psi_k(a_k, a'_k) p(a_k, a'_k) da'_k = 0, \text{ for all } a_k \in \mathcal{A}_k, \quad k = 1, \dots, K_C, \\ \int p(w, v) dv = p_W(w) \text{ for all } w \in \mathcal{W}. \quad (54)$$

Now I show that (52) is the dual of (54). I introduce additional notation. Let  $L^2(\mathcal{W} \times \mathcal{V})$  be the space of all  $L^2$  functions on  $\mathcal{W} \times \mathcal{V}$ , and let  $L^2(\mathcal{W})$  be the space of all  $L^2$  functions on  $\mathcal{W}$ . We also let  $L^2(\mathcal{A}_k)$  be the space of all  $L^2$  functions on  $\mathcal{A}_k$ .

Let  $\mathcal{G}$  and  $\mathcal{H}$  be  $\mathcal{G} = \mathcal{H} = \mathbb{R}^K \times L^2(\mathcal{A}_1) \times \dots \times L^2(\mathcal{A}_{K_C}) \times L^2(\mathcal{W})$ . Denote their generic elements as  $g = (g_1, \dots, g_{K_U}, \bar{g}_1, \dots, \bar{g}_{K_C}, f_g)$  and  $h = (\lambda_1, \dots, \lambda_{K_U}, \mu_1, \dots, \mu_{K_C}, f_h)$ , respectively. Note that  $\mathcal{H}$  is a dual space of  $\mathcal{G}$ .

Define the linear map  $A : \mathcal{M}_{\mathcal{W} \times \mathcal{V}} \mapsto \mathcal{G}$ :

$$A(p) = \left( \int \phi_1 p d(w, v), \dots, \int \phi_K p d(w, v), \int \psi_k p da'_1, \dots, \int \psi_k p da'_{K_C}, \int p dv \right).$$

$A$  is a bounded (hence continuous) linear operator because  $\phi_k$ s and  $\psi_k$ s are assumed to be bounded. Define the dual pairing:

$$\langle A(P), h \rangle = \sum_{k=1}^{K_U} \lambda_k \int \phi_k p d(w, v) + \sum_{k=1}^{K_C} \iint \psi_k p da'_k \mu_k da_k + \int f_h \int p dv dw.$$

It is straightforward to show:

$$\iint \psi_k p da'_k \mu_k da_k = \int \psi_k \mu_k p d(w, v)$$

and

$$\int f_h \int p dv dw = \int f_h p d(w, v).$$

Then:

$$\langle A(P), h \rangle = \int \left[ \sum_{k=1}^{K_U} \lambda_k \phi_k + \sum_{k=1}^{K_C} \mu_k \psi_k + f_h \right] p(w, v) d(w, v). \equiv \langle p, A^*(h) \rangle, \quad (55)$$

where  $A^*(h) : \mathcal{H} \mapsto L^2(\mathcal{W} \times \mathcal{V})$  is defined as

$$A^*(h) = \sum_{k=1}^{K_U} \lambda_k \phi_k + \sum_{k=1}^{K_C} \mu_k \psi_k + f_h.$$

(55) shows that  $A^*$  is the adjoint of  $A$ .

I then apply the strong duality theorem under Assumption 7 and the continuity of  $A$  (Anderson, 1983, Theorem 6), which implies that the optimal solution to (54) equals to the

solution to:

$$\max_{\lambda_1, \dots, \lambda_{K_U}, \mu_1, \dots, \mu_{K_C}, f_h} \int f_h(w) p_w(w) dw \quad \text{subject to} \quad \sum_{k=1}^{K_U} \lambda_k \phi_k + \sum_{k=1}^{K_C} \mu_k \psi_k + f_h \leq m. \quad (56)$$

Simplifying (56) as in the proof of Theorem 2 yields the expression in (52).  $\square$

