

# Heterogeneity, Uncertainty and Learning: Semiparametric Identification and Estimation\*

Jackson Bunting<sup>†</sup>      Paul Diegert<sup>‡</sup>      Arnaud Maurel<sup>§</sup>

December 29, 2023

## Abstract

We provide semiparametric identification results for a broad class of learning models in which outcomes of interest depend on i) predictable heterogeneity, ii) initially unpredictable heterogeneity that may be revealed over time, and iii) transitory uncertainty. We consider a common environment where the researcher only has access to a short panel on choices and realized outcomes. We establish point-identification of the outcome equation parameters and the distribution of the three types of unobservables, under the standard assumption that unpredictable heterogeneity and uncertainty are normally distributed. We also show that, in the absence of predictable heterogeneity, the model is identified without making any distributional assumption. We then derive the asymptotic properties of a sieve MLE estimator for the model parameters, and devise a tractable profile likelihood based estimation procedure. Monte Carlo simulation results indicate that our estimator exhibits good finite-sample properties.

---

\*We thank seminar participants at CREST-PSE, LMU Munich, TSE, UC Davis, UT-Austin, conference participants at the 34<sup>th</sup> EC<sup>2</sup> conference, the 2023 IAAE, SETA and SOLE meetings, and Victor Aguirregabiria, Peter Arcidiacono, Stéphane Bonhomme, Xavier D'Haultfoeulle, Yuichi Kitamura, Mauricio Olivares, Yuya Sasaki, Chris Taber and Daniel Wilhelm for useful comments. We thank Zhangchi Ma for capable research assistance.

<sup>†</sup>Texas A&M University, [jbunting@tamu.edu](mailto:jbunting@tamu.edu).

<sup>‡</sup>Toulouse School of Economics, [paul.diegert@tse-fr.eu](mailto:paul.diegert@tse-fr.eu).

<sup>§</sup>Duke University, NBER and IZA, [arnaud.maurel@duke.edu](mailto:arnaud.maurel@duke.edu).

# 1 Introduction

Learning models, in which agents have imperfect information about their environment and update their beliefs over time, are frequently used in economics. These models have received particular interest in various subfields in empirical microeconomics, including industrial organization and health (see, e.g., Akerberg, 2003; Coscelli and Shum, 2004; Crawford and Shum, 2005; Abbring and Campbell, 2005; Chan and Hamilton, 2006; Yang, 2020; Aguirregabiria and Jeon, 2020, for a survey in the context of oligopoly competition), as well as in labor economics (see, e.g. Miller, 1984; Antonovics and Golan, 2012; Pastorino, 2015; Hincapié, 2020; Pastorino, 2022) and economics of education (see, e.g. Arcidiacono, 2004; Zafar, 2011; Stinebrickner and Stinebrickner, 2012; Stange, 2012; Thomas, 2019; Kinsler and Pavan, 2021; Arcidiacono et al., 2023). Since the seminal work of Erdem and Keane (1996), learning models have also been popular in the marketing literature (see Ching et al., 2013, for a survey). However, while learning models are often structurally estimated, much remains to be known about the identification of this important class of models.

In this paper we provide new semiparametric identification results for a general class of learning models. Importantly, we consider an environment where the researcher has access to a short panel on choices and realized outcomes only. As such, our results are widely applicable, including in frequent situations where one does not have access to elicited beliefs data, or to a vector of selection-free measurements of unobserved individual heterogeneity. Specifically, we consider throughout our analysis a potential outcome model where unit  $i$ 's potential outcome in period  $t$  from assignment  $D_{it} = d$  is

$$Y_{it}(d) = X_{it}^\top \beta_t(d) + (X_i^*)^\top \lambda_t(d) + \epsilon_{it}(d), \quad (1)$$

where  $X_{it}$  is a vector of explanatory variables associated with individual  $i$  in period  $t$ ,  $X_i^*$  denotes a vector of latent individual effects,  $\epsilon_{it}(d)$  is an idiosyncratic shock, and

$(\beta_t(d), \lambda_t(d))$  are unknown parameters. While interactive fixed effects models of this kind have been the object of much interest in econometrics, a key distinctive feature of the setup considered in this paper is the existence of two different types of individual effects. Namely, we assume the individual effect  $X_i^*$  consists of two components:  $X_{k,i}^*$  which are supposed to be initially known by the agent, and  $X_{u,i}^*$  which are initially unknown but may be learned over time. We complement this outcome model with a flexible choice model, in which agent  $i$ 's assignment in period  $t$  can depend arbitrarily on contemporaneous and lagged explanatory variables, assignments and realized outcomes. This framework encompasses most of the decision models that have been considered in the learning literature.

We first show that the model is point-identified under two alternative sets of conditions. Our first and main identification result applies to a setup where, consistent with most of the Bayesian learning models that have been considered in the literature, we assume that the idiosyncratic shocks from the outcome equations  $(\epsilon_{it}(d))$ , as well as the unknown heterogeneity component  $(X_{u,i}^*)$ , are normally distributed. In contrast, the distribution of the known heterogeneity component  $(X_{k,i}^*)$  is left unspecified. From the key observation that the distribution of current realized outcomes conditional on past choices and outcomes is a mixture of normal distributions, we leverage results from Bruni and Koch (1985) to establish identification of the joint distribution of realized outcomes, choices and known heterogeneity component  $X_{k,i}^*$ .

We then also show that a pure learning model with only one type of permanent unobserved heterogeneity  $(X_{u,i}^*)$  actually remains point-identified without making any distributional assumption. A crucial distinction from the general case is that this model is one of selection on observables only, as individual choices depend on beliefs about  $X_{u,i}^*$  only through prior outcomes, choices and covariates. This feature allows us to build on insights from the interactive fixed effects literature, in particular Freyberger (2018), to establish identification in the pure learning case.

We propose to estimate the model parameters  $\theta$  via sieve maximum likelihood

estimation. We focus on a particular class of functionals of  $\theta$ , which includes as special cases economically relevant quantities, such as the predictable and unpredictable outcome variances. These variances can in turn be used to evaluate the relative importance of, e.g., uncertainty vs. heterogeneity in the overall lifecycle earnings variability - a question that has been the object of much interest in labor economics (see, e.g., Cunha et al., 2005; Huggett et al., 2011; Cunha and Heckman, 2016). We show that, under mild regularity conditions, the resulting estimators are consistent and asymptotically normal. Monte Carlo simulation results indicate that our estimator exhibits good finite-sample properties. Importantly for practical purposes, our proposed estimator only involves a modest computational cost.

## Related literatures

Our paper contributes to several strands of the literature. First and foremost, we add to a set of papers that study the identification of learning models, generally in the context of specific applications (Abbring and Campbell, 2005; Arcidiacono et al., 2023; Gong, 2019; Pastorino, 2022). A key difference with most of the papers in this literature is that we only impose mild restrictions on the choice process. In particular, we remain agnostic about how choices depend on individual beliefs about  $X_{u,i}^*$ , while allowing these beliefs to depend arbitrarily on past choices and realized outcomes. Particularly relevant for us is recent work by Pastorino (2022), which establishes formal identification results for a econometric learning model. However, beyond the fact that Pastorino (2022) focus on the context of workers' and firms' learning, there are two main differences relative to our paper. First, unobserved heterogeneity in that paper is restricted to be discrete, while we allow for both continuous and multivariate unobserved heterogeneity. Second, and importantly, the known component of unobserved heterogeneity is excluded from the observed outcomes which form the basis of learning. In contrast, in this paper, outcomes may depend on both known and unknown components.

Our paper also fits into a literature that focuses on the identification of Markovian dynamic discrete choice models in the presence of persistent unobserved heterogeneity (see, among others, Heckman and Navarro, 2007; Hu and Schennach, 2008; Kasahara and Shimotsu, 2009; Hu and Shum, 2012; Sasaki, 2015; Hu and Sasaki, 2018; Aguirregabiria et al., 2021; Bunting, 2022). Unlike these papers, our learning model lacks a Markov structure, since beliefs and decisions are allowed to depend on the entire history of past outcomes and choices.<sup>1</sup> More broadly, our analysis is related to the literature that deals with the identification of mixture models (see, for example, Compiani and Kitamura, 2016; Kitamura and Laage, 2018, and references therein). In particular, central to our main identification result is the observation that the distribution of current outcomes conditional on the sequence of past choices and outcomes is a mixture of normal distributions.

Finally, since the outcome equation in our model involves interactions between unobserved individual- and time-specific effects, our paper also fits into the literature that deals with the identification and estimation of panel data models with interactive fixed effects (see, e.g., Bai, 2009; Gobillon and Magnac, 2016; Freyberger, 2018). Among these papers, our identification strategy is most closely related to Freyberger (2018). A fundamental distinction though comes from the fact that Freyberger (2018) considers a selection-free environment. In contrast, individual choices, along with the associated selection issues affecting the potential outcomes, play a central role in our analysis.

---

<sup>1</sup>Although our framework is more general, Bayesian learning models often naturally possess a first order Markov structure. There are, however, a number of additional significant differences between our paper and the listed literature. For example, Hu and Shum (2012) focus on scalar unobserved heterogeneity, whereas multivariate unobserved heterogeneity is fundamental to our main setting (i.e., since both  $X_i^*$  and beliefs are latent). Beyond this, several of their assumptions may fail to hold in our setting. For instance, since the support of the latent beliefs is larger than the support of the choices, the requirement that the observed variables be noisy “measurements” of the latent variables (Hu and Shum, 2012, Assumption 2) may fail.

## Organization of the paper

The remainder of the paper is organized as follows. Section 2 presents the set-up of the model. Section 3 contains our main identification results, both for the general case and for the case of a pure learning model. We discuss in Section 4 the estimation and inference on the parameters of interest, before turning in Section 5 to the implementation of our estimator. We study in Section 6 its finite-sample performances. Finally, Section 7 concludes. The appendix gathers all the proofs.

*Notation:*  $\mathcal{S}(A)$  indicates the support of random variable  $A$ .  $F_A$  indicates the cumulative distribution function of random variable  $A$ , whereas  $f_A$  indicates the probability mass or density function. For any sequence  $(a_1, a_2, \dots, a_S)$  and  $s \leq S$ , we let  $a^s = (a_1, a_2, \dots, a_s)$ . Upper case letters represent random variables, lower case represent realized values.  $A \perp\!\!\!\perp B \mid C$  indicates that  $A$  and  $B$  are statistically independent conditional upon  $C$ . Finally, to save on notations, we suppress the individual subscript  $i$  from all random variables in the remainder of the article.

## 2 Set-up

Throughout the paper we consider a setup where potential outcomes have an interactive fixed effect structure of the following form:

$$Y_t(d) = X_t^\top \beta_t(d) + X_k^* \lambda_{kt}(d) + (X_u^*)^\top \lambda_{ut}(d) + \epsilon_t(d), \quad (2)$$

where  $d$  represents a possible value of individual  $i$ 's assignment in period  $t$ ,  $Y_t(d)$  is a scalar potential outcome variable associated with assignment  $d$ ,  $X_t$  is a vector of observed explanatory variables,  $X^* = (X_k^*, (X_u^*)^\top)^\top$  are unobserved explanatory variables, and  $(\beta_t(d)^\top, \lambda_t(d)^\top)^\top$  with  $\lambda_t(d) := (\lambda_{kt}(d), \lambda_{ut}(d)^\top)^\top$  is a vector of unknown parameters, and  $\epsilon_t(d)$  is an unobserved random variable. For example,  $Y_t(d)$  may represent potential log-wages in occupation  $d$ .  $Y_t(d)$  may depend on some observed

individual and possibly time-varying characteristics ( $X_t$ ) as well as on multiple dimensions of unobserved abilities ( $X^*$ ), which might play different roles in different occupations (see, e.g., Arcidiacono et al., 2023; Hincapié, 2020). This setup is fairly general and can be applied in a wide range of contexts. For instance,  $Y_t(d)$  may also represent the potential log-quantity of a particular product sold by a firm in a given market  $d$  (see, e.g., Berman et al., 2019). This framework can also be used in the health context, where,  $Y_t(d)$  may correspond to a health outcome measure associated with a certain drug (e.g. CD4 cell counts associated with HIV drug treatment  $d$ , as in Chan and Hamilton, 2006), or the body mass index associated with a certain type of diet  $d$ .

Importantly, we allow for two distinct types of latent individual effects. Namely,  $X_k^*$  is assumed to be known by the agent, while  $X_u^*$  is initially unknown but may be gradually revealed over time. For example, worker  $i$ 's log-wage in occupation  $d$  at time  $t$ ,  $Y_t(d)$ , may depend on her unobserved (to the econometrician) occupation specific productivity,  $X_k^* \lambda_{kt}(d) + (X_u^*)^\top \lambda_{ut}(d)$ . As the worker accumulates more experience, she may update her belief about  $X_u^*$ , and thus about the initially unknown portion of productivity in each of the possible occupations.

Turning to the choice and learning process, the only restriction placed on an individual's assignment in period  $t$ , which we denote by  $D_t$ , is that it does not directly depend on the unknown component of latent heterogeneity. Specifically, we impose the following restriction:

$$D_t \perp\!\!\!\perp X_u^* \mid X^t, Y^{t-1}, D^{t-1}, X_k^*. \quad (3)$$

The above conditional independence condition highlights the asymmetry between the two types of latent effects: assignments may directly depend on the known component of the latent effect  $X_k^*$ , but not on the unknown component of the latent effect  $X_u^*$ . However, we allow the assignment rule to depend arbitrarily on lagged covariates, outcomes and choices. As a result, we do not restrict how agents form their beliefs

about  $X_u^*$ , provided that such beliefs are a measurable function of  $X^t, Y^{t-1}, D^{t-1}$  and  $X_k^*$ . We also remain agnostic about how assignments depend on agents' beliefs over  $X_u^*$ .

In particular, this framework is consistent with a setup where agents are rational and Bayesian updaters, so that beliefs coincide with the objective distribution of  $X_u^*$  conditional on their information set at a given point in time, which may include all realized variables and model parameters. Alternatively, this also accommodates situations where individual decisions may not involve beliefs over the distribution of  $X_u^*$ , or depend instead on myopic beliefs that are formed based on the prior-period choice and outcome. This further allows for heterogeneous beliefs formation, where, for instance, some agents may have rational expectations about their unobserved characteristic  $X_u^*$ , while others may have biased beliefs.

Finally, we denote the conditional choice probability (CCP) function as

$$h_t(d^t, x^t, y^{t-1}, v_k) := \Pr(D_t = d \mid X^t = x^t, Y^{t-1} = y^{t-1}, D^{t-1} = d^{t-1}, X_k^* = v_k).$$

These CCPs play a central role in our identification analysis. In the following section, we provide sufficient conditions under which the CCPs - which are latent objects because of the conditioning on  $X_k^* = v_k$  - are identified. In empirical applications it is very common to impose some structure on the choice process. For example, in a dynamic discrete choice framework it is standard to assume (e.g., Arcidiacono et al., 2023) that

$$D_t = \arg \max_{d \in S(D_t)} \{v_t(d, X_t, X_k^*, S_t) + \eta_t(d)\},$$

where  $v_t$  is known up to a finite-dimensional vector of parameters,  $S_t$  are sufficient statistics for the conditional distribution of  $X_u^*$  at time  $t$ , and  $\eta_t$  follows a known distribution. Having identified the CCPs, one can then apply standard identification arguments from the dynamic discrete choice literature (Magnac and Thesmar, 2002) to identify the primitives of the choice model.



## 2.1 Uncertainty and Learning

The key feature of the model is the distinction between the three forms of unobserved heterogeneity: (1) permanent heterogeneity that is known to the agent,  $X_k^*$ , (2) permanent heterogeneity that is initially unknown to the agent,  $X_u^*$ , and (3) idiosyncratic time-varying shocks,  $\epsilon$ . This provides a framework for quantifying the importance of uncertainty in outcomes. At  $t = 1$ , the variance in future outcomes can be decomposed orthogonally into a component that depends on  $(X_u^*, \epsilon)$  and a component that depends on  $X_k^*$ . Cunha et al. (2005); Cunha and Heckman (2016) consider this decomposition in the context of educational choice, decomposing the variance in lifetime earnings into a component that is predictable when deciding to go to college and a component that is not.

In our framework, the importance of uncertainty can change over time as agents learn about  $X_u^*$  by observing realized outcomes and covariates and use this information to make choices. We provide in Appendix C a class of variance decomposition parameters which includes both the  $t = 1$  orthogonal decomposition as well as  $t > 1$  decompositions that incorporate these learning and selection effects. These decompositions each provide different ways of quantifying the importance of uncertainty to future outcomes. Identification of the model implies the identification of these parameters. After establishing identification of the model, we pay special attention to estimation and inference of a broad class of functionals that encompasses these kinds of variance decompositions.

## 3 Identification

We provide in Subsection 3.1 a high-level overview of the proposed identification strategies. We then discuss identification in the leading case with both known and unknown unobserved heterogeneity (Subsection 3.2), before turning to the pure learning case where the only source of permanent unobserved heterogeneity is assumed to

be initially unknown to the agent (Subsection 3.3).

### 3.1 Overview

At its essence, the main identification problem analyzed in this paper is how to recover the distributions of potential outcomes (i.e.,  $f_{Y_t(d_t)|X^*X_t}$  for each  $t$  and  $d_t$ ) and selection probabilities (i.e.,  $f_{D_t|X_k^*X_t}$  for each  $t$ ) from the censored population data (i.e.,  $f_{Y^T D^T X^T}$ ). We discuss below how the learning structure is useful to solve this selection problem.

To illustrate the problem, consider a simplified version of our model with a binary choice in each period (i.e.,  $\mathcal{S}(D_t) = \{0, 1\}$ ) and without covariates. Let  $D := \prod_{t=1}^T D_t$ , and  $Y(1) := (Y_1(1), \dots, Y_T(1))$ , and focus on identification of the distribution of  $Y(1)$ , which is censored for  $D = 0$ . By Bayes' rule, the relationship between the target and censored distributions can be characterized as follows:

$$f_{Y|D}(y|1) \frac{f_D(1)}{f_{D|Y(1)}(1|y)} = f_{Y(1)}(y)$$

where the conditional density  $f_{Y|D}(y|1)$ , which is directly identified from the data, is weighted by the term  $\frac{f_D(1)}{f_{D|Y(1)}(1|y)}$ , which reflects selection.

Our learning framework provides one strategy for identifying these selection weights. To explain, first assume that all components of the latent effect are initially unknown (i.e.,  $X^* = X_u^*$ ). In a learning model where the decision makers' actions depend on beliefs on  $X^*$ , it is often natural to assume that beliefs depend only on past realized outcomes and choices, and that:

$$f_{D_t|Y(1)D^{t-1}}(1|y, 1) = f_{D_t|Y^{t-1}(1)D^{t-1}}(1|y^{t-1}, 1). \quad (4)$$

Intuitively, equation (4) states that once beliefs are controlled for, contemporaneous choices and outcomes are independent. The right hand side is identified from the joint distribution of  $(D^t, Y^{t-1})$  conditional on  $D^{t-1} = 1$ . Applying this idea recursively, it

follows that the component of the selection weight  $f_{D|Y(1)}(1|y)$  is identified as follows:

$$f_{D|Y(1)}(1|y) = f_{D_T|Y^{T-1}(1)D^{T-1}}(1|y^{T-1}, 1)f_{D_{T-1}|Y^{T-2}(1)D^{T-2}}(1|y^{T-2}, 1) \cdots f_{D_1}(1)$$

We pursue this identification approach in Section 3.3 in a version of the model we call *pure learning*. The exclusion restriction in (4) will generally break down, however, when agents also possess persistent private information that affect their decision (i.e.,  $X_k^*$ ). We propose in Section 3.2 an identification strategy that can be used in such situations. In particular, we show that maintaining a normality assumption commonly made in the learning literature is, in fact, sufficient to identify the joint distribution of  $(Y^T, D^T, X_k^*)$  (and any potential covariates  $X^T$ ) in a first step. One can then identify the model parameters in a second step, along the lines of the reweighting strategy discussed above.