**Lemma 3.** *Suppose that assumptions of Theorem 3 hold. Suppose also that the joint density of  $(W_i, V_i)$  is strictly positive on  $\mathcal{W} \times \mathcal{V}$ . Then  $\theta$  is point-identified if and only if there exists a function  $S^*$  which is a linear functional on the projection of  $\mathcal{M}_{\mathcal{W} \times \mathcal{V}}$  onto  $\mathcal{W}$ , real numbers  $\lambda_1^*, \dots, \lambda_K^* \in \mathbb{R}$  and functions  $\mu_1^*, \dots, \mu_K^*$  which are  $L^2(\mathcal{A}_1), \dots, L^2(\mathcal{A}_{K_C})$  functions, respectively, such that:*

$$m(W_i, v) + \sum_{k=1}^{K_U} \lambda_k \phi_k(W_i, v) + \sum_{k=1}^{K_C} \mu_k(A_k(W_i, v)) \psi_k(W_i, v) = S^*(W_i)$$

almost surely on  $\mathcal{W} \times \mathcal{V}$ . When such  $S^*$  exists,  $\theta$  is identified by  $\theta = \mathbb{E}(S^*(W_i))$ .

*Proof.* As in (56) during the proof of Theorem 3, the sharp lower bound of  $\theta$  is given by

$$\max_{\lambda_1, \dots, \lambda_{K_U}, \mu_1, \dots, \mu_{K_C}, f_h} \int f_h(w) p_w(w) dw \quad \text{subject to} \quad \sum_{k=1}^{K_U} \lambda_k \phi_k + \sum_{k=1}^{K_C} \mu_k \psi_k + f_h \leq m.$$

where all notation follow the proof of Theorem 3. Similarly, the sharp upper bound of  $\theta$  is given by

$$\min_{\lambda_1, \dots, \lambda_{K_U}, \mu_1, \dots, \mu_{K_C}, f_h} \int f_h(w) p_w(w) dw \quad \text{subject to} \quad \sum_{k=1}^{K_U} \lambda_k \phi_k + \sum_{k=1}^{K_C} \mu_k \psi_k + f_h \geq m.$$

Lemma 3 can then be proved by replicating the proof of Lemma 2.  $\square$

### B.3 Estimation and inference under over-identification

In practice, the plug-in bound may yield an empty set, in which case  $\hat{L}$  diverges to  $+\infty$  and  $\hat{U}$  diverges to  $-\infty$ . This happens when the empirical distribution  $\hat{P}_W$  does not satisfy the moment restrictions, which may happen even if the population distribution  $P_W$  satisfies them. This can be compared to over-identification in the generalized method of moments (GMM) estimation, where the GMM objective may be strictly positive in the sample even

if the moments are correctly specified. In this case, the empirical version of (10) (where  $P_W$  is replaced with  $\hat{P}_W$ ) does not have a feasible solution, giving an empty plug-in bound.

There are two approaches for dealing with this issue. First, a researcher may obtain a point estimate that minimizes distance between the model and the data<sup>11</sup>. Second, a researcher may directly obtain a confidence interval without insisting on a point estimate, assuming that the model is correctly specified. For the first approach, consider the following relaxation of the moment restrictions:

$$|\mathbb{E}(\phi_k(W_i, V_i))| \leq \delta, \quad k = 1, \dots, K, \quad (57)$$

where  $\delta \geq 0$ , which reduces to Assumption 4 when  $\delta = 0$ . This can be thought of as an absolute-value GMM criterion. The following proposition explains how to compute the smallest  $\delta$  that makes (20) hold with the empirical distribution.

**Proposition 5.** *Given the sample  $(W_1, \dots, W_N)$ , consider linear programming problem:*

$$\min_{P \in \mathcal{M}_{W \times V}, P \geq 0, \delta \geq 0} \delta \quad \text{subject to} \quad \begin{cases} \left| \int \phi_k(W_i, V_i) dP \right| \leq \delta, & k = 1, \dots, K, \\ \int P(w, dV_i) = \hat{P}_W(w) & \text{for all } w \in \mathcal{W}, \end{cases} \quad (58)$$

where  $\hat{P}_W$  is the empirical distribution of  $W_i$  constructed from  $(W_1, \dots, W_N)$ . Its solution then equals the solution to:

$$\max_{\lambda_1, \dots, \lambda_K} \frac{1}{N} \sum_{i=1}^N \min_{v \in \mathcal{V}} \left\{ \sum_{k=1}^K \lambda_k \phi_k(W_i, v) \right\} \quad \text{subject to} \quad \sum_{k=1}^K |\lambda_k| \leq 1. \quad (59)$$