### 3.2 Known and unknown heterogeneity

This section provides sufficient conditions for identification of the model discussed in Section 2. The first assumption (KL1) imposes that any correlation in the observed outcomes and choices over time and across assignments is due to the known latent effect  $X_k^*$ . It also imposes that the transition of the control variables,  $X_t$ , does not depend on unobservables.

**Assumption KL1.** Equation (2) holds,  $X_t$  includes a constant term, and  $\forall d \in \mathcal{S}(D_t)$

$$dF_{\epsilon_t(d)D_tX_t|Y^{t-1}D^{t-1}X^{t-1}X^*} = dF_{\epsilon_t(d)}dF_{D_t|Y^{t-1}D^{t-1}X^tX_k^*}dF_{X_t|Y^{t-1}D^{t-1}X^{t-1}}.$$

Assumption KL2 imposes that the unknown component of the individual effect,  $X_u^*$ , is drawn from a multivariate normal distribution, and that the random shock in the outcome equation is normally distributed.

**Assumption KL2.**  $X_u^* \mid (X_1 = x_1, X_k^* = v_k) \sim N(0, \Sigma_u(x_1))$  and  $\epsilon_t(d) \sim N(0, \sigma_t(d)^2)$ .

Assumption **KL2** leads to a specific functional form for the posterior distribution, namely the Gaussian conjugate distribution. We summarize this result in Lemma **1**. To do so, we define  $(E_t, \Sigma_t)$  recursively as follows. First,  $(E_1, \Sigma_1) = (0, \Sigma_u(X_1))$ . Second,

$$\begin{aligned}\Sigma_{t+1} &= (\Sigma_t^{-1} + \lambda_{ut}(D_t)\lambda_{ut}(D_t)^\top \sigma_t^{-2}(D_t))^{-1} \\ E_{t+1} &= \Sigma_{t+1} \left( \Sigma_t^{-1} E_t + \lambda_{ut}(D_t) \frac{Y_t - X_t^\top \beta_t(D_t) - X_k^* \lambda_{kt}(D_t)}{\sigma_t^2(D_t)} \right).\end{aligned}$$

**Lemma 1.** Let Assumptions **KL1** and **KL2** hold. Then  $X_u^*$  conditional upon  $(D^{t-1}, Y^{t-1}, X^t, X_k^*)$  is distributed  $N(E_t, \Sigma_t)$ .

Since  $X_u^*$  conditional on  $(Y^{t-1}, D^{t-1}, X^t, X_k^*)$  follows a normal distribution with mean  $E_t$  and variance-covariance matrix  $\Sigma_t$ ,  $(E_t, \Sigma_t)$  is a sufficient statistic for  $X_u^*$  at time  $t$ . Notice that  $(E_t, \Sigma_t)$  is a deterministic function of  $(D^{t-1}, Y^{t-1}, X^t, X_k^*)$  and  $\theta_1 = ((\beta_t, \lambda_{kt}, \lambda_{ut}, \sigma_t)_{t=1}^T, \Sigma_u(X_1)) \in \Theta_1$ . Furthermore, we can express  $(E_t, \Sigma_t)$  non-recursively<sup>2</sup> as:

$$\begin{aligned}\Sigma_{t+1} &= \left( \Sigma_u^{-1}(X_1) + \sum_{s=1}^t \lambda_{us}(D_s)\lambda_{us}(D_s)^\top \sigma_s^{-2}(D_s) \right)^{-1} \\ E_{t+1} &= \Sigma_{t+1} \left( \sum_{s=1}^t \lambda_{us} \frac{Y_s - X_s^\top \beta_s(D_s) - X_k^* \lambda_{ks}(D_s)}{\sigma_s^2(D_s)} \right)\end{aligned}$$

Suppose  $X_u^* \in \mathbb{R}^p$ . Our three remaining assumptions are as follows.

**Assumption KL3.** (A) For some  $d_1$ , the element of  $\beta_1(d_1)$  associated with

---

<sup>2</sup>Our identification result would go through if one replaces the first part of Assumption **KL2** with  $X_u^* \mid (X_1 = x_1, X_k^* = v_k) \sim N(0, \Sigma_u(v_k, x_1))$  under some regularity conditions on  $v_k \mapsto \Sigma_u(v_k, x_1)$ , including for each  $v_k - \tilde{v}_k > 0$ ,  $\Sigma_u(x_1, v_k) - \Sigma_u(x_1, \tilde{v}_k)$  is positive (or negative) semi-definite. For simplicity, we maintain the stronger Assumption **KL2** when establishing identification in Theorem 1 below.

the constant term is zero, and  $\lambda_{k1}(d_1) = 1$ . (B) For some  $(d_1, d_2, \dots, d_p)$ ,  $(\lambda_{u1}(d_1)\lambda_{u2}(d_2)\dots F_{up}(d_p)) = I_{p \times p}$ .

Assumption **KL3** is a normalization on the finite dimensional parameters. This type of assumption is standard in interactive fixed effect models (Freyberger, 2018), since no scale assumption is placed on the distribution of the unknown latent effects. For example, one could replace Assumption **KL3** (A) by a zero mean restriction on the latent individual effect  $X^*$ , and a unitary variance assumption on the known component of the latent effect  $X_k^*$ . We also impose Assumption **KL3** (B) since the unknown latent effect is inherently scale free.

**Assumption KL4.** (A) For each  $x_1 \in \mathcal{S}(X_1)$ ,  $\Theta_1$  is a compact set. (B)  $\mathcal{S}(X_k^*)$  is a compact set. (C) For each  $t$ ,  $\lambda_{ut}^\top(d_t)\Sigma_t\lambda_{ut}(d_t) + \sigma_t^2(d_t) \neq 0$ ,  $\sigma_t(d_t) \neq 0$  and  $\Sigma_u(x_1)$  is non-singular. (D)  $dF_{X_k^*|Y^{t-1}, X^t, D^t}(v_k; y^{t-1}, x^t, d^t) > 0$  for all  $t$  and  $v_k$  in the support of  $X_k^*$ . (E) For each  $t$  and  $d_t$ , the variance-covariance matrix of  $X_t$  conditional on  $D_t = d_t$  is non-singular.

Assumption **KL4** places support restrictions on various objects of the model. In particular, Part (B) imposes that the known latent factor  $X_k^*$  has compact support. This holds if the distribution of  $X_k^*$  has discrete support although this obviously applies to a broader set of distributions. We return to this compactness condition in Remark 1 below. Part (C) requires that the distribution of  $Y_t(d) \mid (X_t, D_t, X_k^*)$  is non-degenerate. Part (D) is a “rectangular” support assumption on  $X_k^*$ . This assumption is typically satisfied in dynamic discrete choice models, as they generally impose a large support assumption on the random utility shocks. Finally, Part (E) imposes that there exists sufficient variation in  $X_t$  conditional on  $D_t$ .

**Assumption KL5.** (A) For each  $d_t$  there are sequences  $d^{t-1}, \tilde{d}^{t-1}$  such that  $\lambda_{ut}(d_t)^\top \Sigma_t \sum_{s=1}^{t-1} \left( \lambda_{us}(d_s) \frac{\lambda_{ks}(d_s)}{\sigma_s^2(d_s)} - \lambda_{us}(\tilde{d}_s) \frac{\lambda_{ks}(\tilde{d}_s)}{\sigma_s^2(\tilde{d}_s)} \right) \neq 0$ . (B) For all  $d_t$ ,  $\lambda_{kt}(d_t) \neq 0$ . (C) For all  $d^t$ ,  $\lambda_{kt}(d_t) - \lambda_{ut}(d_t)^\top \Sigma_t \sum_{s=1}^{t-1} \lambda_{us}(d_s) \frac{\lambda_{ks}(d_s)}{\sigma_s^2(d_s)} \neq 0$ . (D) For each  $(d_2, d_1)$ ,  $\lambda_{u2}(d_2)^\top \Sigma_2 \lambda_{u1}(d_1) \frac{\lambda_{k1}(d_1)}{\sigma_1^2(d_1)} \neq 0$  (E) There are sets  $\{d_{2,i} \in \mathcal{S}(D_2) : i = 1, 2, \dots, p\}$ ,

$\{\tilde{d}_{2,i} \in \mathcal{S}(D_2) : i = 1, 2, \dots, p\}$  which satisfy

$$\begin{aligned} & (\lambda_{u2}(d_{2,1})\lambda_{u2}(d_{2,2}) \dots \lambda_{u2}(d_{2,p}))^{-\top} \text{vec}(\lambda_{k2}(d_{2,1}), \dots, \lambda_{k2}(d_{2,p})) \\ & \neq \left( \lambda_{u2}(\tilde{d}_{2,1})\lambda_{u2}(\tilde{d}_{2,2}) \dots \lambda_{u2}(\tilde{d}_{2,p}) \right)^{-\top} \text{vec}(\lambda_{k2}(\tilde{d}_{2,1}), \dots, \lambda_{k2}(\tilde{d}_{2,p})). \end{aligned}$$

(F) For any  $d^T$ ,  $\{\lambda_{ut}(d_t) : t = 1, \dots, T\}$  is linearly independent.

Assumption **KL5** is a regularity condition that ensures that the latent individual effect  $X^*$  alters outcomes sufficiently differently across time and assignments. This condition is relatively mild as it primarily rules out knife-edge cases where the cumulative effect of different elements of the individual effect perfectly offset each other.<sup>3</sup> More specifically, Part (A) requires that the aggregate effect of  $X_k^*$  on outcomes for choice  $d_t$  is different for at least two histories  $(d^{t-1}, \tilde{d}^{t-1})$ . Part (B) assumes that the direct effect of  $X_k^*$  is non-zero in each period for each assignment. Part (C) states the aggregate effect of  $X_k^*$  on outcomes must be non-zero—that is, that the direct effect  $\lambda_{kt}(d_t)$  is not perfectly offset by the effect mediated through previous choices. Part (D) ensures that there is a non-zero effect of previous choices in  $t = 2$ . Part (E) requires that in  $t = 2$  the relative effect of known and unknown  $X^*$  changes across choices. In the special case where  $X_u^* \in \mathbb{R}$ , the condition reduces to  $\frac{\lambda_{k2}(d_2)}{\lambda_{u2}(d_2)} \neq \frac{\lambda_{k2}(\tilde{d}_2)}{\lambda_{u2}(\tilde{d}_2)}$ , i.e., that the ratio of factor loadings is non-constant across assignments. More generally, for  $X_u^* \in \mathbb{R}^p$ , this condition implies that, at least for  $t = 2$ , the set of assignments must contain at the minimum  $p + 1$  elements. Finally, part (F) requires that the unknown factor affects each outcome via a different linear combination.

Define  $g_t = dF_{X_t|Y^{t-1}D^{t-1}X^{t-1}}$ . We are now in a position to state our main identification result for the model parameters  $\theta = ((\beta_t, \lambda_t, \sigma_t, h_t, g_t)_{t=1}^T, \Sigma_u, F_{X_k^*X_1}) \in \Theta$ .

**Theorem 1.** Suppose the distribution of  $(Y_t, D_t, X_t)_{t=1}^T$  is observed for  $T = 2p + 1$  periods, and that Assumptions **KL1-KL5** hold. Then  $\theta$  is point identified.

---

<sup>3</sup>This type of assumption is similarly required in latent factor models without selection or learning in order to rule out degeneracies (see, e.g., Freyberger, 2018, Assumption L4).

The proof of this theorem relies on the normality of the error term  $\epsilon_t(d)$ . The first step is to show that  $Y_t$  is normally distributed conditional upon lagged outcomes  $Y^{t-1}$ , assignments  $D^t$ , covariates  $X^t$  and the known component of the latent individual effect  $X_k^*$ . This implies that  $Y_t$  conditional upon  $(Y^{t-1}, D^t, X^t)$  is a mixture distribution parameterized by  $X_k^*$ . Then under the compact support and non-degeneracy assumptions (Assumptions [KL4](#) (A)-(C)), one can apply a result from [Bruni and Koch \(1985\)](#) to identify the aforementioned mixture distribution up to an affine transformation of  $X_k^*$ . Next, the normalization and regularity assumptions (Assumptions [KL3-KL5](#)) are used to pin down the affine transformation, leading to identification of the joint distribution of  $(Y^T, D^T, X^T, X_k^*)$ . Knowledge of this distribution identifies the components of the model related to the known component of the latent individual effect, namely  $((\beta_t, \lambda_{kt}, h_t)_{t=1}^T, F_{X_k^* X_1})$ . Thus it remains to disentangle the effect of the learned component (i.e.,  $X_u^*$ ) and uncertainty (i.e.,  $\epsilon_t(d)$ ) in order to identify  $((\lambda_{ut}, \sigma_t)_{t=1}^T, \Sigma_u)$ . To do so, we show that the joint distribution of  $(Y^T, D^T, X^T)$  conditional upon  $X_k^*$ , suitability weighted by the assignment probabilities, is a normal-weighted mixture of normal distributions. This observation leads to identification  $((\lambda_{ut}, \sigma_t)_{t=1}^T, \Sigma_u)$  from the second moments of the reweighted distribution. See [Section A.1](#) for the formal argument.

*Remark 1* (Compact support assumption). Assumption [KL4](#) (B) imposes that the known component of the latent individual effect has bounded support. In applications, it is common to assume  $X_k^*$  has finite support with known cardinality. Assumption [KL4](#) (B) relaxes this assumption in the sense that the number of support points of  $X_k^*$  need not be known a priori, and indeed may be infinite. The assumption that the support of the mixing distribution is compact plays an important role in establishing identification.<sup>4</sup>

*Remark 2* (Normality of unknown factor). As summarized in [Lemma 1](#), an impor-

---

<sup>4</sup>Compactness is used in particular to apply the Stone - Weierstrass approximation theorem, which is a central argument in [Bruni and Koch \(1985, Theorem 1\)](#).

tant advantage of the normality assumptions (Assumption [KL2](#)) is the resulting conjugate prior with a tractable closed form. For this reason, these assumptions are very common in the applied literature. In the context of our analysis, the most important implication of these assumptions is to enable identification of the (latent) distribution of  $Y_t \mid (X_k^*, Y^{t-1}, D^t, X^t)$  from variation in the realized outcome  $Y_t$  only. First, the normality assumptions on  $\epsilon_t$  and  $X_u^*$  lead to normality of  $Y_t \mid (X_k^*, Y^{t-1}, D^t, X^t)$ , using standard Bayesian arguments. It follows that, for any given  $(Y^{t-1}, D^t, X^t) = (y^{t-1}, d^t, x^t)$ , the distribution of  $Y_t \mid (Y^{t-1}, D^t, X^t)$  is a mixture of normal distributions with mixture weights given by the distribution of  $X_k^* \mid (Y^{t-1}, D^t, X^t)$ . Identification then follows from existing results for mixtures of normal distribution (Bruni and Koch, 1985).

This discussion also highlights why we restrict  $X_k^*$  to be a scalar random variable. Namely, that identification of its distribution is coming from variation in the scalar outcome variable  $Y_t$ . If a vector of outcomes were available—that is, if  $Y_t$  was vector-valued—then we expect our arguments to easily extend to multivariate  $X_k^*$ .

*Remark 3* (Invariance to normalization). The normalization assumption (Assumption [KL3](#)) is a true normalization in the sense that particular meaningful economic parameters are invariant to the assumption. In particular, we can show that average and quantile structural functions are identified without the normalization assumption. To formalize this notion, define  $C_{kt}(d) := X_k^* \lambda_{kt}(d)$ ,  $C_{ut}(d) := (X_u^*)^\top \lambda_{ut}(d)$  and let  $Q_\alpha[X]$  be the  $\alpha$ -quantile of the random variable  $X$ . Let  $x \in \mathcal{S}(X_t)$  and define the quantile structural functions

$$s_{1,t}(x, \alpha) = x^\top \beta_t(d) + Q_\alpha[C_{kt}(d) + C_{ut}(d) + \epsilon_t(d)],$$

$$s_{2,t}(x, \alpha_1, \alpha_2, \alpha_3) = x^\top \beta_t(d) + Q_{\alpha_1}[C_{kt}(d)] + Q_{\alpha_2}[C_{ut}(d)] + Q_{\alpha_3}[\epsilon_t(d)],$$

and the average structural function as  $s_{3,t}(x) = X_t^\top \beta_t(d) + \int e dF_{C_{kt}+C_{ut}+\epsilon_t}(e)$ .

In Online Appendix [B](#) we prove the following corollary:



*Corollary 1.* Suppose the Assumptions **KL1**, **KL4** and **KL5** hold and that  $(X_u^* \mid X_1 = x_1, X_k^* = v_k) \sim N(\mu_u, \tilde{\Sigma}_u(x_1))$  and  $\epsilon_t(d) \sim N(c_t(d), \sigma_t(d)^2)$ . Furthermore, suppose for some  $(d_1, d_2, \dots, d_p)$ ,  $(\lambda_{u1}(d_1)\lambda_{u2}(d_2) \dots F_{up}(d_p))$  is full rank. Then  $s_{1,t}$ ,  $s_{2,t}$  and  $s_{3,t}$  are identified on the support of  $X_t$ .

### 3.3 Pure learning model

This section considers a special case of the model of Section 2, in which all components of the latent individual effect are initially unknown to the decision making agent. That is,  $X^* = X_u^*$ . Without needing to distinguish initially known and unknown heterogeneity, a stronger identification result is achieved. In particular, no parametric restrictions on the distribution of the unobservables are required.

**Assumption L1.**  $Y_t(d) = X_t^\top \beta_t(d) + (X_u^*)^\top \lambda_{ut}(d) + \epsilon_t(d)$ , where  $X_t$  includes a constant term. For any  $d \in \mathcal{S}(D_t)$ ,

$$dF_{\epsilon_t(d)D_tX_t|Y^{t-1}D^{t-1}X^{t-1}X^*} = dF_{\epsilon_t(d)}dF_{D_t|Y^{t-1}D^{t-1}X^t}dF_{X_t|Y^{t-1}D^{t-1}X^{t-1}}.$$

Assumption **L1** adapts Assumption **KL1** to reflect there is no initially-known component of unobserved heterogeneity.

**Assumption L2.** (A) The joint density of  $Y, X^*$  and  $D, X$  admits a bounded density with respect to the product measure of the Lebesgue measure on  $\mathcal{S}(Y) \times \mathcal{S}(X^*)$  and some dominating measure on  $\mathcal{S}(D) \times \mathcal{S}(X)$ . All marginal and conditional densities are bounded. (B)  $X^* \mid X_1$  has full support. (C) The characteristic function of  $\epsilon_t(d)$  is non-vanishing,  $\mathbb{E}[\epsilon_t] = 0$ .

Assumption **L2** substantially weakens Assumption **KL2** by replacing the normality assumption with a full support assumption. Note that a full support assumption on  $Y_t(d)$  is implied by Assumption **KL2**. Let  $X^* \in \mathbb{R}^p$ .