*Proof.* I can rewrite (58) as:

$$\min_{P \in \mathcal{M}_{W \times V}, P \geq 0, \delta \geq 0} \delta \quad \text{subject to} \quad \begin{cases} \int dP = 1, \\ \int \phi_k(W_i, V_i) dP \leq \delta, & k = 1, \dots, K, \\ \int \phi_k(W_i, V_i) dP \geq -\delta, & k = 1, \dots, K, \\ \int P(w, dV_i) = \hat{P}_W(w) & \text{for all } w \in \mathcal{W}. \end{cases}$$

I can then replicate the argument of Theorem 2, obtaining (59) as the simplified dual of the above.  $\square$

---

<sup>11</sup>Andrews and Kwon (2019) study this approach for standard GMM estimation of moment equality models without unobservables.

Proposition 5 shows that a researcher can find the smallest  $\delta$  by solving (59), which is similar to the plug-in bound problem. A difference is that (59) is a constrained optimization problem, but the constraint has a simple structure whose Jacobian can be derived in closed-form.

Let  $\delta^*$  be the solution to (59), and let  $\delta \geq \delta^*$ . I then compute the lower bound  $\tilde{L}$  under the relaxation in (57) by computing the plug-in lower bound with a negative  $L^1$  penalty on  $\lambda$ , with  $\delta$  being the penalty multiplier:

$$\tilde{L} = \max_{\lambda_1, \dots, \lambda_K} \left[ \frac{1}{N} \sum_{i=1}^N \min_{v \in \mathcal{V}} \left\{ m(W_i, v) + \sum_{k=1}^K \lambda_k \phi_k(W_i, v) \right\} - \delta \sum_{k=1}^K |\lambda_k| \right]. \quad (60)$$

The following proposition justifies use of the  $L^1$  penalty. I compute the upper bound  $\tilde{U}$  similarly with a positive  $L^1$  penalty.

**Proposition 6.** *Given the sample  $(W_1, \dots, W_N)$  and given  $\delta \in \mathbb{R}$ , consider the linear programming problem that finds the smallest value of  $\theta = \mathbb{E}(m(W_i, V_i))$  that satisfies (57):*

$$\begin{aligned} \min_{P \in \mathcal{M}_{W \times V}, P \geq 0} \int m(W_i, V_i) dP \quad \text{subject to} \quad & \left| \int \phi_k(W_i, V_i) dP \right| \leq \delta, \quad k = 1, \dots, K, \\ & \int P(w, dv) = \hat{P}_W(w) \text{ for all } w \in \mathcal{W}, \end{aligned} \quad (61)$$

where  $\hat{P}_W$  is the empirical distribution of  $W_i$  constructed from  $(W_1, \dots, W_N)$ . Its solution then equals  $\tilde{L}$ , defined in (60).

*Proof.* I can rewrite (61) as:

$$\begin{aligned} \min_{P \in \mathcal{M}_{W \times V}, P \geq 0} \int m(W_i, V_i) dP \quad \text{subject to} \quad & \int \phi_k(W_i, V_i) dP \leq \delta^*, \quad k = 1, \dots, K, \\ & \int \phi_k(W_i, V_i) dP \geq -\delta^*, \quad k = 1, \dots, K, \\ & \int P(w, dv) = \hat{P}_W(w) \text{ for all } w \in \mathcal{W}. \end{aligned}$$

I can then replicate the argument of Theorem 2, obtaining (60) as the simplified dual of the above.  $\square$

Proposition 6 shows that (60) equals the smallest value of  $\theta$  among the distributions whose absolute-value GMM criterion defined in (57) is at most  $\delta$ . In principle, such a distribution is not necessarily unique even when  $\delta = \delta^*$ . If it is unique, the modified plug-in bound becomes a point.

In practice, due to machine precision or the stopping criterion of the optimization algorithm, the numerical solution to (59) might be strictly smaller than the analytical solution  $\delta^*$ . To resolve this problem, a researcher may choose  $\delta$  to be larger than  $\delta^*$ , in which case (60) computes the smallest value of  $\theta$  among the distributions that attain the *near-minimum* of the absolute-value GMM criterion. In the case that the minimizer distribution is unique, the modified plug-in bounds with  $\delta > \delta^*$  becomes a small interval instead of a point.

Although (60) resolves the empty set problem, it has two drawbacks. First, it is an ad-hoc approach, having no formal justification for why relaxation of moment conditions is a good idea. Second, the procedure may yield a point identified set (or a small interval) even if the model is partially identified. The literature dealt with the second problem by choosing  $\delta$  that is reasonably larger than  $\delta^*$  (Mogstad, Santos, and Torgovitsky, 2018), but how much larger it should be remains a question. In the remainder of this subsection, I discuss a more principled approach, which is directly computing a confidence interval without insisting on a point estimate.

Note that the inference procedure that tests (18) does not involve the plug-in bound per se. The plug-in bound is involved only in the step of choosing  $\Lambda_F$ , which I propose to be the set of  $\lambda$ s that are close to the solutions to the plug-in bound problems. The inference procedure is valid regardless of whether the plug-in bound is empty or not; the issue is that there is no guidance for choosing  $\Lambda_F$  when the plug-in bound is empty. In what follows, I propose a strategy for choosing  $\Lambda_F$  when the plug-in bound is empty.

I propose to use the relaxed plug-in bounds for choosing  $\Lambda_F$ . The procedure consists of three steps. The first step solves (59) and finds minimum  $\delta^*$ . The second step considers a grid of positive real numbers  $\{\delta_1, \dots, \delta_M\}$  such that  $\delta_m \geq \delta^*$  for all  $m \in \{1, \dots, M\}$ . Then, for each  $\delta_m$ , compute the relaxed plug-in bounds:

$$\begin{aligned}\tilde{\lambda}_L(\delta_m) &= \operatorname{argmax}_{\lambda_1, \dots, \lambda_K} \left[ \frac{1}{N} \sum_{i=1}^N \min_{v \in \mathcal{V}} \left\{ m(W_i, v) + \sum_{k=1}^K \lambda_k \phi_k(W_i, v) \right\} - \delta_m \sum_{k=1}^K |\lambda_k| \right]. \\ \tilde{\lambda}_U(\delta_m) &= \operatorname{argmin}_{\lambda_1, \dots, \lambda_K} \left[ \frac{1}{N} \sum_{i=1}^N \max_{v \in \mathcal{V}} \left\{ m(W_i, v) + \sum_{k=1}^K \lambda_k \phi_k(W_i, v) \right\} + \delta_m \sum_{k=1}^K |\lambda_k| \right].\end{aligned}\tag{62}$$

The third step then chooses points in the neighborhoods of *every*  $\tilde{\lambda}_L(\delta_m)$  and  $\tilde{\lambda}_U(\delta_m)$ . In the simulation and application, I choose points by adding Gaussian noise to every  $\tilde{\lambda}_L(\delta_m)$  and  $\tilde{\lambda}_U(\delta_m)$  whose standard deviations are inversely proportional to the gradients of the  $\tilde{\lambda}$ s. I check performance of this approach by simulation in Section 7.

When  $\delta^* = 0$ , i.e. when the plug-in bound is non-empty, a researcher may choose

$M = 1$  with  $\delta_1 = 0$ , in which case the procedure reduces to the previous procedure in Section 5.2. This means that the inference procedure with relaxed bounds generalizes the procedure discussed in Section 5.2.

## B.4 Construction of the dataset

I use data on U.S. households from the Panel Study of Income Dynamics (PSID) dataset. I use the dataset of Guvenen (2009), who estimated RIP and HIP processes using PSID earnings data of male head of households collected annually from 1968 to 1993. The dataset contains male head of households who are not in the poverty (SEO) subsample and who consecutively reported positive hours (between 520 and 5110 hours a year) and earnings (between a preset minimum and maximum wage). I also follow Guvenen (2009) and use potential experience as a measure of experience:

$$h = \text{age} - \max\{\text{years of schooling}, 12\} - 6,$$

Note that potential experience is a strictly exogenous variable since it is a deterministic function of age. I collect individuals with consecutive waves of data from 1976 to 1991, which yields  $N = 800$  and  $T = 15$  taking the first wave as an initial value of earnings.