**Assumption L3.** For some choice sequence  $(d_t: t = 1, 2, \dots, p)$ , (A)  $(\lambda_1(d_1) \dots F_p(d_p)) = I_{p \times p}$  and (B) The element of  $\beta_t(d_t)$  associated with the constant component of  $X_t$  is zero.

**Assumption L4.** (A) For each  $(y^{t-1}, x^t) \in \mathcal{S}(Y^{t-1}, X^t)$ ,  $\Pr(D_t = d \mid Y^{t-1} = y^{t-1}, X^t = x^t) > 0$  for all  $d \in \mathcal{S}(D_t)$ . (B) The variance-covariance matrix of  $X^* \mid X_1$  is full rank. (C) The variance-covariance matrix of  $(X_t)$  conditional upon  $D_t = d_t$  is non-singular.

Assumption L3 are normalization assumptions, which are standard in interactive fixed effect models. An alternative normalization could be placed on the expectation of  $X^*$  conditional upon  $X_1$ . Assumption L4 (A) is similar to Assumption KL4 (D). It requires that for each history  $(y^{t-1}, d^{t-1}, x^t)$ , some units are assigned to  $D_t = d_t$  for each  $d_t \in \mathcal{S}(D_t)$ . This assumption is satisfied in many standard parametric discrete choice models (see, e.g., Keane and Wolpin, 1997). At the cost of notational burden, this assumption could be weakened to hold for certain sequences of choices. In particular, that for each  $d_t \in \mathcal{S}(D_t)$ , there is a finite sequence of choice sequences whose first element is the choice sequence of Assumption L3 (A), whose adjacent elements are equal on at least  $p$  points of their domain, and whose final element maps  $t$  to  $d_t$ .

**Assumption L5.** For any  $d^T$ ,  $\{\lambda_{ut}(d_t) : t = 1, \dots, T\}$  are linearly independent.

Assumption L5 is a standard assumption in the interactive fixed effect literature (Assumption N6, Freyberger, 2018). Similar to Assumption KL5, it rules out degeneracies by ensuring that the outcome in each period  $Y_t(d_t)$  depends on a distinct linear combination of  $X_u^*$ .

We now define the period  $t$  conditional choice probability function as

$$h_t(y^{t-1}, d^t, x^t) := \Pr(D_t = d_t \mid Y^{t-1} = y^{t-1}, D^{t-1} = d^{t-1}, X^t = x^t).$$

Unlike in Section 3.2, the CCP function does not depend on any latent variable and is thus identified directly from the data. As in Section 3.2, we place very little structure on the learning process of decision making agents. This highlights that the core identification results do not rely on structure imposed on the belief formation process. However it is worth emphasizing that, should there be such structure, our identification results would enable identification of the belief formation process. To illustrate this, consider the case that the decision making agents are rational and Bayesian updaters and that the sufficient statistics for  $X_u^*$  at time  $t$  are a known function of the information set and the model parameters. That is, that there is a known function  $g$  such that the sufficient statistics equal  $g(Y^{t-1}, D^{t-1}, X^{t-1}, \theta)$ , where  $\theta$  are the model parameters. In this case, identification of  $\theta$  is sufficient for identification of the beliefs.

Define  $f_{\epsilon_t} = \{f_{\epsilon_t(d)} : d \in \mathcal{S}(D_t)\}$ . Let the model parameter vector be  $\theta = ((\beta_t, \lambda_t, \sigma_t, h_t, f_{\epsilon_t}, g_t)_{t=1}^T, \Sigma_u, f_{X_k^* X_1}) \in \Theta$ . The following theorem states that the preceding conditions are sufficient for point identification of  $\theta$ .

**Theorem 2.** Suppose the distribution of  $(Y_t, D_t, X_t)_{t=1}^T$  is observed for  $T = 2p + 1$  and that Assumptions L1-L5 hold. Then  $\theta$  is point identified.

The key insight that enables identification of  $\theta$  is that, under Assumption L1, this is a model of selection on observables. That is, although assignment probabilities depend on unobserved beliefs over  $X^*$ , they do not depend on the unobserved factor  $X^*$  itself. It follows that one can control for beliefs at time  $t$  by conditioning upon prior outcomes, choices and covariates. This in turn allows us to express the joint distribution of  $(Y^t, D^t, X^t)$ , suitably weighted by the assignment probabilities, as a mixture model over the potential outcomes  $Y^t(d_t)$  conditional upon the latent factor  $X^*$  and exogenous covariates  $X$ . From here the arguments of Freyberger (2018) yield identification of the mixture and component distributions. See Section A.2 for details.

*Remark 4* (Auxiliary measurements). In some cases, additional unselected noisy measurements of known abilities are available. See, for instance, Cunha et al. (2005)

and Heckman and Navarro (2007). With this additional data, sufficient conditions for identification of the distribution of the latent effect are well known in the literature (Hu and Schennach, 2008; Cunha et al., 2010). If the sufficient conditions are satisfied conditional on each  $(Y_t, D_t, X_t)_{t=1}^T$ , then the joint distribution of  $((Y_t, D_t, X_t)_{t=1}^T, X_k^*)$  is identified from the auxiliary measurements. From here, one can redefine  $X_t = (X_t, X_k^*)$  and the conditions of Theorem 2 are sufficient for distribution-free identification of the model with known and unknown heterogeneity.

## 4 Estimation

We propose to estimate the model parameters via sieve maximum likelihood. Let  $W_i = (Y_{it}, D_{it}, X_{it} : t = 1, \dots, T)$  and  $\theta^* \in \Theta$  be the true value of the parameters. In the following we focus on the model of Section 3.2, although we conjecture related conditions could be presented for the model of Section 3.3. The log-likelihood contribution of  $W_i = w$  is

$$\begin{aligned} \ell(w; \theta) = & \log \int \int \prod_{t=1}^T \frac{1}{\sigma_t(d_t)} \phi_1 \left( \frac{y_t - x_t^\top \beta_t(d_t) - v_k F_{kt}(d_t) - v_u^\top F_{ut}(d_t)}{\sigma_t(d_t)} \right) \\ & \times \prod_{t=1}^T h_t(d^t, x^t, y^{t-1}, v_k) \times \prod_{t=1}^{T-1} g_t(x_{t+1}; y^t, d^t, x^t) \\ & \times \frac{1}{\sqrt{|\Sigma_u(x_1)|}} \phi_p \left( \Sigma_u^{-\frac{1}{2}}(x_1) x_u^* \right) \times dv_u dF_{X_k^* X_1}(v_k, x_1) \end{aligned} \quad (5)$$

where  $\phi_s$  is the probability distribution function of the standard multivariate normal distribution with  $s$  components,  $g_t$  is the distribution of  $X_{t+1}$  conditional upon  $(Y^t, D^t, X^t) = (y^t, d^t, x^t)$ . There are four components of the likelihood function: the outcomes, the assignment probabilities, the distribution of the covariates, and the joint distribution of  $(X_1, X^*)$ .

To estimate  $\theta$ , let  $\Theta_n$  be a finite dimensional sieve space that serves as an approx-

imation to  $\Theta$ . The sieve maximum-likelihood estimator for  $\theta^*$  is defined as

$$\frac{1}{n} \sum_{i=1}^n \ell(w_i; \hat{\theta}) \geq \sup_{\theta \in \Theta_n} \frac{1}{n} \sum_{i=1}^n \ell(w_i; \theta) - o_p(1/n) \quad (6)$$

The following result states that under standard conditions (stated in Appendix D.1), if our model is identified,  $\hat{\theta}$  is consistent for  $\theta^*$ ,

**Theorem 3.** Let  $(W_i)_{i=1}^n$  be i.i.d. data where  $T \geq 2p + 1$  and Assumptions KL1-KL5 and Assumptions E1-E5 hold. Then  $\hat{\theta}$  as defined in Equation (6) is consistent for  $\theta^*$ .

Researchers are often interested in functionals of the model parameters, such that the variance decompositions discussed in section 2.1. These variance decompositions involve both the finite dimensional parameters of the model as well as the distribution of  $X_k^*$  and the CCPs. Therefore, many of the existing results on inference on the finite dimensional parameters of a semiparametric model (e.g. Ai and Chen, 2003) do not directly apply to this setting.

Instead, we next provide an inference result for a plug-in estimator of a more general class of functionals of the model parameters. For a functional  $f$ , under a set of smoothness and regularity conditions similar to those given in Chen and Liao (2014), we show establish that the plug-in estimator  $f(\hat{\theta})$  has an asymptotically normal distribution and characterize its asymptotic variance.

**Theorem 4.** Let  $(W_i)_{i=1}^n$  be i.i.d. data where  $T \geq 2p + 1$  and that Assumptions KL1-KL5 and Assumptions E1-E13 hold. Then  $\sqrt{n} \frac{f(\hat{\theta}) - f(\theta^*)}{\|v_n^*\|} \xrightarrow{d} \mathcal{N}(0, 1)$  where  $v_n^*$  is the sieve Riesz representer of  $f(\theta)$  and  $\|\cdot\|$  is defined in Equation (13) in the online appendix.

The asymptotic variance and rate of convergence of the plug in sieve estimator depend on the  $\|v_n^*\|$ . For *regular* functionals,  $\|v_n^*\|$  converges to a constant, which implies that the plug-in estimator has a root-n convergence rate. Note, however, that Theorem 4 also allows for the possibility that the sieve variance  $v_n^*$  may diverge—that

is, that  $f$  is an *irregular* functional. In either case, consistent estimators for the sieve variance are available (Chen and Liao, 2014, Section 3).

We leave it to future work to derive primitive conditions under which functionals such as the variances decompositions discussed in section 2.1 satisfy the high level conditions of Theorem 4. However, we do provide in Appendix D.2 lower level conditions under which Theorem 4 holds.

## 5 Implementation and Monte Carlo Experiments

In this section we show how the sieve MLE estimator can be tractably implemented and perform a Monte Carlo experiment illustrating the good finite sample performance of the estimator.

### 5.1 Implementation

To implement the sieve maximum likelihood estimator proposed in the previous section, first partition  $\theta$  into  $\{F_{X_k^*|X_1}\}$  and  $\theta_2 := \theta \setminus \{F_{X_k^*|X_1}\}$ .<sup>5</sup> Integrating out  $X_u^*$  in (5), we obtain  $\ell(w; \theta_2, F_{X_k^*|X_1}) = \log \int f(w, x_k^*; \theta_2) dF_{X_k^*|X_1}(x_k^*; x_1)$  where,

$$\begin{aligned} f(w, x_k^*; \theta_2) := & \frac{1}{\sqrt{|V(w, x_k^*; \theta_1)|}} \phi_T \left( V(w, x_k^*; \theta_1)^{-\frac{1}{2}} (y^T - m(w, x_k^*; \theta_1)) \right) \\ & \times \prod_{t=1}^T h_t(d^t, x^t, y^{t-1}, x_k^*) \times \prod_{t=1}^{T-1} g_t(z_{t+1}; x_t, y_t, d_t), \end{aligned}$$

and  $m(w, x_k^*; \theta_1)$  and  $V(w, x_k^*; \theta_1)$  are the  $T$ -dimensional vector and  $T \times T$  matrix giving the expected mean and variance of  $Y^T$  conditional on  $(D^T, X^T, X_k^*) = (d^T, x^T, x_k^*)$ . They are defined as follows. Let  $\beta(w) = [\beta_1(d_1) \ \cdots \ \beta_T(d_T)]$ ,

---

<sup>5</sup>Note that  $\theta_1 \subset \theta_2$ , where as previously defined,  $\theta_1 = ((\alpha_t, \beta_t, F_{kt}, F_{ut}, \sigma_t)_{t=1}^T, \Sigma_u(x_1))$

$\Lambda_k(w) = [\lambda_{k1}(d_1) \ \cdots \ \lambda_{kT}(d_T)]$ , and  $\Lambda_u(w) = [\lambda_{u1}(d_1) \ \cdots \ \lambda_{uT}(d_T)]$ . Then,

$$m(w, x_k^*; \theta) = \beta(w)^\top x + \Lambda_k(w)^\top x_k^*,$$

$$V(w, x_k^*; \theta) = \Lambda_u(w)^\top \Sigma_u(x_1) \Lambda_u(w) + \text{diag}(\sigma_1^2(d_1), \dots, \sigma_T^2(d_T)),$$

There are three non-parametric objects in the likelihood function:  $h$ ,  $g$ , and  $F_{X_k^*|X_1}$ . These can be estimated non-parametrically using a sieve space, or a parametric form can be imposed on any of them. The choice of a model or sieve spaces for  $h$  and  $g$  are typically context specific. For  $F_{X_k^*|X_1}$ , we propose using a sieve space closely related to estimator discussed in Koenker and Mizera (2014) and Fox et al. (2016). Assume that  $X_1$  has finite support,  $(\bar{x}_1, \dots, \bar{x}_R)$ . For each  $n$ , fix a grid of support for  $X_k^*$  with  $q_n < \infty$  points,  $\mathcal{S}_n = \{\bar{x}_{1n}^*, \dots, \bar{x}_{q_n n}^*\}$ . We can then use following sieve space,

$$\mathcal{F}_n = \left\{ (x^*; \bar{x}_r) \mapsto \sum_{s=1}^{q_n} \omega_{sr} \mathbf{1}\{x^* \leq \bar{x}_{sn}^*\} \mid \omega \in \Delta^R(q_n) \right\}$$

where  $\Delta^k(m) = \{\omega \in [0, 1]^{m \times k} : \sum_{s=1}^m \omega_{si} = 1, 1 \leq i \leq k\}$  is the  $k$ -product of  $m$ -dimensional simplexes<sup>6</sup>. Notice that  $\mathcal{F}_n$  is simply the space of conditional distributions with support contained in  $\mathcal{S}_n$ . If  $\mathcal{S}_n$  becomes dense in  $\mathbb{R}$  and the number of points grows at a suitable rate, this sieve space satisfies the conditions of Theorems 3 and 4.

This sieve space is particularly convenient computationally. Each  $F \in \mathcal{F}_n$  corresponds to an  $\omega \in \Delta^R(q_n)$  such that,

$$\sum_{i=1}^n \ell(w_i; \theta_2, F) = \log \sum_{s=1}^{q_n} \omega_{sr} \sum_{i=1}^n \sum_{r=1}^R \mathbf{1}(x_{i1} = \bar{x}_r) f(w_i, \bar{x}_{sr}^*; \theta_2).$$

Holding  $\theta_2$  fixed, maximizing over  $F \in \mathcal{F}_n$  therefore amounts to maximizing a concave

---

<sup>6</sup>If  $X_1$  is independent of  $X_k^*$ , then this reduces to a  $\Delta^1(q_n)$ .

objective function over the linearly constrained set,  $\omega \in \Delta^R(q_n)$ . This problem can be solved efficiently and reliably using standard software for convex optimization. For example, the algorithm proposed in Kim et al. (2020), is specialized for this setting and implemented in the R package *mixsqp*.

This allows us to calculate the profile log likelihood,  $\theta_2 \mapsto \max_{F \in \mathcal{F}_n} \sum_{i=1}^n \ell(w_i; \theta_2, F)$ , and the full MLE problem can be solved by maximizing this function in  $\theta_2$ . In appendix E.1 we show how the gradient of the profile log likelihood function can be calculated implicitly making it feasible to use first order optimization algorithms to maximize the profile log likelihood function over  $\theta_2$  efficiently.

## 5.2 Monte Carlo simulations

Next, we present results from Monte Carlo simulations which illustrate the computational tractability and finite-sample performance of the proposed estimator. We focus here on a simple model with a parametric assignment model.

The data generating process (DGP) used in the simulations is based on the model in Section 3.2 with both known and unknown heterogeneity. We include two time-invariant covariates, one continuous and one discrete, which are independent of  $X^*$ . Assignment probabilities are derived from a model in which agents maximize the following utility function,

$$u_t(d, X_k^*, Y^{t-1}, X, D^{t-1}) = \rho \mathbb{E}(Y_t(d) | X_k^*, Y^{t-1}, X, D^{t-1}) + \rho \gamma \mathbf{1}(D_t = 2) X_k^* + \nu_t(d),$$

where  $\{\nu_t(d) : t = 1, 2, 3, d = 1, 2\}$  are mutually independent with an Extreme Value Type 1 distribution. This utility function puts a weight on the expected outcome of their choice and another term which depends on  $X_k^*$ . This additional term can reflect biased beliefs, heterogeneity in preferences, or a combination of both.  $X_k^*$  is distributed as a mixture of truncated normal random variables. The parameter values



used in the simulations are reported in Appendix E.2.

Table 1: Average computational times by sample size

	N = 250	N = 500	N = 1000	N = 2000	N = 4000
Computational Time	0:24	0:31	0:55	2:15	3:32

We perform a Monte Carlo experiment, estimating parameters of the model with 200 simulations and sample sizes of 250, 500, 1000, 2000 and 4000. We use the sieve MLE estimator described in Section 4, maintaining the parametric structure on the assignment probabilities but estimating  $F_{X_k^*}$ <sup>7</sup> nonparametrically using the sieve space described in Section 5.1. The number of support points in the estimated distributions,  $q_n$ , grows at a rate of  $n^{1/3}$ , from 62 to 158. There are 32 parameters in  $\theta_2$ .<sup>8</sup>

With this specification, computation remains tractable for these sample sizes. Average computational times to solve the profile MLE problem reported in Table 1 run from half a minute to around three and half minutes.

---

<sup>7</sup>Since  $X_1$  is independent of  $X_k^*$ ,  $F_{X_k^*|X_1} = F_{X_k^*}$

<sup>8</sup>We include the finite parameters of the assignment model in  $\theta_2$  since it is modeled parametrically. Since we assume that  $Z$  is time invariant,  $g$  does not need to be estimated.