Recall that my method requires that there is no multicollinearity in each individual time series (Assumption 2). To maintain this assumption, I removed 40 individuals (that is, five percent of data) who have smallest variations in their reported incomes, giving a dataset of  $N = 760$  and  $T = 15$ . Estimation results with this removal does not qualitatively differ from the results with full data. Tables 6 and 7 present confidence intervals for  $\mathbb{E}(\rho_i)$ ,  $\text{Var}(\rho_i)$  and the CDF of  $\rho_i$  without removing observations. The confidence intervals do not qualitatively differ to Tables 4 and 5 in the main text, in the sense that upper confidence limit of  $\mathbb{E}(\rho_i)$  is significantly less than 1, lower confidence limit of  $\text{Var}(\rho_i)$  is strictly positive for the RIP process and the CDF of  $\rho_i$  has confidence limits away from 0 and 1.

## B.5 Simulation based de-noising method

Before describing how I separated  $\varepsilon_{it}$  from  $Y_{it}$ , I first describe the simulation based de-noising method proposed in Arellano and Bonhomme (2018). They considered a model

$$Z = X + \varepsilon$$

where all variables are scalar<sup>12</sup> and  $Z$  is observed but  $X$  and  $\varepsilon$  are not. Instead, distribution of  $\varepsilon$  is known. A researcher's objective is to obtain the distribution of  $X$ , or the data drawn from it.

The idea of Arellano and Bonhomme (2018) is then to find a distribution of  $X$  that minimizes the second order Wasserstein distance ( $W_2$ ) between  $Z$  and  $X + \varepsilon$ . It is known that  $W_2$  metrizes convergence in distribution plus the convergence of the second moment among distributions with compact support, which means that convergence of  $W_2$  to zero is equivalent to convergence in the distribution plus the convergence of the second moment. They proposed a procedure for obtaining i.i.d. draws of  $X$  given the observed data of  $Z$  and the simulated i.i.d. draws of  $\varepsilon$ .

I use their idea to separate transitory income from the observed income data, obtaining simulated draws of permanent income. I assume that the transitory income process in the application,  $\varepsilon_{it}$ , follows i.i.d. zero-mean Normal distribution whose variance equals to the variance estimate of the transitory income in Guvenen (2009). I then simulate  $K = 200$  i.i.d. draws of  $\varepsilon_k = (\varepsilon_{k1}, \dots, \varepsilon_{kT})'$ . Then, given the permanent income data  $\tilde{Y}_j = (\tilde{Y}_{j1}, \dots, \tilde{Y}_{jT})'$ ,  $j = 1, \dots, N$ , I obtain the *simulated* income data  $\hat{Y}_{jk}$  defined by:

$$\hat{Y}_{jk} = \tilde{Y}_j + \varepsilon_k,$$

giving the simulated income data of size  $NK$ .

I compare this to the observed income data  $Y_i = (Y_{i1}, \dots, Y_{iT})'$ ,  $i = 1, \dots, N$ . I compute the distance between the simulated and the observed data using the second order Wasserstein distance:

$$W_2^2(\{Y_i\}, \{\hat{Y}_{ij}\}) = \min_{0 \leq p_{ijk} \leq 1} \sum_{i=1}^N \sum_{j=1}^N \sum_{k=1}^K p_{ijk} \|Y_i - \hat{Y}_{jk}\|^2$$

subject to  $\sum_{i=1}^N p_{ijk} = 1, \quad \sum_{j=1}^N \sum_{k=1}^K p_{ijk} = 1.$

I then obtain the simulated permanent income  $\tilde{Y}_j = (\tilde{Y}_{j1}, \dots, \tilde{Y}_{jT})'$  by solving:

$$\underset{\tilde{Y}_j}{\operatorname{argmin}} W_2^2(\{Y_i\}, \{\hat{Y}_{ij}\}),$$

giving a dataset with  $N = 800$  individuals and  $T = 15$  waves.