Table 2: Bias and Variance ( $\times 1,000$ ) of Finite Parameter Estimators

	N = 250		N = 500		N = 1,000		N = 2,000		N = 4,000	
	sq bias	var	sq bias	var	sq bias	var	sq bias	var	sq bias	var
$\alpha_1(2)$	71.722	87.915	34.056	60.969	12.915	47.132	0.729	19.025	0.042	5.697
$\alpha_2(1)$	0.148	27.977	0.264	12.384	0.116	7.391	0.002	2.878	0.010	1.376
$\alpha_2(2)$	73.516	108.963	34.175	74.416	12.407	57.187	0.463	25.797	0.027	8.115
$\alpha_3(1)$	0.005	36.562	0.455	13.824	0.197	5.311	0.000	2.237	0.013	0.963
$\alpha_3(2)$	47.840	163.156	32.091	82.423	12.025	62.314	0.593	25.979	0.037	7.322
$\beta_{z1,1}(1)$	0.513	10.084	0.399	5.220	0.137	3.172	0.016	1.486	0.000	0.721
$\beta_{z1,1}(2)$	0.852	15.221	0.304	6.753	0.045	3.349	0.007	1.744	0.002	0.801
$\beta_{z2,1}(1)$	0.837	16.296	0.659	7.859	0.392	4.464	0.038	1.849	0.006	0.803
$\beta_{z2,1}(2)$	1.378	20.814	0.601	12.059	0.093	5.618	0.003	2.687	0.006	1.215
$\beta_{z3,1}(1)$	0.408	9.298	0.243	3.879	0.156	1.886	0.028	1.030	0.007	0.569
$\beta_{z3,1}(2)$	0.383	19.188	0.402	9.105	0.083	4.201	0.010	2.096	0.002	0.861
$\beta_{z1,2}(1)$	0.606	58.909	0.362	23.238	0.361	11.159	0.027	4.772	0.004	2.295
$\beta_{z1,2}(2)$	0.187	46.657	0.218	25.403	0.022	11.158	0.002	5.124	0.010	2.611
$\beta_{z2,2}(1)$	0.005	40.411	0.001	19.844	0.001	9.046	0.002	4.347	0.041	2.476
$\beta_{z2,2}(2)$	0.038	57.755	0.055	26.567	0.000	12.370	0.000	6.761	0.012	3.295
$\beta_{z3,2}(1)$	0.495	40.189	0.079	19.943	0.017	7.637	0.000	3.943	0.022	2.046
$\beta_{z3,2}(2)$	0.102	65.654	0.332	32.107	0.014	15.179	0.024	7.111	0.000	3.438
$F_{k1}(1)$	2.746	27.521	1.701	12.890	0.624	7.268	0.009	3.684	0.001	1.468
$F_{k2}(1)$	1.148	25.977	0.558	10.825	0.226	4.777	0.004	2.594	0.000	1.089
$F_{k2}(2)$	0.869	10.978	0.254	5.825	0.073	2.654	0.007	1.383	0.000	0.743
$F_{k3}(1)$	3.986	33.663	0.873	13.719	0.178	5.683	0.003	3.069	0.000	1.330
$F_{k3}(2)$	5.702	36.861	0.674	12.556	0.224	5.305	0.011	2.408	0.005	1.080
$F_{u1}(2)$	0.979	13.945	0.306	4.733	0.170	2.438	0.015	1.330	0.000	0.605
$F_{u2}(1)$	0.040	8.317	0.027	5.139	0.036	1.947	0.014	1.003	0.002	0.481
$F_{u2}(2)$	1.478	14.880	0.494	6.218	0.130	3.324	0.009	1.522	0.004	0.643
$F_{u3}(1)$	0.446	9.912	0.093	5.003	0.062	2.187	0.030	0.968	0.023	0.469
$F_{u3}(2)$	0.106	21.919	0.101	8.903	0.112	4.148	0.004	2.140	0.005	0.936
$\sigma^2(1)$	0.453	2.477	0.091	1.241	0.030	0.672	0.007	0.298	0.001	0.136
$\sigma^2(2)$	1.228	4.449	0.242	2.237	0.027	1.058	0.016	0.701	0.008	0.332
$\sigma_u^2$	0.017	72.902	0.052	41.174	0.043	17.906	0.012	9.337	0.010	4.335
$\gamma$	3.791	103.555	0.562	34.226	0.192	18.098	0.357	10.611	0.043	5.779
$\rho$	3.317	121.851	0.165	51.065	0.031	25.777	0.012	11.329	0.077	4.724

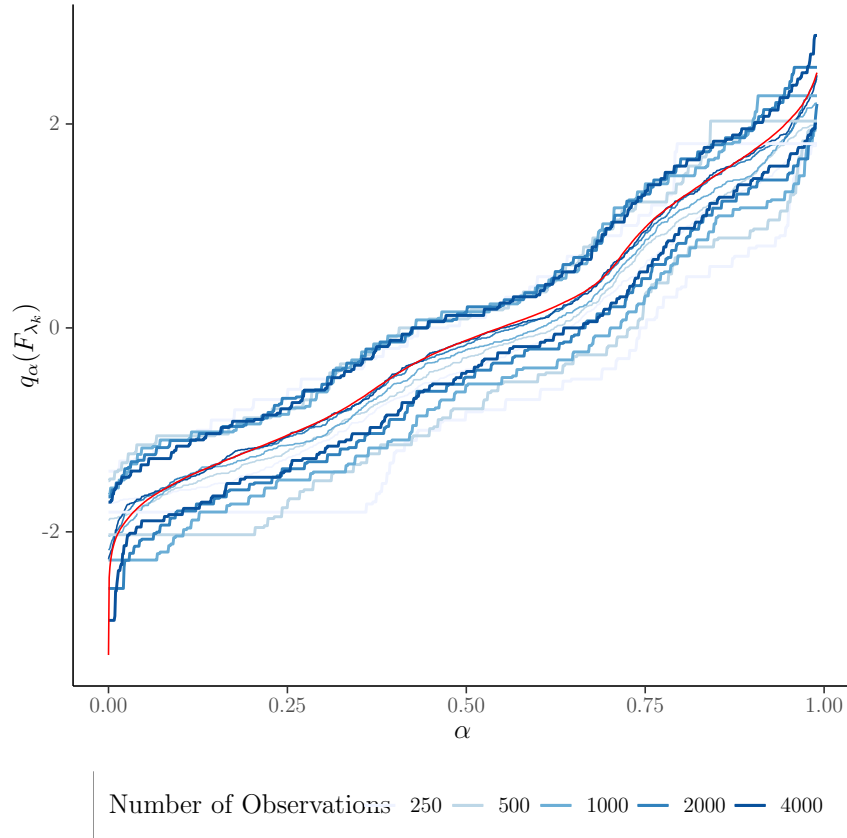
All calculations are based on 200 Monte Carlo simulations of the DGP described in the main text. Squared bias and variance of finite parameter estimates are multiplied times 1,000

The squared bias and variance of the sieve estimator of  $\theta_2$  are presented in Table 2. (Note that all values in Table 2 are multiplied by 1,000.) For each of the parameters, the bias becomes negligible for all parameters relative to variance as sample size grows. For most parameter except the intercepts, this bias is small even for small sample sizes. The variance declines at a rate consistent with  $\sqrt{n}$  convergence of the mean squared error. This is consistent with Theorem 4 since the functional mapping projecting the full parameter space into the finite dimensional part of the parameter space is known to be a regular functional.

To present results for the nonparametric estimator of the distribution of known unobserved heterogeneity  $F_{X_k^*}$ , we focus on the quantiles of  $F_{X_k^*}$ . Let  $q_\alpha(F)$  be the  $\alpha$

quantile of a random variable with the distribution  $F$ . For each value of  $\alpha \in [0, 1]$ , we calculate the mean and the 5th and 95th percentile of the simulated distribution of the estimator of  $q_\alpha(F_{X_k^*})$ . The results are presented in Figure 1. The red line shows the CDF of the true distribution of  $X_k^*$ , while the blue lines that closely follow the red line are the mean of the simulated distribution of the quantile estimators for each sample size. Darker blue lines represent larger sample sizes. The blue lines above and below the CDF are the 5th and 95th percentiles of the simulated distribution of the quantile estimators.

Figure 1: Quantiles of Estimator of  $X_k^*$ : 95% Coverage Intervals



Note: The red line shows the true distribution of  $X_k^*$ . The blue lines show the mean, and the 5th and 95th percentiles of the simulated distribution of the estimate of  $q_\alpha(f_{X_k^*})$ .

The results indicate that the bias of the quantile estimators becomes negligible in moderate sample sizes. The estimator broadly captures the shape of the true

distribution of  $X_k^*$ , and also appears to converge toward the true distribution as the sample size grows. We do not provide a formal result on the rate of convergence of this parameter, but we expect this nonparametric estimator to converge at rate a slower than  $\sqrt{n}$ . At a sample size of  $n = 4,000$ , the simulated distribution of this estimator is still relatively disperse.

Finally, we consider the plug-in estimator for one of the functionals discussed in Section 2.1 and Appendix C. We focus on the decomposition of the present value of a stream of outcomes into forecastable and non-forecastable components at  $t = 1$ . With a discount rate of 0.95, for each choice sequence  $(d_1, d_2, d_3)$ , the variance of the unknown and known components are,

$$\begin{aligned} \text{unknown :} \quad & \sigma_u^2 \sum_{1 \leq t_1, t_2 \leq 3} (.95)^{t_1+t_2-2} X_{ut_1}^*(d_{t_1}) X_{ut_2}^*(d_{t_2}) + \sum_{1 \leq t \leq 3} (.95)^{2t-2} \sigma_t^2(d_t) \\ \text{known :} \quad & \text{Var}(X_k^*) \sum_{1 \leq t_1, t_2 \leq 3} (.95)^{t_1+t_2-2} X_{kt_1}^*(d_{t_1}) X_{kt_2}^*(d_{t_2}) \end{aligned}$$

We estimate these functionals, which involve both  $\theta_2$  and  $F_{X_k^*}$  using the plug-in estimator described in section 4. The results are presented in Table 3. For moderate sample sizes, the squared bias is small relative to the variance, and like the estimators of  $\theta_2$ , the variance appears to be consistent with a  $\sqrt{n}$  convergence rate.

Table 3: Bias and Variance of the Period 0 Lifetime Earnings Estimators

		N = 250		N = 500		N = 1000		N = 2000		N = 4000	
		sq bias	var	sq bias	var	sq bias	var	sq bias	var	sq bias	var
(1, 1, 1)	known	0.007	0.994	0.000	0.451	0.001	0.214	0.002	0.136	0.001	0.067
(1, 1, 1)	unknown	0.001	3.060	0.001	1.512	0.000	0.677	0.001	0.332	0.000	0.154
(1, 1, 2)	known	0.003	1.456	0.012	0.697	0.005	0.380	0.001	0.225	0.001	0.095
(1, 1, 2)	unknown	0.002	2.324	0.000	1.134	0.001	0.516	0.000	0.268	0.000	0.117
(1, 2, 1)	known	0.313	1.774	0.130	0.931	0.041	0.530	0.000	0.282	0.000	0.111
(1, 2, 1)	unknown	0.031	1.723	0.003	0.853	0.001	0.368	0.000	0.192	0.000	0.087
(1, 2, 2)	known	0.218	3.134	0.173	1.533	0.047	0.876	0.000	0.412	0.000	0.155
(1, 2, 2)	unknown	0.011	1.201	0.006	0.599	0.003	0.284	0.000	0.154	0.000	0.062
(2, 1, 1)	known	0.235	1.492	0.065	0.819	0.018	0.363	0.000	0.223	0.000	0.097
(2, 1, 1)	unknown	0.028	1.755	0.003	0.850	0.007	0.363	0.002	0.165	0.001	0.084
(2, 1, 2)	known	0.143	2.432	0.075	1.125	0.021	0.560	0.000	0.324	0.001	0.137
(2, 1, 2)	unknown	0.009	1.228	0.007	0.595	0.010	0.265	0.001	0.132	0.000	0.066
(2, 2, 1)	known	1.051	3.036	0.295	1.556	0.072	0.734	0.002	0.380	0.000	0.174
(2, 2, 1)	unknown	0.103	1.097	0.019	0.450	0.011	0.191	0.002	0.095	0.001	0.046
(2, 2, 2)	known	0.482	5.843	0.247	2.771	0.032	1.562	0.003	0.761	0.000	0.330
(2, 2, 2)	unknown	0.062	0.786	0.027	0.318	0.015	0.170	0.001	0.088	0.000	0.040

All calculations are based on 200 Monte Carlo simulations of the DGP described in the main text.

## 6 Conclusion

We provide new identification results for a general class of learning models, that encompasses many of the models that have been considered in the applied literature. We consider an environment where the researcher has access to panel data on choices and realized outcomes only. As such, our results are widely applicable, including in frequent environments where one does not have access to elicited beliefs data or auxiliary selection-free measurements. We show that the model is point-identified under two alternative sets of conditions. Our first set of conditions applies to a version of the learning where we assume that the idiosyncratic shocks from the outcome equations are normally distributed, a restriction that is very commonly imposed in empirical Bayesian learning models. We also show that normality can be relaxed in the case of a pure learning model, and establish identification for this class of models.

We then derive a sieve MLE estimator for the model parameters and a particular class of functionals, which includes as a leading special cases the predictable and unpredictable outcome variances. Notably, these variances can in turn be used to evaluate the relative importance of uncertainty versus heterogeneity in lifecycle earnings variability (Cunha et al., 2005). Under certain regularity conditions, the resulting

estimators are consistent and asymptotically normal. Importantly for practical purposes, the profile likelihood based estimation procedure proposed in this paper can be implemented at a modest computational cost.

## Bibliography

Abbring, J. and Campbell, J. (2005), A firm’s first year, Technical report, Tinbergen Institute Discussion Paper 05-046/3.

Ackerberg, D. A. (2003), ‘Advertising, learning, and consumer choice in experience good markets: an empirical examination’, *International Economic Review* **44**(3), 1007–1040.

Aguirregabiria, V., Gu, J. and Luo, Y. (2021), ‘Sufficient statistics for unobserved heterogeneity in structural dynamic logit models’, *Journal of Econometrics* **223**(2), 280–311.

Aguirregabiria, V. and Jeon, J. (2020), ‘Firms’ beliefs and learning: Models, identification, and empirical evidence’, *Review of Industrial Organization* **56**, 203–235.

Ai, C. and Chen, X. (2003), ‘Efficient estimation of models with conditional moment restrictions containing unknown functions’, *Econometrica* **71**(6), 1795–1843.

Antonovics, K. and Golan, L. (2012), ‘Experimentation and job choice’, *Journal of Labor Economics* **30**(2), 333–366.

Arcidiacono, P. (2004), ‘Ability sorting and the returns to college major’, *Journal of Econometrics* **121**(1-2), 343–375.

Arcidiacono, P., Aucejo, E., Maurel, A. and Ransom, T. (2023), College attrition and the dynamics of information revelation, Technical report, Duke University.

- Bai, J. (2009), ‘Panel data models with interactive fixed effects’, *Econometrica* **77**(4), 1229–1279.
- Berman, N., Rebeyrol, V. and Vicard, V. (2019), ‘Demand learning and firm dynamics: evidence from exporters’, *Review of Economics and Statistics* **101**(1), 91–106.
- Bruni, C. and Koch, G. (1985), ‘Identifiability of continuous mixtures of unknown gaussian distributions’, *The Annals of Probability* pp. 1341–1357.
- Bunting, J. (2022), ‘Continuous permanent unobserved heterogeneity in dynamic discrete choice models’, *arXiv preprint arXiv:2202.03960* .
- Chan, T. Y. and Hamilton, B. H. (2006), ‘Learning, private information, and the economic evaluation of randomized experiments’, *Journal of Political Economy* **114**(6), 997–1040.
- Chen, X. (2007), ‘Large sample sieve estimation of semi-nonparametric models’, *Handbook of econometrics* **6**, 5549–5632.
- Chen, X. and Liao, Z. (2014), ‘Sieve m inference on irregular parameters’, *Journal of Econometrics* **182**(1), 70–86.
- Chen, X., Liao, Z. and Sun, Y. (2014), ‘Sieve inference on possibly misspecified semi-nonparametric time series models’, *Journal of Econometrics* **178**, 639–658.
- Ching, A. T., Erdem, T. and Keane, M. P. (2013), ‘Learning models: An assessment of progress, challenges, and new developments’, *Marketing Science* **32**(6), 913–938.
- Compiani, G. and Kitamura, Y. (2016), ‘Using mixtures in econometric models: a brief review and some new results’, *Econometrics Journal* **19**(3), C95–C127.
- Coscelli, A. and Shum, M. (2004), ‘An empirical model of learning and patient spillovers in new drug entry’, *Journal of Econometrics* **122**(2), 213–246.

- Crawford, G. and Shum, M. (2005), ‘Uncertainty and learning in pharmaceutical demand’, *Econometrica* **73**(4), 1137–1173.
- Cunha, F. and Heckman, J. J. (2016), ‘Decomposing trends in inequality in earnings into forecastable and uncertain components’, *Journal of Labor Economics* **34**(S2), S31–S65.
- Cunha, F., Heckman, J. J. and Navarro, S. (2005), ‘Separating uncertainty from heterogeneity in life cycle earnings’, *Oxford Economic Papers* **57**(2), 191–261.
- Cunha, F., Heckman, J. J. and Schennach, S. M. (2010), ‘Estimating the technology of cognitive and noncognitive skill formation’, *Econometrica* **78**(3), 883–931.
- D’Haultfoeuille, X. (2011), ‘On the completeness condition in nonparametric instrumental problems’, *Econometric Theory* **27**(3), 460–471.
- Erdem, T. and Keane, M. P. (1996), ‘Decision-making under uncertainty: Capturing dynamic brand choice processes in turbulent consumer goods’, *Marketing Science* **15**(1), 1–20.
- Fox, J. T., il Kim, K. and Yang, C. (2016), ‘A simple nonparametric approach to estimating the distribution of random coefficients in structural models’, *Journal of Econometrics* **195**(2), 236–254.
- Freyberger, J. (2018), ‘Non-parametric panel data models with interactive fixed effects’, *The Review of Economic Studies* **85**(3), 1824–1851.
- Gobillon, L. and Magnac, T. (2016), ‘Regional policy evaluation: Interactive fixed effects and synthetic controls’, *The Review of Economics and Statistics* **98**(3), 535–551.
- Gong, Y. (2019), Signal-based learning models without the rational expectations assumption: Identification and counterfactuals, Technical report, Technical report, mimeo.



- Heckman, J. J. and Navarro, S. (2007), ‘Dynamic discrete choice and dynamic treatment effects’, *Journal of Econometrics* **136**(2), 341–396.
- Hincapié, A. (2020), ‘Entrepreneurship over the life cycle: Where are the young entrepreneurs?’, *International Economic Review* **61**(2), 617–681.
- Hu, Y. and Sasaki, Y. (2018), ‘Closed-form identification of dynamic discrete choice models with proxies for unobserved state variables’, *Econometric Theory* **34**(1), 166–185.
- Hu, Y. and Schennach, S. M. (2008), ‘Instrumental variable treatment of nonclassical measurement error models’, *Econometrica* **76**(1), 195–216.
- Hu, Y. and Shum, M. (2012), ‘Nonparametric identification of dynamic models with unobserved state variables’, *Journal of Econometrics* **171**(1), 32–44.
- Huggett, M., Ventura, G. and Yaron, A. (2011), ‘Sources of lifetime inequality’, *American Economic Review* **101**(7), 2923–2954.
- Kasahara, H. and Shimotsu, K. (2009), ‘Nonparametric identification of finite mixture models of dynamic discrete choices’, *Econometrica* **77**(1), 135–175.
- Keane, M. P. and Wolpin, K. I. (1997), ‘The career decisions of young men’, *Journal of political Economy* **105**(3), 473–522.
- Kim, Y., Carbonetto, P., Stephens, M. and Anitescu, M. (2020), ‘A fast algorithm for maximum likelihood estimation of mixture proportions using sequential quadratic programming’, *Journal of Computational and Graphical Statistics* **29**(2), 261–273.
- Kinsler, J. and Pavan, R. (2021), ‘Local distortions in parental beliefs over child skill’, *Journal of Political Economy* **129**(1), 81–100.
- Kitamura, Y. and Laage, L. (2018), ‘Nonparametric analysis of finite mixtures’, *arXiv preprint arXiv:1811.02727v1* .

- Koenker, R. and Mizera, I. (2014), ‘Convex optimization, shape constraints, compound decisions, and empirical bayes rules’, *Journal of the American Statistical Association* **109**(506), 674–685.
- Magnac, T. and Thesmar, D. (2002), ‘Identifying dynamic discrete decision processes’, *Econometrica* **70**(2), 801–816.
- Miller, R. A. (1984), ‘Job matching and occupational choice’, *Journal of Political Economy* **92**(6), 1086–1120.
- Pastorino, E. (2015), ‘Job matching within and across firms’, *International Economic Review* **56**(2), 647–671.
- Pastorino, E. (2022), ‘Careers in firms: The role of learning about ability and human capital acquisition’. Forthcoming in the *Journal of Political Economy*.
- Sasaki, Y. (2015), ‘Heterogeneity and selection in dynamic panel data’, *Journal of Econometrics* **188**(1), 236–249.
- Stange, K. M. (2012), ‘An empirical investigation of the option value of college enrollment’, *American Economic Journal: Applied Economics* **4**(1), 49–84.
- Stinebrickner, T. and Stinebrickner, R. (2012), ‘Learning about academic ability and the college dropout decision’, *Journal of Labor Economics* **30**(4), 707–748.
- Thomas, J. (2019), ‘The signal quality of grades across academic fields’, *Journal of Applied Econometrics* **34**(4), 566–587.
- Yang, N. (2020), ‘Learning in retail entry’, *International Journal of Research in Marketing* **37**(2), 336–355.
- Zafar, B. (2011), ‘How do college students form expectations?’, *Journal of Labor Economics* **29**(2), 301–348.

## A Proofs for identification section

In this section, we use the following notations:  $\phi$  denotes the standard normal p.d.f.;  $\mathcal{S}(X)$  represents the support of a random variable  $X$ .

### A.1 Proofs for Section 3.2

*Proof of Lemma 1.* We proceed inductively. First, by Assumption KL2 and the definition of  $(E_1, \Sigma_1)$ ,  $X_u^* \mid (X_1, X_k^*) \sim \mathcal{N}(E_1, \Sigma_1)$ . Second, for  $t > 1$  suppose  $X_u^* \mid (Y^{t-2}, D^{t-2}, X^{t-1}) \sim \mathcal{N}(E_{t-1}, \Sigma_{t-1})$  and that

$$\begin{aligned}
& f_{X_u^* | Y^{t-1} D^{t-1} X^t X_k^*}(X_u^*; Y^{t-1}, D^{t-1}, X^t, X_k^*) \\
& \propto_{(1)} f_{X_u^* | Y^{t-2} D^{t-2} X^{t-1} X_k^*}(X_u^*; Y^{t-2}, D^{t-2}, X^{t-1}, X_k^*) \\
& \quad \times f_{Y_{t-1} D_{t-1} X_t | Y^{t-2} D^{t-2} X^{t-1} X^*}(Y_{t-1}, D_{t-1}, X_t; Y^{t-2}, D^{t-2}, X^{t-1}, X^*) \\
& = f_{X_u^* | Y^{t-2} D^{t-2} X^{t-1} X_k^*}(X_u^*; Y^{t-2}, D^{t-2}, X^{t-1}, X_k^*) f_{X_t | Y^{t-1} D^{t-1} X^{t-1} X^*}(X_t; Y^{t-1}, D^{t-1}, X^{t-1}, X^*) \\
& \quad \times f_{Y_{t-1}(D_{t-1}) | Y^{t-2} D^{t-1} X^{t-1} X^*}(Y_{t-1}; Y^{t-2}, D^{t-1}, X^{t-1}, X^*) f_{D_{t-1} | Y^{t-2} D^{t-2} X^{t-1} X^*}(D_{t-1}; Y^{t-2}, D^{t-2}, X^{t-1}, X^*) \\
& \propto_{(2)} f_{X_u^* | Y^{t-2} D^{t-2} X^{t-1} X_k^*}(X_u^*; Y^{t-2}, D^{t-2}, X^{t-1}, X_k^*) f_{Y_{t-1}(D_{t-1}) | X_{t-1} X^*}(Y_{t-1}; X_{t-1}, X^*) \\
& \propto_{(3)} \exp\left(-\frac{1}{2}(X_u^* - E_t)^\top \Sigma_t^{-1}(X_u^* - E_t)\right) \phi\left(\frac{Y_t - X_t^\top \beta_t(D_t) - X_k^* \lambda_{kt}(D_t) - (X_u^*)^\top \lambda_{ut}(D_t)}{\sigma_t(D_t)}\right) \\
& \propto \exp\left(-\frac{1}{2}(X_u^* - E_t)^\top \Sigma_t^{-1}(X_u^* - E_t)\right) \\
& \quad \times \exp\left(-\frac{1}{2}(X_u^* - \lambda_{ut}(D_t)(\lambda_{ut}(D_t)^\top \lambda_{ut}(D_t))^{-1}(Y_t - X_t^\top \beta_t(D_t) - X_k^* \lambda_{kt}(D_t)))^\top\right. \\
& \quad \times \left.\frac{\lambda_{ut}(D_t) \lambda_{ut}(D_t)^\top}{\sigma_t^2(D_t)}(X_u^* - \lambda_{ut}(D_t)(\lambda_{ut}(D_t)^\top \lambda_{ut}(D_t))^{-1}(Y_t - X_t^\top \beta_t(D_t) - X_k^* \lambda_{kt}(D_t)))\right) \\
& =_{(4)} \exp\left(-\frac{1}{2}(X_u^* - E_{t+1})^\top \Sigma_{t+1}^{-1}(X_u^* - E_{t+1})\right).
\end{aligned}$$

Display (1) follows from Bayes' theorem. Display (2) holds since Assumption KL1 has the following three implications: first  $X_t \perp\!\!\!\perp X^* \mid (Y^{t-1}, D^{t-1}, X^{t-1})$ ; second  $\epsilon_{t-1}(d_{t-1}) \perp\!\!\!\perp (Y^{t-2}, D^{t-1}, X^{t-1}, X^*) \Rightarrow \epsilon_{t-1}(d_{t-1}) \perp\!\!\!\perp (Y^{t-2}, D^{t-1}, X^{t-2}) \mid (X_{t-1}, X^*) \Rightarrow Y_{t-1}(d_{t-1}) \perp\!\!\!\perp (Y^{t-2}, D^{t-1}, X^{t-2}) \mid (X_{t-1}, X^*)$ ; third  $D_{t-1} \perp\!\!\!\perp X_u^* \mid$

$(Y^{t-2}, D^{t-2}, X^{t-1}, X_k^*)$ . Display (3) holds from the induction assumption and Assumptions **KL1** and **KL2**. Display (4) follows from the definitions in Lemma 1.  $\square$

The following corollary of Lemma 1 will be used in the proof to Theorem 1.

**Corollary 2.** Let Assumptions **KL1** and **KL2** hold. Then  $Y_t$  conditional upon  $(Y^{t-1}, D^t, X^t, X_k^*)$  is distributed

$$N(X_t^\top \beta_t(D_t) + X_k^* \lambda_{kt}(D_t) + E_t^\top \lambda_{ut}(D_t), \lambda_{ut}(D_t)^\top \Sigma_t \lambda_{ut}(D_t) + \sigma_t^2(D_t))$$

*Proof of Corollary 2.* For  $t > 1$ ,

$$\begin{aligned} & f_{Y_t | Y^{t-1} D^t X^t X_k^*}(Y_t; Y^{t-1}, D^t, X^t, X_k^*) \\ &= \int f_{Y_t(D_t) | Y^{t-1} D^t X^t X_k^*}(Y_t; Y^{t-1}, D^t, X^t, X_k^*) f_{X_u^* | Y^{t-1} D^t X^t X_k^*}(X_u^*; Y^{t-1}, D^t, X^t, X_k^*) dX_u^* \\ &=_{(1)} \int f_{Y_t(D_t) | X_t X^*}(Y_t; X_t, X^*) f_{X_u^* | Y^{t-1} D^{t-1} X^t X_k^*}(X_u^*; Y^{t-1}, D^{t-1}, X^t, X_k^*) dX_u^* \\ &\propto_{(2)} \int \phi\left(\frac{Y_t - X_t^\top \beta_t(D_t) - X_k^* \lambda_{kt}(D_t) - (X_u^*)^\top \lambda_{ut}(D_t)}{\sigma_t(D_t)}\right) \\ &\quad \times \exp\left((X_u^* - E_t)^\top \Sigma_t^{-1} (X_u^* - E_t)\right) dX_u^* \\ &= \phi\left(\frac{Y_t - X_t^\top \beta_t(D_t) - X_k^* \lambda_{kt}(D_t) - E_t^\top \lambda_{ut}(D_t)}{\sqrt{\lambda_{ut}^\top(D_t) \Sigma_t \lambda_{ut}(D_t) + \sigma_t^2(D_t)}}\right) \end{aligned}$$

Equality (1) holds because Assumption **KL1** implies  $Y_t(D_t) \perp\!\!\!\perp (X^{t-1}, D^t, Y^{t-1}) \mid (X_t, X^*)$  and  $D_t \perp\!\!\!\perp X_u^* \mid (X_k^*, X^t, Y^{t-1}, D^{t-1})$ . Equality (2) holds because Assumption **KL1** and **KL2** imply Lemma 1 and  $\epsilon_t(d) \mid (X_t, X^*) \sim N(0, \sigma_t^2(D))$ . A similar argument applies for  $t = 1$ .  $\square$

*Proof of Theorem 1.* The proof is in four parts. First, we use Corollary 2 and Bruni and Koch (1985) to identify the distribution of  $Y_t \mid (D^t, Y^{t-1}, X^t, X_k^*)$  up to an affine transformation of  $X_k^*$ . Second, we show that the support of  $X_k^*$  is identified. The third part is to use the normalization (Assumption **KL3**) to show that the affine transformation is the identity function. Finally, we use identification of the distribution of

$(Y^t, D^t, X^t, X_k^*)$  to identify the distribution of  $(Y^t, D^t, X^t, X^*)$ .

*Part 1.* By Corollary 2,  $f_{Y_t|Y^{t-1}D^tX^t}(y_t; y^{t-1}, d^t, x^t) =$

$$\int f_{Y_t|Y^{t-1}D^tX^tX_k^*}(y_t; y^{t-1}, d^t, x^t, x_k^*) df_{X_k^*|Y^{t-1}D^tX^t}(x_k^*; y^{t-1}, d^t, x^t) dx_k^*.$$

I.e., the nonparametrically identified  $f_{Y_t|Y^{t-1}D^tX^t}(y_t; y^{t-1}, d^t, x^t)$  is a mixture of Gaussians. To identify the component and mixture distributions, we will apply Bruni and Koch (1985, Theorem 3). First, for any  $t$ ,  $d^t$  and  $x_1$ , define

$$\Lambda = \{x_k^* \mapsto (x_t^\top \beta(d_t) + x_k^* \mu_1(\theta^t) + \mu_2(\theta^t), \sigma(\theta^t)) : \theta^t \in \Theta^t\},$$

where  $\theta^t = (\beta^t, \lambda_k^t, \lambda_u^t, \sigma^t, \Sigma_u)$ ,  $\Theta^t$  is the corresponding subset of  $\Theta$ , and

$$\begin{aligned} \mu_1(\theta^t) &= \left( \lambda_{kt}(d_t) - \lambda_{ut}^\top(d_t) \Sigma_t \sum_{s=1}^{t-1} \lambda_{us}(d_s) \frac{\lambda_{ks}(d_s)}{\sigma_s^2(d_s)} \right), \\ \mu_2(\theta^t) &= \lambda_{ut}(d_t)^\top \Sigma_t \sum_{s=1}^{t-1} \lambda_{us}(d_s) \frac{y_{is} - x_{is}^\top \beta(d_s)}{\sigma_s^2(d_s)}, \\ \sigma(\theta^t) &= \lambda_{ut}(d_t)^\top \Sigma_t \lambda_{ut}(d_t) + \sigma_t^2(d_t). \end{aligned}$$

For example, for  $t = 1$ ,  $\sigma(\theta^1) = \lambda_{u1}(d_1)^\top \Sigma_u(x_1) \lambda_{u1}(d_1) + \sigma_1^2(d_1)$ . Notice that  $X_k^* \lambda_{kt}(d_t) + E_t^\top \lambda_{ut}(d_t) = X_k^* \mu_1(\theta^t) + \mu_2(\theta^t)$ . Under Assumptions KL4(A,B,C) and KL5(C), 4,  $\Lambda \subset \Lambda_4$  where  $\Lambda_4$  is defined in Bruni and Koch (1985, p. 1344). Thus Bruni and Koch (1985, Theorem 3) applies and

$$((x_t^\top \beta(d_t) + \pi(v_k) \mu_1(\theta^t) + \mu_2(\theta^t), \sigma(\theta^t), df_{X_k^*|Y^{t-1}D^tX^t}(\pi(v_k); y^{t-1}, d^t, x^t)) \quad (7)$$

is identified with  $\pi$  an unknown non-constant affine function which may depend on the history  $(y^{t-1}, d^t, x^t)$ .

*Part 2.* In this part, we show  $\pi$  is identity for the normalized choice  $D_1 = d_1$ , which provides identification of  $\mathcal{S}(X_k^*)$  by Assumption KL4 (D). In this part, it will be

useful to denote  $\beta_t(d) = (\alpha_t(d), \gamma_t(d)^\top)^\top$ , where  $\alpha_t(d)$  is the coefficient on the constant term in  $X_t$ .

Let  $t = 1$  and  $d_1$  as in Assumption **KL3**(A), then since  $\mu_1(\theta^1) = \lambda_{k1}(d_1)$  and  $\mu_2(\theta^1) = 0$ , from Part 1 we have identified:

$$(x_1^\top \gamma(d_1) + \pi(v_k), \sigma(\theta^1), df_{X_k^* | D^1 X^1}(\pi(v_k); d^1, x^1)),$$

with  $\pi(v_k) = \pi_0 + \pi_1 v_k$ . Since  $\lambda_{k1}(d_1) = 1$ ,  $\pi_1 = 1$ . We now show  $\pi_0 = 0$ . First notice that  $\pi_0$  does not depend on  $(d_1, x_1)$  since the support of  $X_k^* | (D_1 = d_1, X_1 = x_1)$  is the same for each  $(d_1, x_1)$ . Now suppose that for any  $x_1$ ,  $x_1^\top \gamma(d_1) + \pi_0 = x_1^\top \tilde{\gamma}(d_1) + \tilde{\pi}_0$ . In particular for  $\tilde{x}_1 \neq x_1$ ,  $(x_1 - \tilde{x}_1)^\top (\gamma(d_1) - \tilde{\gamma}(d_1)) = 0$ . By Assumption **KL4**(E), we conclude  $\gamma(d_1) - \tilde{\gamma}(d_1) = 0$ . This in conjunction with  $\alpha_1(d_1) = 0$  gives  $\pi_0 = \tilde{\pi}_0 = 0$ .

To conclude this part, we show that if  $\pi$  is identity for each history  $(d^s, y^{s-1}, x^s)$   $s = 1, 2, \dots, t$ , then  $f_{Y^t D^t X^t X_k^*}(y^t, d^t, x^t, v_k)$  is point identified (i.e.,  $(\alpha_t(d_t), \gamma_t(d_t), \mu_1(\theta^t), \mu_2(\theta^t), \sigma(\theta^t), df_{X_k^* | Y^{t-1} D^t X^t}(v_k; y^{t-1}, d^t, x^t))$  is point identified). For  $t = 1$ , as  $(\mu_1(\theta^1), \mu_2(\theta^1)) = (\lambda_{k1}(d_1), 0)$  identification follows immediately from display (7) and Assumption **KL4**(E).

Now suppose  $(\alpha_s(d_s), \gamma_s(d_s), \mu_1(\theta^s), \mu_2(\theta^s), \sigma(\theta^s), df_{X_k^* | Y^{s-1} D^s X^s}(v_k; y^{s-1}, d^s, x^s))$  is point identified for each  $s < t$ . From equation (7),

$$((\alpha_t(d_t) + x_t^\top \gamma(d_t) + v_k \mu_1(\theta^t) + \mu_2(\theta^t), \sigma(\theta^t), df_{X_k^* | Y^{t-1} D^t X^t}(v_k; y^{t-1}, d^t, x^t))$$

is identified for every  $(d^t, y^{t-1}, x^t)$ .  $\mu_1(\theta^t)$  is identified from variation in  $v_k$ . Assumption **KL4**(E) implies identification of

$$(\alpha_t(d_t) + \mu_2(\theta^t), \gamma(d_t)).$$

Then  $\mu_2(\theta^t) = \sum_{s=1}^{t-1} (y_s - \alpha_s(d_s) - x_s^\top \gamma_s(d_s)) \frac{\partial}{\partial y_s} (\alpha_t(d_t) + \mu_2(\theta^t))$ , from which follows identification of  $\alpha_t(d_t)$ .

*Part 3.* In this part we show that  $\pi$  is the identity function. First, we use knowledge of  $\mathcal{S}(X_k^*)$  to prove the affine function must satisfy  $|\frac{\partial}{\partial v}\pi(v)| = 1$  for any history  $(y^{t-1}, d^t, x^t)$ . Then we use restrictions on the panel dimension to conclude  $\pi$  is identity for each history  $(y^{t-1}, d^t, x^t)$ .

Second, for each fixed  $(y^{t-1}, d^t, x^t)$ ,  $df_{X_k^*|Y^{t-1}D^tX^t}(\pi(v_k); y^{t-1}, d^t, x^t)$  is identified from part 1. Then, by Assumption **KL4**(D),

$$\mathcal{S}(X_k^*) = df_{X_k^*|Y^{t-1}D^tX^t}^{-1}[\mathbb{R}_+] = (df_{X_k^*|Y^{t-1}D^tX^t} \circ \pi)^{-1}[\mathbb{R}_+],$$

where  $R_+ = \{x \in \mathbb{R} : x > 0\}$ . And since  $\pi$  is bijective,

$$(\pi \circ df_{X_k^*|Y^{t-1}D^tX^t}^{-1})[\mathbb{R}_+] = df_{X_k^*|Y^{t-1}D^tX^t}^{-1}[\mathbb{R}_+].$$

In particular

$$\begin{aligned} \pi(\sup df_{X_k^*|Y^{t-1}D^tX^t}^{-1}[\mathbb{R}_+]) &= \sup df_{X_k^*|Y^{t-1}D^tX^t}^{-1}[\mathbb{R}_+], \\ \pi(\inf df_{X_k^*|Y^{t-1}D^tX^t}^{-1}[\mathbb{R}_+]) &= \inf df_{X_k^*|Y^{t-1}D^tX^t}^{-1}[\mathbb{R}_+]. \end{aligned}$$

The only affine functions that satisfy these identities are  $\pi^+(v) = v$  and  $\pi^-(v) = (\bar{v} + \underline{v}) - v$  for  $\underline{v} = \inf df_{X_k^*|Y^{t-1}D^tX^t}^{-1}[\mathbb{R}_+]$  and  $\bar{v} = \sup df_{X_k^*|Y^{t-1}D^tX^t}^{-1}[\mathbb{R}_+]$ .

Third, it remains to show that  $\pi = \pi^+$ . To do so, it will be useful to define:

$$\tilde{\mu}_{ts}(d^{t-1}) = \Sigma_t \frac{\lambda_{us}(d_s)}{\sigma_s^2(d_s)}$$

It will also be useful to denote  $\mu_j(d^t) = \mu_j(\theta^t)$ , to emphasize the dependence of  $\mu_j$  on  $d^t$ . Then notice  $\mu_1(d^t) = \lambda_{kt}(d_t) - \lambda_{ut}(d_t) \sum_{s=1}^{t-1} \tilde{\mu}_{ts}(d^{t-1}) \lambda_{ks}(d_s)$  and  $\mu_2(d^t) = \lambda_{ut}(d_t)^\top \sum_{s=1}^{t-1} \tilde{\mu}_{ts}(d^{t-1}) (Y_{is} - X_{is}^\top \beta_s(d_s))$ .