Estimation results in Section 8.4 of the main text are not qualitatively affected by this

---

<sup>12</sup>They also consider a more general case of multivariate factor models.



de-noising procedure. Tables 8 and 9 present estimation results without the de-noising step, which are qualitatively similar to estimation results with de-noising. In particular, upper confidence limit of  $\mathbb{E}(\rho_i)$  is significantly less than 1, lower confidence limit of  $\text{Var}(\rho_i)$  is strictly positive for the RIP process and the CDF of  $\rho_i$  has confidence limits away from 0 and 1.

	Confidence interval of $\mathbb{E}(\rho_i)$	Confidence interval of $\text{Var}(\rho_i)$
RIP process	[0.415, 0.652]	[0.073, 0.235]
HIP process	[0.262, 0.692]	[0.000, 0.659]

Table 6: Confidence interval for  $\mathbb{E}(\rho_i)$  and  $\text{Var}(\rho_i)$ , for the dataset of  $N = 800$  that does not remove individuals with small variation in their reported incomes and with additional moment restrictions. Nominal coverage probability is 0.9.

$\mathbb{P}(\rho_i \leq r)$	RIP process	HIP process
$r = 0.0$	[0.000, 0.422]	[0.000, 0.665]
$r = 0.1$	[0.001, 0.482]	[0.014, 0.766]
$r = 0.2$	[0.024, 0.613]	[0.075, 0.804]
$r = 0.3$	[0.054, 0.674]	[0.104, 0.827]
$r = 0.4$	[0.090, 0.770]	[0.204, 0.872]
$r = 0.5$	[0.125, 0.845]	[0.217, 0.935]
$r = 0.6$	[0.188, 0.930]	[0.208, 0.979]
$r = 0.7$	[0.256, 0.959]	[0.250, 1.000]
$r = 0.8$	[0.330, 0.987]	[0.310, 1.000]
$r = 0.9$	[0.411, 0.998]	[0.369, 1.000]
$r = 1.0$	[0.498, 1.000]	[0.428, 1.000]

Table 7: Confidence intervals for  $\mathbb{P}(\rho_i \leq r)$ , for the dataset of  $N = 800$  that does not remove individuals with small variation in their reported incomes. Nominal coverage probability is 0.9.

	Confidence interval of $\mathbb{E}(\rho_i)$	Confidence interval of $\text{Var}(\rho_i)$
RIP process	[0.451, 0.615]	[0.050, 0.293]
HIP process	[0.242, 0.633]	[0.000, 0.700]

Table 8: Confidence interval for  $\mathbb{E}(\rho_i)$  and  $\text{Var}(\rho_i)$ , for the dataset without the de-noising step of size  $N = 760$  with removing individuals with small variation in their reported incomes. Nominal coverage probability is 0.9.

$\mathbb{P}(\rho_i \leq r)$	RIP process	HIP process
$r = 0.0$	[0.000, 0.364]	[0.000, 0.715]
$r = 0.1$	[0.012, 0.413]	[0.006, 0.765]
$r = 0.2$	[0.026, 0.510]	[0.037, 0.845]
$r = 0.3$	[0.084, 0.584]	[0.122, 0.843]
$r = 0.4$	[0.118, 0.725]	[0.152, 0.882]
$r = 0.5$	[0.150, 0.826]	[0.170, 0.936]
$r = 0.6$	[0.232, 0.879]	[0.216, 0.983]
$r = 0.7$	[0.309, 0.934]	[0.283, 1.000]
$r = 0.8$	[0.393, 0.982]	[0.324, 1.000]
$r = 0.9$	[0.493, 0.990]	[0.359, 1.000]
$r = 1.0$	[0.558, 1.000]	[0.384, 1.000]

Table 9: Confidence intervals for  $\mathbb{P}(\rho_i \leq r)$ , for the dataset without the de-noising step of size  $N = 760$  with removing individuals with small variation in their reported incomes . Nominal coverage probability is 0.9.