The proof is inductive. First consider  $t = 1$ . From Assumption **KL3**(A),  $\lambda_{k1}(d_1) = 1$ . For  $\tilde{d}_1 \neq d_1$ , given part 1 and  $\frac{\partial}{\partial x}|\pi(x)| = 1$ ,  $\lambda_{k1}(\tilde{d}_1)$  is iden-

tified up to sign as  $\frac{\partial}{\partial v_k} \left( x_1^\top \beta(\tilde{d}_1) + \lambda_{k1}(\tilde{d}_1) \pi(v_k) \right)$ . Similarly, for  $d^2 = (d_2, d_1)$ ,  $\mu_1(d^2) = \lambda_{k2}(d_2) - \lambda_{u2}(d_2)^\top \tilde{\mu}_{21}(d^1) \lambda_{k1}(d_1)$  is identified up to sign and  $\lambda_{u2}(d_2)^\top \tilde{\mu}_{21}(d^1)$  are identified as  $\frac{\partial}{\partial x} (\beta(d_2)' x_2 + \pi(v_k) \mu_1(d^2) + \mu_2(d^2))$  for  $x = v_k$  and  $x = y_1$ , respectively (since  $\mu_2(d^2) = \lambda_{u2}(d_2)^\top \tilde{\mu}_{21}(d^1) (y_1 - x_1^\top \beta_1(d_1))$  and  $\mu_1(d^2)$  does not depend on  $y_1$ ). Repeating this argument for the choice sequence  $(\tilde{d}_1, d_2)$  yields identification of  $(\lambda_{k2}(d_2) - \lambda_{u2}(d_2)^\top \tilde{\mu}_{21}(\tilde{d}^1) \lambda_{k1}(\tilde{d}_1))$  up to sign and  $\lambda_{u2}(d_2)^\top \tilde{\mu}_{21}(\tilde{d}^1)$ .

Summarizing, we have identification of  $\lambda_{u2}(d_2)^\top \tilde{\mu}_{21}(d^1)$ ,  $\lambda_{u2}(d_2)^\top \tilde{\mu}_{21}(\tilde{d}^1)$ , and  $(-1)^{j_1} \lambda_{k1}(\tilde{d}_1)$ ,  $(-1)^{j_{d_2}} (\lambda_{k2}(d_2) - \lambda_{u2}(d_2)^\top \tilde{\mu}_{21}(d^1))$ , and  $(-1)^{\tilde{j}_{d_2}} (\lambda_{k2}(d_2) - \lambda_{u2}(d_2)^\top \tilde{\mu}_{21}(\tilde{d}^1) \lambda_{k1}(\tilde{d}_1))$  with  $(j_1, \tilde{j}_{d_2}, j_{d_2}) \in \{0, 1\}^3$ . We show only the correct choice of sign will satisfy

$$\begin{aligned} & \lambda_{u2}(d_2)^\top \tilde{\mu}_{21}(\tilde{d}^1) (-1)^{j_1} \lambda_{k1}(\tilde{d}_1) + (-1)^{\tilde{j}_{d_2}} (\lambda_{k2}(d_2) - \lambda_{u2}(d_2)^\top \tilde{\mu}_{21}(\tilde{d}^1) \lambda_{k1}(\tilde{d}_1)) \\ &= \lambda_{u2}(d_2)^\top \tilde{\mu}_{21}(d^1) \lambda_{k1}(d_1) + (-1)^{j_{d_2}} (\lambda_{k2}(d_2) - \lambda_{u2}(d_2)^\top \tilde{\mu}_{21}(d^1) \lambda_{k1}(d_1)). \end{aligned}$$

Suppose  $j_{d_2} = 0$ . It is straightforward to show the following implications:

$$\begin{aligned} (j_1, \tilde{j}_{d_2}) = (1, 1) &\implies \lambda_{k2}(d_2) = 0, \\ (j_1, \tilde{j}_{d_2}) = (0, 1) &\implies \lambda_{k2}(d_2) - \lambda_{u2}(d_2)^\top \tilde{\mu}_{21}(\tilde{d}^1) \lambda_{k1}(\tilde{d}_1) = 0, \\ (j_1, \tilde{j}_{d_2}) = (1, 0) &\implies \lambda_{u2}(d_2)^\top \tilde{\mu}_{21}(\tilde{d}^1) \lambda_{k1}(\tilde{d}_1) = 0. \end{aligned}$$

The three implications contradict Assumptions [KL5](#) (B), (C) and (D), respectively.

Now suppose  $j_{d_2} = 1$ , then

$$\begin{aligned} (j_1, \tilde{j}_{d_2}) = (0, 0) &\implies \lambda_{k2}(d_2) - \lambda_{u2}(d_2)^\top \tilde{\mu}_{21}(d^1) \lambda_{k1}(d_1) = 0, \\ (j_1, \tilde{j}_{d_2}) = (1, 1) &\implies \lambda_{u2}(d_2)^\top \tilde{\mu}_{21}(d^1) \lambda_{k1}(d_1) = 0, \\ (j_1, \tilde{j}_{d_2}) = (0, 1) &\implies \lambda_{u2}(d_2)^\top \tilde{\mu}_{21}(\tilde{d}^1) \lambda_{k1}(\tilde{d}_1) - \lambda_{u2}(d_2)^\top \tilde{\mu}_{21}(d^1) \lambda_{k1}(d_1) = 0, \\ (j_1, \tilde{j}_{d_2}) = (1, 0) &\implies \lambda_{k2}(d_2) - \lambda_{u2}(d_2)^\top \tilde{\mu}_{21}(\tilde{d}^1) \lambda_{k1}(\tilde{d}_1) - \lambda_{u2}(d_2)^\top \tilde{\mu}_{21}(d^1) \lambda_{k1}(d_1) = 0. \end{aligned}$$



The first three implications contradict Assumptions **KL5** (C), (D) and (A), respectively. Finally, we show that the final equality contradicts Assumption **KL5** (E). For each  $d \in \{d_{2,i} \in \mathcal{S}(D_2) : i = 1, 2, \dots, p\} \cup \{\tilde{d}_{2,i} \in \mathcal{S}(D_2) : i = 1, 2, \dots, p\}$  of Assumption **KL5** (E), by considering the sequences  $(d_1, d)$ ,  $(\tilde{d}_1, d)$ ,  $(-1)^{j_d}(\lambda_{k2}(d) - \lambda_{u2}(d)^\top \tilde{\mu}_{21}(d^1) \lambda_{k1}(d_1))$  and  $(-1)^{\tilde{j}_d}(\lambda_{k2}(d) - \lambda_{u2}(d)^\top \tilde{\mu}_{21}(\tilde{d}^1) \lambda_{k1}(\tilde{d}_1))$  is identified with  $(j_d, \tilde{j}_d) \in \{(1, 0), (0, 0)\}$ . Since  $\lambda_{k1}(\tilde{d}_1) \neq 0$  by Assumption **KL5** (B), for the sign of  $\lambda_{k1}(\tilde{d}_1)$  to be constant across sequences, we can rule out all signs except  $(j_1, (j_{d_{2,i}}, \tilde{j}_{d_{2,i}}, j_{\tilde{d}_{2,i}}, \tilde{j}_{\tilde{d}_{2,i}} : i = 1, \dots, p)) \in \{(0, (0, 0, 0, 0)^p), (1, (1, 0, 1, 0)^p)\}$ . If  $(j_1, (j_{d_{2,i}}, \tilde{j}_{d_{2,i}}, j_{\tilde{d}_{2,i}}, \tilde{j}_{\tilde{d}_{2,i}} : i = 1, \dots, p)) = (1, (1, 0, 1, 0)^p)$ , then

$$\begin{aligned} 0 &= \text{vec}(\lambda_{k2}(d_{2,1}), \dots, \lambda_{k2}(d_{2,k})) - (\lambda_{u2}(d_{2,1}) \dots F_{u2}(d_{2,k}))^\top (\tilde{\mu}_{21}(\tilde{d}^1) \lambda_{k1}(\tilde{d}_1) + \tilde{\mu}_{21}(d^1) \lambda_{k1}(d_1)) \\ &= \text{vec}(\lambda_{k2}(\tilde{d}_{2,1}), \dots, \lambda_{k2}(\tilde{d}_{2,k})) - (\lambda_{u2}(\tilde{d}_{2,1}) \dots F_{u2}(\tilde{d}_{2,k}))^\top (\tilde{\mu}_{21}(\tilde{d}^1) \lambda_{k1}(\tilde{d}_1) + \tilde{\mu}_{21}(d^1) \lambda_{k1}(d_1)), \end{aligned}$$

which contradicts Assumption **KL5**(E).

For the induction step, suppose  $\pi$  is identity for each history  $(d^s, y^{s-1}, x^s)$   $s = 1, \dots, t-1$  and consider choice sequences  $d^{t-1} = ((d^{t-2})^\top, d_{t-1})^\top$  and  $\tilde{d}^{t-1} = ((d^{t-2})^\top, \tilde{d}_{t-1})^\top$  for  $d_{t-1} \neq \tilde{d}_{t-1}$ . From part 1,  $(\beta_t(d_t)' x_t + \mu_1(d^{t-2}, d, d_t) \pi(v_k) + \mu_2(d^{t-2}, d, d_t))$  is identified for  $d = d_{t-1}, \tilde{d}_{t-1}$  and  $\pi \in \{\pi^+, \pi^-\}$ . By the preceding arguments,  $(\lambda_{ut}(d_t) \sum_{s=1}^{t-1} \tilde{\mu}_{ts}(d^{t-1}) \lambda_{ks}(d_s))$ ,  $\lambda_{ut}(d_t) \sum_{s=1}^{t-1} \tilde{\mu}_{ts}(\tilde{d}^{t-1}) \lambda_{ks}(\tilde{d}_s)$ ,  $(-1)^{j_1} (\lambda_{kt}(d_t) - \lambda_{ut}(d_t)^\top \sum_{s=1}^{t-1} \tilde{\mu}_{ts}(d^{t-1}) \lambda_{ks}(d_s))$ , and  $(-1)^{j_2} (\lambda_{kt}(d_t) - \lambda_{ut}(d_t)^\top \sum_{s=1}^{t-1} \tilde{\mu}_{ts}(\tilde{d}^{t-1}) \lambda_{ks}(\tilde{d}_s))$  is identified with  $(j_1, j_2) \in \{0, 1\}^2$

As before we show that only that only  $(j_1, j_2) = (0, 0)$  is consistent with the identity

$$\begin{aligned} &(-1)^{j_1} \left( \lambda_{kt}(d_t) - \lambda_{ut}(d_t)^\top \sum_{s=1}^{t-1} \tilde{\mu}_{ts}(d^{t-1}) \lambda_{ks}(d_s) \right) + \lambda_{ut}(d_t) \sum_{s=1}^{t-1} \tilde{\mu}_{ts}(d^{t-1}) \lambda_{ks}(d_s) \\ &= (-1)^{j_2} \left( \lambda_{kt}(d_t) - \lambda_{ut}(d_t)^\top \sum_{s=1}^{t-1} \tilde{\mu}_{ts}(\tilde{d}^{t-1}) \lambda_{ks}(\tilde{d}_s) \right) + \lambda_{ut}(d_t) \sum_{s=1}^{t-1} \tilde{\mu}_{ts}(\tilde{d}^{t-1}) \lambda_{ks}(\tilde{d}_s) \end{aligned}$$

For this, consider

$$\begin{aligned}
(j_1, j_2) = (0, 1) &\implies \left( \lambda_{kt}(d_t) - \lambda_{ut}(d_t) \sum_{s=1}^{t-1} \tilde{\mu}_{ts}(\tilde{d}^{t-1}) \lambda_{ks}(\tilde{d}_s) \right) = 0, \\
(j_1, j_2) = (1, 0) &\implies \left( \lambda_{kt}(d_t) - \lambda_{ut}(d_t) \sum_{s=1}^{t-1} \tilde{\mu}_{ts}(d^{t-1}) \lambda_{ks}(d_s) \right) = 0, \\
(j_1, j_2) = (1, 1) &\implies \lambda_{ut}(d_t) \sum_{s=1}^{t-1} \tilde{\mu}_{ts}(d^{t-1}) \lambda_{ks}(d_s) - \lambda_{ut}(d_t) \sum_{s=1}^{t-1} \tilde{\mu}_{ts}(\tilde{d}^{t-1}) \lambda_{ks}(\tilde{d}_s) = 0,
\end{aligned}$$

which contradict Assumptions **KL5** (C), (C) and (A), respectively. Thus  $\pi$  is the identity function for the history  $(d^t, y^{t-1}, x^t)$ .

*Part 4.* From parts 1 to 3,  $f_{Y^T D^T X^T X_k^*}$ , and thus  $h_t$ , is identified. First,

$$\begin{aligned}
&f_{Y^T D^T X^T X_k^*}(y^T, d^T, x^T, v_k) \\
&= \int f_{Y^T(d^T) D^T X^T X^*}(y^T, d^T, x^T, v) dv_u \\
&= \int f_{Y_T(d_T)|X_T, X^*}(y_T; x_T, v) f_{D_T|Y^{T-1} D^{T-1} X^T X_k^*}(d_T; y^{T-1}, d^{T-1}, x^T, v_k) \\
&\quad \times f_{X_T|Y^{T-1} D^{T-1} X^{T-1}}(x_T; y^{T-1}, d^{T-1}, x^{T-1}) \dots f_{Y_1(d_1)|X_1 X^*}(y_1; x_1, v) \\
&\quad \times f_{D_1|X_1 X_k^*}(d_1; x_1, v_k) f_{X_u^*|X_1 X_k^*}(v_u; x_1, v_k) f_{X_1 X_k^*}(x_1, v_k) dv_u.
\end{aligned}$$

This implies that on the support of  $f_{Y^T D^T X^T X_k^*}$ ,

$$\begin{aligned}
&\frac{f_{Y^T D^T X^T X_k^*}(y^T, d^T, x^T, v_k)}{f_{D_1 X_1 X_k^*}(d_1, x_1, v_k) \prod_{t=2}^T f_{D_t X_t|Y^{t-1} D^{t-1} X^{t-1} X_k^*}(d_t; y^{t-1}, d^{t-1}, x^t, v_k)} \\
&= \int \prod_{t=1}^T f_{Y_t(d_t)|X_t X^*}(y_t; x_t, v) f_{X_u^*|X_k^* X_1}(v_u; v_k, x_1) dv_u.
\end{aligned}$$

I.e., the function is equal to the probability density function of a jointly normal random variable with mean

$$(x_t^\top \beta_t(d_t) + v_k \lambda_{kt}(d_t))_{t=1}^T,$$

and covariance matrix

$$\lambda_u(d)^\top \Sigma_u(x_1) \lambda_u(d) + \text{diag} \left( \sigma_t^2(d_t) : t = 1, \dots, T \right),$$

where  $\lambda_u(d) = (\lambda_{u1}(d_1) \lambda_{u2}(d_2) \dots F_{uT}(d_T))$ . From parts 1 and 2, the components of the mean function are identified. The components of the covariance matrix are identified under Assumptions **KL3** (B) and **KL5** (F).  $\square$

## A.2 Proofs for Section 3.3

In this section denote  $\mathcal{L} = \{m : \mathbb{R}^k \rightarrow \mathbb{R} : \sup_{a \in \mathbb{R}^k} |m(a)| < \infty, \int |m(a)| da < \infty\}$  and  $\mathcal{L}_A = \{m : \mathbb{R}^k \rightarrow \mathbb{R} : \sup_{a \in \mathbb{R}^k} |m(a)| < \infty, \int |m(a)| f_A(a) da < \infty\}$  for a random variable  $A$  with p.d.f.  $f_A$ .

*Proof.* Let  $x \in \mathcal{S}(X)$  be given and fix a choice sequence  $d = (d_1, d_2, \dots, d_T)$  whose first  $p$  elements satisfy Assumption **L3**, and define  $W_1 = (Y_1, \dots, Y_p)$ ,  $W_2 = Y_{p+1}$  and  $W_3 = (Y_{p+2}, \dots, Y_T)$ . Let  $L_{123} : \mathcal{L}_{W_3} \rightarrow \mathcal{L}$  and  $L_{13} : \mathcal{L}_{W_3} \rightarrow \mathcal{L}$  be defined as  $[L_{123}m](w_1) =$

$$\int \frac{f_{YDX}(y, d, x)}{f_{D_1X_1}(d_1, x_1) \prod_{t=2}^T f_{D_tX_t|Y^{t-1}D^{t-1}X^{t-1}}(d_t, x_t; y^{t-1}, d^{t-1}, x^{t-1})} m(w_3) dw_3,$$

and  $[L_{13}m](w_1) = \int [L_{123}m](w_1) dw_2$ . In addition, define

$$\begin{aligned} L_{1X^*} : \mathcal{L} &\rightarrow \mathcal{L} & [L_{1X^*}m](w_1) &= \int \prod_{t=1}^p f_{Y_t(d_t)|X_tX^*}(y_t; x_t, v) m(v) dv, \\ L_{X^*3} : \mathcal{L}_{W_3} &\rightarrow \mathcal{L} & [L_{X^*3}m](v) &= \int \prod_{t=p+2}^T f_{Y_t(d_t)|X_tX^*}(y_t; x_t, v) f_{X^*|X_1}(v; x_1) m(w_1) dw_1, \\ D_{X^*} : \mathcal{L}_{X^*} &\rightarrow \mathcal{L}_{X^*} & [D_{X^*}m](v) &= f_{Y_{p+1}(d_{p+1})|X_{p+1}X^*}(y_{p+1}; x_{p+1}, v) m(v). \end{aligned}$$

The following derivation shows  $L_{123} = L_{1X^*}D_{X^*}L_{X^*3}$ . First,

$$\begin{aligned} f_{YDX}(y, d, x) &= \int f_{YDXX^*}(y, d, x, v)dv \\ &= \int f_{Y_T(d_T)|X_TX^*}(y_T; x_T, v) f_{D_TX_T|Y^{T-1}D^{T-1}X^{T-1}}(d_T, x_T; y^{T-1}, d^{T-1}, x^{T-1}) \\ &\quad \times f_{Y_{T-1}(d_{T-1})|X_{T-1}X^*}(y_{T-1}; x_{T-1}, v) \dots f_{D_1X_1}(d_1, x_1) f_{X^*|X_1}(v; x_1)dv. \end{aligned}$$

Then, by Assumption **L4** (A),

$$\begin{aligned} &\frac{f_{YDX}(y, d, x)}{f_{D_1X_1}(d_1, x_1) \prod_{t=2}^T f_{D_tX_t|Y^{t-1}D^{t-1}X^{t-1}}(d_t, x_t; y^{t-1}, d^{t-1}, x^{t-1})} \\ &= \int \prod_{t=1}^T f_{Y_t(d_t)|X_tX^*}(y_t; x_t, v) f_{X^*|X_1}(v; x_1)dv, \end{aligned}$$

and therefore that

$$\begin{aligned} [L_{123}m](w_1) &= \int \left( \int \prod_{t=1}^T f_{Y_t(d_t)|X_tX^*}(y_t; x_t, v) f_{X^*|X_t}(v; x_t)dv \right) m(w_3)dw_3 \\ &= \int \prod_{t=1}^{p+1} f_{Y_t(d_t)|X_tX^*}(y_t; x_t, v) \left( \int \prod_{t=p+2}^T f_{Y_t(d_t)|X_tX^*}(y_t; x_t, v) f_{X^*|X_t}(v) m(w_3)dw_3 \right) dv \\ &= \int \prod_{t=1}^p f_{Y_t(d_t)|X_tX^*}(y_t; x_t, v) (f_{Y_{p+1}(d_{p+1})|X_{p+1}X^*}(y_{p+1}; x_{p+1}, v) [L_{X^*3}m](v)) dv \\ &= \int \int \prod_{t=1}^p f_{Y_t(d_t)|X_tX^*}(y_t; x_t, v) [D_{X^*}L_{X^*3}m](v)dv \\ &= [L_{1X^*}D_{X^*}L_{X^*3}m](w_1), \end{aligned}$$

and  $L_{123} = L_{1X^*}D_{X^*}L_{X^*3}$ . Similarly,  $L_{13} = L_{1X^*}L_{X^*3}$ .

From here, Assumptions **L1**, **L2**, **L3**, **L4** (B), and **L5** imply the arguments of Theorem 1 Freyberger (2018) apply<sup>9</sup>, so that  $\lambda_t(d_t)$ ,  $f_{Y_t(d_t)|X_tX^*}(\cdot; x_t, \cdot)$  and  $f_{X^*|X_1}(\cdot; x_1)$

---

<sup>9</sup>The listed assumptions imply the assumptions of Freyberger (2018, Theorem 1) with the primary exception of Assumption **L1** that differs from Assumption N5 in Freyberger (2018) by allowing period  $t$  variables to impact the evolution of period  $t'$  covariates for  $t' > t$ . However, since Assumption **L1** implies  $f_{Y_t(d_t)|X_tX^*}(y; x, v) = f_{\epsilon_t(d_t)}(y - \beta_t(d_t)^\top x - \lambda_t^\top v)$ , Freyberger (2018, Lemma 1) and

are identified for each  $t$  for the given  $(d_t, x)$ . Given identification of  $f_{Y_t(d_t)|X_tX^*}(\cdot; x_t, \cdot)$  for each  $x_t \in \mathcal{S}(X_t)$  and  $\lambda_t(d_t)$ , Assumption [L4](#) (C) implies identification of  $\beta_t(d_t)$  and thus  $f_{\epsilon_t(d_t)}$ .

Next, given an arbitrary  $t$  and  $d_t$ , define  $\tilde{d}$  by replacing the  $t$ th element of  $d$  with  $d_t$ . Then let  $\rho$  be a permutation  $(1, 2, \dots, T) \mapsto (t_1, t_2, \dots, t_T)$  such that  $t \mapsto t_1$  and define  $\tilde{W}_1 = (Y_{t_1}, Y_{t_2}, \dots, Y_{t_p})$ ,  $\tilde{W}_2 = (Y_{t_{p+1}}, Y_{t_{p+1}}, \dots, Y_{t_T})$ ,

$$\begin{aligned} \tilde{L}_{2X^*} : \mathcal{L} &\rightarrow \mathcal{L} & [\tilde{L}_{2X^*}m](\tilde{w}_2) &= \int \prod_{i=p+1}^T f_{Y_{t_i}(d_{t_i})|X_{t_i}X^*}(y_{t_i}; x_{t_i}, v) f_{X^*|X_1}(v; x_1) m(v) dv, \\ \tilde{L}_{X^*1} : \mathcal{L}_{\tilde{W}_1} &\rightarrow \mathcal{L} & [\tilde{L}_{X^*1}m](v) &= \int \prod_{i=1}^p f_{Y_{t_i}(d_{t_i})|X_{t_i}X^*}(y_{t_i}; x_{t_i}, v) m(\tilde{w}_1) d\tilde{w}_1, \end{aligned}$$

and  $\tilde{L}_{21} : \mathcal{L}_{\tilde{W}_1} \rightarrow \mathcal{L}$  as

$$[\tilde{L}_{21}m](\tilde{w}_2) = \int \frac{f_{YDX}(y, d, x)}{f_{D_1X_1}(d_1, x_1) \prod_{t=2}^T f_{D_tX_t|Y^{t-1}D^{t-1}X^{t-1}}(d_t, x_t; y^{t-1}, d^{t-1}, x^{t-1})} m(\tilde{w}_1) d\tilde{w}_1.$$

As before,  $\tilde{L}_{21} = \tilde{L}_{2X^*}\tilde{L}_{X^*1}$ . Since  $\tilde{L}_{2X^*}$  and  $\tilde{L}_{21}$  are identified and injective,  $\tilde{L}_{X^*1}$  is identified by  $\tilde{L}_{2X^*}^{-1}\tilde{L}_{21} = \tilde{L}_{X^*1}$  and thus  $\beta_t(d_t), \lambda_t(d_t), f_{\epsilon(d_t)}$ .  $\square$

---

D'Haultfoeuille (2011) may be applied with small modification.

## Online appendix

### B Proof of Corollary 1

In this proof we denote  $\beta_t(d) = (\alpha_t(d), \gamma_t(d)^\top)^\top$ , where  $\alpha_t(d)$  is the coefficient on the constant term in  $X_t$ . Fix  $(d_1, d_2, \dots, d_p)$  as in the statement and define  $\lambda_u = (\lambda_{u1}(d_1)\lambda_{u2}(d_2)\dots F_{up}(d_p))$ ,  $\tilde{X}_u^* = \lambda_u^\top (X_u^* - \mu_u)$ ,  $\tilde{\epsilon}_t(d) = \epsilon_t(d) - c_t(d)$ ,  $\tilde{X}_k^* = b + \lambda_{k1}(d_1)X_k^*$  where  $b = \alpha_1(d_1) + \lambda_{u1}(d_1)^\top \mu_u + c_1(d_1)$ . Finally, define  $\tilde{F}_{kt}(d_t) = \lambda_{k1}(d_1)^{-1}\lambda_{kt}(d_t)$ ,  $\tilde{F}_{ut}(d_t) = \lambda_u^{-1}\lambda_{ut}(d_t)$ , and  $\tilde{\alpha}_t(d) = \alpha_t(d) - \tilde{F}_{kt}(d)b + \lambda_{ut}(d)^\top \mu_u + c_t(d)$ . We then have that

$$Y_t(d) = \tilde{\alpha}_t(d) + X_t^\top \gamma_t(d) + (\tilde{X}_u^*)^\top \tilde{F}_{ut}(d) + \tilde{X}_k^* \tilde{F}_{kt}(d) + \tilde{\epsilon}_t(d),$$

$E[\tilde{\epsilon}_t(d)] = 0$  and  $E[\tilde{X}_u^* \mid X_1 = x, X_k^* = v_k] = 0$  so that the reparameterized model satisfies Assumption **KL2** (with  $\Sigma_u(X_1) = \lambda_u^\top \tilde{\Sigma}_u(X_1)\lambda_u$ ). Also,  $\tilde{F}_{k1}(d_1) = 1$ ,  $\tilde{\alpha}_1(d_1) = 0$  and  $\tilde{F}_p := (\tilde{F}_{u1}(d_1)\tilde{F}_{u2}(d_2)\dots\tilde{F}_{up}(d_p)) = I_{p \times p}$  so the reparameterized model satisfies Assumption **KL3**. By Theorem 1,  $\tilde{\theta} = ((\tilde{\alpha}_t, \gamma_t, \tilde{F}_{kt}, \tilde{F}_{ut}, \sigma_t^2, \tilde{h}_t)_{t=1}^\top, \Sigma_u, F_{\tilde{X}_k^*})$  is identified, where  $\tilde{h}_t$  and  $F_{\tilde{X}_k^*}$  are the CCPs and conditional distribution of  $\tilde{X}_k^*$ , respectively. This, in turn, implies the identification of the distribution of  $C_{jt}$  for  $j = k, u$ . Finally,

$$\begin{aligned} & \tilde{\alpha}_t + x^\top \gamma_t + Q_\alpha[\tilde{C}_{kt} + \tilde{C}_{ut} + \tilde{\epsilon}] \\ &= \alpha_t - \tilde{F}_{kt}b + \lambda_{ut}\mu_u + c_t + x^\top \gamma_t + Q_\alpha[\tilde{C}_{kt} + \tilde{C}_{ut} + \tilde{\epsilon}] \\ &= \alpha_t - \tilde{F}_{kt}b + \lambda_{ut}\mu_u + c_t + x^\top \gamma_t + Q_\alpha[C_{kt} + \tilde{F}_{kt}b + C_{ut} - \lambda_{ut}^\top \mu_u + \epsilon_t - c_t] \\ &= \alpha_t + x^\top \gamma_t + Q_\alpha[C_{kt} + C_{ut} + \epsilon_t] \end{aligned}$$

### C Variance decompositions

As discussed in Section 2.1, an important class of parameters in learning models are terms that decompose the variance of potential outcomes into components that are

predictable and unpredictable given the agents' information. These parameters can be expressed as functionals of the finite- and infinite-dimensional components of the model parameters. Section 4 provides general inference results, which can be applied to a plug-in sieve MLE estimator of these parameters. In this section, we define these parameters and discuss their relevance to quantifying the importance of uncertainty and learning.

To define this class of parameters, consider a weighted sum of potential outcomes,  $Y(\omega, d) = \sum_t \omega_t Y_t(d_t)$  for a sequence of choices  $d = \{d_t : t \leq T\}$  and weights,  $\omega = \{\omega_t : t \leq T\}$ . Cunha and Heckman (2016) consider a special case of this parameter in the context of an educational choice model. In particular, they consider the present value of lifetime earnings, which is defined as  $Y(\omega, d)$ , with  $\omega_t = 1(t \geq t_0)(1 - \rho)^{t_0 - t}$ , for some discount rate  $0 \leq \rho < 1$ .

Next, define the agent's information set,  $\mathcal{I}_t = \{Y^{t-1}, D^{t-1}, X^t, X_k^*\}$  for  $t > 1$  and  $\mathcal{I}_1 = \{X_1, X_k^*\}$ . Restricting attention to weighted sums where  $\omega_s = 0$  for  $s < t$ , the variance of  $Y(\omega, d)$  conditional on  $\mathcal{I}_t$  can be understood as the variance that is due to the agent's uncertainty over  $Y(\omega, d)$  given their information up to period  $t$ . We refer to this as the *posterior variance*, because this is derived from the posterior distribution of  $X_u^*$  after performing a Bayesian update with the information in  $\mathcal{I}_t$ .

In its full generality, the model allows for endogeneity in  $X_t$  as the transition probabilities depend on past choices and outcomes. Therefore, the posterior variance of  $Y(\omega, d)$  includes terms which reflects uncertainty about the future realizations of  $X_t$  conditional on  $X_k^*$ . We will abstract from this assuming that  $X^t = X_1 :- X$  for all  $t$  in order to focus on uncertainty over  $X^*$ .<sup>10</sup>

In particular, with this restriction on  $X_t$ , Lemma 1 implies that the posterior

---

<sup>10</sup>When  $t$  is time varying and transitions depend on  $(D^{t-t}, Y^{t-1})$ , the posterior variance will include the covariances between future realizations of  $X_t$  and between  $X_t$  and  $X_u^*$  conditional on the information set. These terms reflect another channel through which unobserved heterogeneity is related to the agents' uncertainty. In this case, the plug-in estimator of the posterior variance will involve all the infinite dimension parameters of the model  $(f_{X_t|X^{t-1}, D^{t-1}, Y^{t-1}}, f_{D_t|X^t, Y^{t-1}, D^{t-1}, X_k^*}, f_{X_k^*, X_1}, \Sigma_u)$ .

variance has the form,  $\text{Var}(Y(\omega, d) \mid \mathcal{I}_t) = V_t^u(X, D^{t-1}; \omega, d)$  where,

$$V_t^u(X, D^{t-1}; \omega, d) := \sum_{t_1, t_2 \geq t} \omega_{t_1} \omega_{t_2} X_{ut_1}^*(d_{t_1})^\top \Sigma_t(D^{t-1}, X) X_{ut_2}^*(d_{t_2}) + \sum_{t_1 \geq t} \omega_{t_1}^2 \sigma_{t_1}^2(d_{t_1})$$

where  $\Sigma_t(D^{t-1}, X_1)$ <sup>11</sup> is the posterior variance of  $X_u^*$  as written in Lemma 1. When  $t = 1$ ,  $D^{t-1}$  is empty, so for notational convenience, we suppress this argument and write  $V_1^u(X_1; \omega, d)$  and  $\Sigma_t(X_1)$ . Note that  $\Sigma_1(X_1) = \Sigma_u(X_1)$ .

At  $t = 1$ , the following variance decomposition provides a natural way to quantify the relative importance of uncertainty in potential outcomes,

$$\text{Var}(Y(\omega, d) \mid X) = V_1^u(X; \omega, d) + \sum_{t_1, t_2 \geq 1} \omega_{t_1} \omega_{t_2} X_{kt_1}^*(d_{t_1}) X_{kt_2}^*(d_{t_2}) \text{Var}(X_k^* \mid X) \quad (8)$$

This corresponds to the decomposition in Cunha and Heckman (2016) and in that context, has the simple interpretation that the first term is the portion of variance in the lifetime earnings that is due to uncertainty and the second part is due to privately known heterogeneity.

For  $t > 1$ , the analysis is more complicated. For any  $t > 1$ ,  $V_t^u(X, D^{t-1}; \omega, d) < V_1^u(X; \omega, d)$ , because the realized outcomes are informative about the  $X_u^*$ . Agents also select  $D^{t-1}$  based on their private information ( $X_k^*$ ), which induces a selected distribution of  $X_k^*$  conditional on  $X, Y^{t-1}, D^{t-1}$ . Given these contributions of learning and selection to variance of  $Y(\omega, d)$ , there are several possible ways of quantifying the relative importance of uncertainty. The following are three alternative decompositions, which express total variance (conditional on some subset of observables) as the sum of a term that reflects uncertainty and another reflects variance induced by

---

<sup>11</sup>Here, we write  $\Sigma_t$  with the arguments  $X, D^{t-1}$  to be explicit in how this random variable depends on the components of  $\mathcal{I}_t$ .



private information  $(X_k^*)$ ,

$$\begin{aligned} \text{Var}(Y(\omega, d) \mid D^{t-1} = d^{t-1}, X = x) &= V_t^u(d^{t-1}, x; \omega, d) \\ &\quad + \text{Var}(\mathbb{E}(Y(\omega, d) \mid \mathcal{I}_t) \mid D^{t-1} = d^{t-1}, X = x) \end{aligned} \quad (9)$$

$$\text{Var}(Y(\omega, d) \mid X = x) = \mathbb{E}(V_t^u(D^{t-1}, x; \omega, d)) + \text{Var}(\mathbb{E}(Y(\omega, d) \mid \mathcal{I}_t) \mid X = x) \quad (10)$$

$$\text{Var}(Y(\omega, d) \mid X = x) = V_t^u(d^{t-1}, x; \omega, d) + \tilde{\text{Var}}(\tilde{\mathbb{E}}(Y(\omega, d) \mid \mathcal{I}_t) \mid X = x) \quad (11)$$

Decomposition (9) compares the variance of uncertainty to the total variance conditional on choosing the sequence  $d^t$ . These are natural parameters to consider, but the ratio,  $V_t^u(d^t, x; \omega, d) / \text{Var}(Y(\omega, d) \mid D^t = d^t, X = x)$  reflects both the effect of learning in the numerator and selection in the denominator.

Decomposition (10) compares the total variance  $Y(\omega, d)$  to the expected posterior variance of  $Y(\omega, d)$  after  $t$  periods. The expectation of  $V^u(D^t, x; \omega, d)$  can be understood as the uncertainty that a randomly chosen person would have in period  $t$  after observing their outcomes and endogenously choosing actions based on that information and their private information.

Finally decomposition (11) is based on a counterfactual distribution. Here  $\tilde{\mathbb{E}}$  and  $\tilde{\text{Var}}$  represent the expectation and variance in a counterfactual distribution where  $D^t$  is assigned randomly. This decomposition compares the variance in  $Y(\omega, d)$  which is due to uncertainty vs. known heterogeneity among people randomly assigned to the choice sequence  $d^t$ .

## D Proofs for estimation section

### D.1 Consistency of sieve MLE

In this section we introduce conditions for the sieve maximum likelihood estimator defined in Equation (6) to be consistent for the true model parameter  $\theta^* \in \Theta$ . We begin by imposing smoothness restrictions on the unknown functions. To do so, given  $\gamma > 0$ ,  $\omega \geq 0$  and  $\mathcal{X}$  a subset of a Euclidean space, let  $\Lambda^\gamma(\mathcal{X})$  denote a Hölder space equipped with the Hölder norm  $\|h\|_{\Lambda^\gamma}$  (that is, for  $k$  the largest integer smaller than  $\gamma$ ,  $\Lambda^\gamma(\mathcal{X})$  is a space of functions  $h: \mathcal{X} \rightarrow \mathbf{R}$  having at least  $k$  continuous derivatives, the  $k$ th of which is Hölder continuous with exponent  $\gamma - k$ ). Then define a weighted Hölder ball with radius  $c \in (0, \infty)$  as  $\Lambda_c^{\gamma, \omega}(\mathcal{X}) = \{h \in \Lambda^\gamma(\mathcal{X}): \|h(\cdot)[1 + \|\cdot\|_E^2]^{-\omega}\|_{\Lambda^\gamma} \leq c\}$ , where  $\|\cdot\|_E$  is the Euclidean norm.

Without loss of generality, suppose the CCP function  $h_t(d^t, x^t, y^{t-1}, v_k)$  depends on  $(d^t, x^t, y^{t-1})$  via some measurable vector-valued function  $(d^t, x^t, y^{t-1}) \mapsto j_t$  which is known up to  $((\beta_{st}, \lambda_{kst}, \lambda_{ust}, \sigma_{st})_{st=1}^T, \Sigma_u)$ . This is without loss of generality since the function may be identity. Other examples include rational learning where  $j_t \in \mathbb{R}^{p(p+3)/2+2}$  includes sufficient statistics for  $X_u^*$  (i.e, the mean and variance), and a sort of myopia where  $j_t \in \mathbb{R}^{3+2}$  depends on the history only via the previous period  $(d_{t-1}, x_{t-1}, y_{t-1})$ . Write  $J_t = (J_{1,t}^\top, J_{2,t}^\top)^\top$  and  $X_t = (X_{1,t}^\top, X_{2,t}^\top)^\top$  where  $J_{1,t}, X_{1,t}$  are continuous random variables and  $J_{2,t}, X_{2,t}$  are random variables with finite support and, with some abuse of notation, redefine the CCP function as  $h_t(j_{1,t}, j_{2,t}, v_k)$ . Define

$$\begin{aligned}\mathcal{H}_t &= \Lambda_c^{\gamma_1, \omega_1}(\mathcal{S}(X_k^*) \times \mathcal{S}(J_{1,t})), \\ \mathcal{F} &= \{f: \mathcal{S}(X_k^*, X_{1,1}) \rightarrow \mathbb{R} \mid f(\cdot, x_1) \text{ is càdlàg}, f(v, \cdot) \in \Lambda_c^{\gamma_2, \omega_2}(\mathcal{S}(X_{1,1}))\} \\ \mathcal{G}_t &= \Lambda_c^{\gamma_3, \omega_3}(\mathcal{S}(X_{1,t+1}) \times \mathcal{S}(Y_t) \times \mathcal{S}(X_{1,t})).\end{aligned}$$

The use of a weighted Holder space enables us to allow the support of the continuous random variables to be unbounded. Though not required for consistency,

Assumption **E6** places restrictions on  $(\gamma_1, \gamma_2, \gamma_3)$ , the parameters that govern the smoothness of the function classes. Next, to simplify notation we make the following assumption which strengthens Assumption **KL1**:

**Assumption E1.** For any  $t$ ,  $F_{X_{t+1}|Y^t D^t X^t} = F_{X_{t+1}|Y_t D_t X_t}$ .  $F_{X_U^*|X_1} = F_{X_U^*}$ .

Define  $k_{1,t} = |\mathcal{S}(J_{2,t})|$ ,  $k_2 = |\mathcal{S}(X_{2,1})|$ , and  $k_{3,t} = |\mathcal{S}(X_{2,t+1}, D_t, X_{2,t})|$ . Notice that  $\Theta = \Theta_1 \times \mathcal{H}_1^{k_{1,1}} \times \dots \times \mathcal{H}_T^{k_{1,T}} \times \mathcal{F}^{k_2} \times \mathcal{G}_1^{k_{3,1}} \times \dots \times \mathcal{G}_{T-1}^{k_{3,T-1}}$  and we denote an element of  $\Theta$  as  $\theta = (\theta_1, h_1, \dots, h_T, f_{X^*}, g_1, \dots, g_{T-1})$ . Define the norms on  $\mathcal{H}_t^{k_{1,t}}$ ,  $\mathcal{F}^{k_2}$  and  $\mathcal{G}_t^{k_{3,t}}$  as follows:

$$\begin{aligned} \|h_t\|_{\infty, \omega_1} &= \sup_{j_2 \in \mathcal{S}(J_{2,t})} \|h_t(\cdot, j_2, \cdot) [1 + \|\cdot\|_E^2]^{-\omega_1}\|_{\infty}, \\ \|f_{X^*}\|_{\infty, \omega_2} &= \sup_{x_2 \in \mathcal{S}(X_{2,1})} \|f_{X^*}(\cdot, (\cdot, x_2)) [1 + \|\cdot\|_E^2]^{-\omega_2}\|_{\infty}, \\ \|g_t\|_{\infty, \omega_3} &= \sup_{(x'_2, d, x_2) \in \mathcal{S}(X_{2,t+1}, D_t, X_{2,t})} \|g_t((\cdot, x'_2); \cdot, d, (\cdot, x_2)) [1 + \|\cdot\|_E^2]^{-\omega_3}\|_{\infty}, \end{aligned}$$

where  $\|\cdot\|_{\infty}$  is the uniform norm. Finally, define a metric  $d$  on  $\Theta$  as

$$d(\theta, \tilde{\theta}) = \|\theta_1 - \tilde{\theta}_1\|_E + \sum_{t=1}^T \|h_t - \tilde{h}_t\|_{\infty, \tilde{\omega}_1} + \|f_{X^*} - \tilde{f}_{X^*}\|_{\infty, \tilde{\omega}_2} + \sum_{t=1}^{T-1} \|g_t - \tilde{g}_t\|_{\infty, \tilde{\omega}_3},$$

for scalars  $\tilde{\omega}_1, \tilde{\omega}_2, \tilde{\omega}_3$ . Now, let  $\mathcal{H}_{n,t}$ ,  $\mathcal{F}_n$  and  $\mathcal{G}_{n,t}$  be sieve spaces for  $\mathcal{H}_t$ ,  $\mathcal{F}$  and  $\mathcal{G}_t$  respectively. Then  $\Theta_n = \Theta_1 \times \mathcal{H}_{n,1}^{k_{1,1}} \times \dots \times \mathcal{H}_{n,T}^{k_{1,T}} \times \mathcal{F}_n^{k_2} \times \mathcal{G}_{n,1}^{k_{3,1}} \times \dots \times \mathcal{G}_{n,T-1}^{k_{3,T-1}}$  and

$$\frac{1}{n} \sum_{i=1}^n \ell(w_i; \hat{\theta}) \geq \sup_{\theta \in \Theta_n} \frac{1}{n} \sum_{i=1}^n \ell(w_i; \theta) - o_p(1/n).$$

**Assumption E2.**  $\theta^* \in \Theta$  and  $(\Theta, d)$  is compact.

**Assumption E3.** For each  $n \geq 1$ ,  $\Theta_n \subseteq \Theta_{n+1} \subseteq \Theta$  and  $\Theta_n$  is compact under  $d$ . As  $n \rightarrow \infty$ ,  $\min_{\theta \in \Theta_n} d(\theta, \theta_0) \rightarrow 0$ .

**Assumption E4.**  $E[\ell(W, \theta)]$  is continuous at  $\theta = \theta^*$

**Assumption E5.**

- (i) For each  $n$ ,  $E[\sup_{\theta \in \Theta_n} |\ell(W, \theta)|]$  is finite.
- (ii) There is a non-zero  $s < \infty$  and integrable random variable  $g(W)$  such that
$$\forall \theta, \tilde{\theta} \in \Theta_n, d(\theta, \tilde{\theta}) < \delta \implies |\ell(W, \theta) - \ell(W, \tilde{\theta})| \leq \delta^s g(W).$$
- (iii) For all  $\delta > 0$ ,  $\log N(\delta^{1/s}, \Theta_n, d) = o(n)$ .

The identification assumptions imply  $\theta^* = \arg \max_{\theta \in \Theta} E[\ell(W, \theta)]$  and for all  $\theta \in \Theta \setminus \{\theta^*\}$ ,  $E[\ell(W, \theta^*)] \geq E[\ell(W, \theta)]$ . By assuming compactness of  $\Theta$ , we ensure that  $\theta^*$  is a well-separated maximum of  $E[\ell(W, \theta)]$ . Assumption E3 requires the sieve space  $\Theta_n$  to be a good approximation to  $\Theta$ . Assumption E4 requires the population criterion to be continuous. Finally, Assumption E5 is similar to Condition 3.5M in Chen (2007).

Theorem 3 follows from Remark 3.3 in Chen (2007), so its proof is omitted.

## D.2 Plug-in sieve estimator

We assume a linear sieve space and limit its complexity.

**Assumption E6.** (1)  $\mathcal{H}_{n,t}$ ,  $\mathcal{F}_n$  and  $\mathcal{G}_{n,t}$  are linear sieves of length  $M_{Hn,t}$ ,  $M_{Fn}$  and  $M_{Gn,t}$  respectively, where  $M_{Hn,t} = O(n^{\frac{1}{2\gamma_1/(1+\dim(J_{1,t}))+1}})$ ,  $M_{Fn} = O(n^{\frac{1}{2\gamma_2/(1+\dim(X_{1,1}))+1}})$ , and  $M_{Gn,t} = O(n^{\frac{1}{2\gamma_3/(\dim(X_{1,t+1})+1+\dim(X_{1,t}))+1}})$ . (2)  $\min \left\{ \frac{\gamma_1}{1+\dim(J_{1,t})}, \frac{\gamma_2}{1+\dim(X_{1,1})}, \frac{\gamma_3}{\dim(X_{1,t+1})+1+\dim(X_{1,t})} \right\} > 1/2$ .

Assumption E6 controls the rate at which the number of sieve terms grow. To achieve this, part (2) of Assumption E6 requires that the nonparametric functions have adequate smoothness. In applied work, one may focus on discrete  $X_t$  and posit a parametric model for  $h_t$ , in which case the above restrictions are milder.

The next assumption strengthens E3 and ensures the number of sieve terms grows sufficiently fast.

**Assumption E7.**  $\min_{\theta \in \Theta_n} d(\theta, \theta^*) = o(n^{-1/4})$ .

Assume  $\ell$  is pathwise differentiable and define an inner product on  $\Theta$  as

$$\langle \theta_1 - \theta^*, \theta_2 - \theta^* \rangle = -\frac{\partial^2}{\partial \tau_1 \partial \tau_2} E[\ell(W, \theta^* + \tau_1(\theta_1 - \theta^*) + \tau_2(\theta_2 - \theta^*))] \Big|_{\tau_1=0, \tau_2=0}, \quad (12)$$

for  $\theta_1, \theta_2 \in \Theta$ . the corresponding norm for  $\theta \in \Theta$  as

$$\|\theta - \theta^*\|^2 := -\frac{\partial^2}{\partial \tau^2} E[\ell(W, \theta^* + \tau(\theta - \theta^*))] \Big|_{\tau=0}. \quad (13)$$

**Assumption E8.** There is  $C_1 > 0$  such that for all small  $\varepsilon > 0$

$$\sup_{\{\theta \in \Theta_n : \|\theta - \theta^*\| \leq \varepsilon\}} \text{Var}(\ell(W, \theta) - \ell(W, \theta^*)) \leq C_1 \varepsilon^2$$

**Assumption E9.** For any  $\delta > 0$ , there exists a constant  $s \in (0, 2)$  such that

$$\sup_{\{\theta \in \Theta_n : \|\theta - \theta^*\| \leq \delta\}} |\ell(W, \theta) - \ell(W, \theta^*)| \leq \delta^s U(W)$$

with  $E([U(W)]^\gamma) \leq C_2$  for some  $\gamma \geq 2$ .

The following theorem is now a consequence of Theorem 3.2 in Chen (2007) or Theorem 1 in Shen and Wong (1994).

**Theorem 5.** Let  $(Y_{it}, D_{it}, X_{it} : t = 1, \dots, T)_{i=1}^n$  be i.i.d. data where  $T \geq 2p + 1$  and Assumptions **KL1-KL5** and Assumptions **E1-E9** hold. Then  $\|\hat{\theta} - \theta^*\| = o_p(n^{-1/4})$ .

Given the preceding result, we focus on a shrinking neighborhood of  $\theta^*$ . Let

$$\mathcal{N}_0 := \{\theta \in \Theta : \|\theta - \theta^*\| = o(n^{-1/4}), d(\theta, \theta^*) = o(1)\},$$

and  $\mathcal{N}_n := \mathcal{N}_0 \cap \Theta_n$ . Define  $\theta_n^* = \text{argmin}_{\theta \in \mathcal{N}_n} \|\theta - \theta^*\|$ . Let  $\mathcal{V}$  denote the closed (under  $\|\cdot\|$ ) linear span of  $\mathcal{N}_0$  centered at  $\theta^*$ , and define  $\mathcal{V}_n$  as the analogous closure of  $\mathcal{N}_n$ .

Then we define a linear approximation to  $\ell(W, \theta) - \ell(W, \theta^*)$  as the directional

derivative of  $\ell$  at  $(W, \theta^*)$  in the direction  $(\theta - \theta^*)$ :

$$\frac{\partial \ell(W, \theta^*)}{\partial \theta}[\theta - \theta^*] := \left. \frac{\partial \ell(W, \theta^* + \tau(\theta - \theta^*))}{\partial \tau} \right|_{\tau=0}.$$

Likewise, let  $\frac{\partial f(\theta^*)}{\partial \theta}[v] = \left. \frac{\partial f(\theta^* + \tau v)}{\partial \tau} \right|_{\tau=0}$  for any  $v \in \mathcal{V}$ .

**Assumption E10.** Let  $\mathcal{T}$  be an epsilon ball about  $0 \in \mathbb{R}$ . (1) for all  $\theta \in \mathcal{N}_0$  and  $W$ , the derivative  $\partial \ell(W, \theta^* + \tau(\theta - \theta^*)) / \partial \tau$  exists for all  $\tau \in \mathcal{T}$ ; (ii) for all  $\theta \in \mathcal{N}_0$ ,  $E[\ell(W, \theta^* + \tau(\theta - \theta^*))]$  is finite for each  $\tau \in \mathcal{T}$ ; (3) for all  $\theta \in \mathcal{N}_0$ ,  $E[\sup_{\tau \in \mathcal{T}} |\frac{\partial}{\partial \tau} \ell(W, \theta^* + \tau[\theta - \theta^*])|] < \infty$ .

Assumption E10 provides sufficient conditions for the set  $\mathcal{V}$  to be a Hilbert space under  $\langle \cdot, \cdot \rangle$ <sup>12</sup>. Define  $v_n^*$  to be the Riesz representer of  $\frac{\partial f(\theta^*)}{\partial \theta}[\cdot]$  on  $\mathcal{V}_n$ , which exists under Assumption E11.

**Assumption E11.** (1)  $v \mapsto \frac{\partial f(\theta^*)}{\partial \theta}[v]$  is a linear functional. (2) If  $\lim_{n \rightarrow \infty} \|v_n^*\|$  is finite then  $\|v_n^* - v^*\| \times \|\theta_n^* - \theta^*\| = o(n^{-1/2})$  where  $v^*$  is the limit of  $v_n^*$ . Otherwise  $\left| \frac{\partial f(\theta^*)}{\partial \theta}[\theta_n^* - \theta^*] \right| / \|v_n^*\| = o(n^{-1/2})$ . (3)  $\sup_{\theta \in \mathcal{N}_0} \frac{|f(\theta) - f(\theta^*) - \frac{\partial f(\theta^*)}{\partial \theta}[\theta - \theta^*]|}{\|v_n^*\|} = o(n^{-1/2})$ .

Assumption E11 imposes some restrictions on the functional of interest  $\theta \mapsto f(\theta)$ . Part (1) imposes that the directional derivative is a linear functional, a mild condition that is satisfied by our examples in Section 4. Part (2) is a restriction on the growth rate of the dimension of the sieve space. Part (3) restricts the linear approximation error of  $f(\cdot)$  in a neighborhood of  $\theta^*$ , for which sufficient conditions could be stated in terms of the smoothness of  $f(\cdot)$  and the growth rate of the dimension of the sieve space. See Chen et al. (2014) for further discussion.

Let  $u_n^* := \frac{v_n^*}{\|v_n^*\|}$ ,  $\varepsilon_n = o(n^{-1/2})$  and  $\mu_n\{g(\mathbf{W})\} := n^{-1} \sum_{i=1}^n [g(W_i) - E[g(W_i)]]$  denote the centered empirical process indexed by the function  $g$ .

---

<sup>12</sup>See Chen and Liao (2014, p. 642).

**Assumption E12.**  $\mu_n \left\{ \frac{\partial \ell(\mathbf{W}, \theta^*)}{\partial \theta} [v] \right\}$  is linear in  $v \in \mathcal{V}$ .

$$\sup_{\theta \in \mathcal{N}_n} \mu_n \left\{ \ell(\mathbf{W}, \theta \pm \varepsilon_n u_n^*) - \ell(\mathbf{W}, \theta) - \frac{\partial \ell(\mathbf{W}, \theta^*)}{\partial \theta} [\pm \varepsilon_n u_n^*] \right\} = O_p(\varepsilon_n^2).$$

For some positive sequence  $\eta_n \rightarrow 0$ ,

$$\sup_{\theta \in \mathcal{N}_n} \left| E[\ell(W, \theta) - \ell(W, \theta \pm \varepsilon_n u_n^*)] - \frac{\|\theta \pm \varepsilon_n u_n^* - \theta^*\|^2 - \|\theta - \theta^*\|^2}{2} (1 + O(\eta_n)) \right| = O(\varepsilon_n^2).$$

**Assumption E13.**  $\sqrt{n} \mu_n \left\{ \frac{\partial \ell(\mathbf{W}, \theta^*)}{\partial \theta} [u_n^*] \right\} \rightarrow_d N(0, 1)$

Theorem 4 is a direct application of Lemma 2.1 in Chen and Liao (2014) so its proof is omitted.

## E Appendix to Monte Carlo simulations and implementation section

### E.1 Implicit Differentiation

Here we show how to calculate derivative the profile likelihood function. Let  $\mathcal{N}(r)$  be the set of observations with  $z_i = z^{(r)}$ , and let  $\omega_{\cdot, r} = (\omega_{1r}, \dots, \omega_{q_n, r})$ . The log likelihood for this set of observations is,

$$\mathcal{L}_r(\omega_{\cdot, r}, \theta_1) = \sum_{i \in \mathcal{N}(r)} \ell(w_i; \theta_1, \omega_{\cdot, r}) = \sum_{i \in \mathcal{N}(r)} \log \sum_{s=1}^{q_n} \omega_{sr} f(w_i, \bar{v}_{s, r}; \theta_1)$$

Let  $\omega_{\cdot, r}^*(\theta)$  be the solution that maximizes  $\mathcal{L}_r(\omega_{\cdot, r}, \theta_1)$  with respect to  $\omega_{\cdot, r}$  subject to the constraint that  $\omega_{\cdot, r} \in \Delta(q_n)$ , and let  $\mathcal{L}_r^*(\theta_1) = \mathcal{L}(\omega_{\cdot, r}^*(\theta_1), \theta_1)$  be the profile likelihood function. The gradient of the profile likelihood function is,

$$\frac{d\mathcal{L}_r^*(\theta_1)}{d\theta_1} = \sum_{i \in \mathcal{N}(r)} \frac{\partial \ell(w_i; \theta_1, \omega_{\cdot, r}^*(\theta_1))}{\partial \theta_1} + \frac{\partial \ell(w_i; \theta_1, \omega_{\cdot, r}^*(\theta_1))}{\partial \omega_{\cdot, r}^*(\theta_1)} \frac{d\omega_{\cdot, r}^*(\theta_1)}{d\theta_1}$$

The derivatives of the likelihood function can be calculated directly. The derivative  $\frac{d\omega_{\cdot,r}^*(\theta_1)}{d\theta_1}$  is defined implicitly by the Karush-Kuhn-Tucker (KKT) conditions of the inner optimization step. We next show how to calculate it.

Proposition 3.3 in Kim et al. (2020) shows that for each  $\theta_1$ , maximizing  $\mathcal{L}_r(\omega_{\cdot,r}, \theta_1)$  subject to  $\omega_{\cdot,r} \in \Delta(q_n)$  is equivalent to maximizing  $\mathcal{L}_r(\omega_{\cdot,r}, \theta_1) + \sum_{s=1}^{q_n} \omega_{s,r}$  subject to  $\omega_{s,r} \geq 0$  for all  $s$ .

The equality constraints in the KKT conditions of this equivalent problem are  $G_r(\omega_{\cdot,r}, \lambda_r; \theta_1) = 0$  where  $\lambda_r \in \mathcal{R}^{q_n}$  are the dual parameters, and,

$$G_r(\omega_{\cdot,r}, \lambda_r, \theta_1) = \begin{bmatrix} \sum_{i \in \mathcal{N}(r)} \frac{f(w_i; \bar{v}_{\cdot,r}; \theta_1)}{\sum_{s=1}^{q_n} \omega_{s,r} f(w_i; \bar{v}_{\cdot,r}; \theta_1)} + \iota_{q_n} + \text{diag}(\lambda_r) \\ \lambda_r \circ \omega_{\cdot,r} \end{bmatrix}$$

Since  $G_r$  is constant along the solution  $(\omega_{\cdot,r}^*, \lambda_r^*)(\theta_1) := (\omega_{\cdot,r}^*(\theta_1), \lambda_r^*(\theta_1))$ , for all  $\theta_1$ , we have,

$$0 = \frac{\partial G_r(\omega^*(\theta_1), \lambda^*(\theta_1), \theta_1)}{\partial(\omega_{\cdot,r}^*(\theta_1), \lambda_r^*(\theta_1))} \frac{d(\omega_{\cdot,r}^*, \lambda_r^*)(\theta_1)}{d\theta_1} + \frac{\partial G_r(\omega^*(\theta_1), \lambda^*(\theta_1), \theta_1)}{\partial \theta_1}.$$

This implicitly defines  $\frac{d\omega_{\cdot,r}^*(\theta_1)}{d\theta_1}$  assuming that the partial derivative of  $G_r$  with respect to its first two arguments is invertible.

The partial derivatives of  $G_r$  can be calculated as follows,

$$\begin{aligned} \frac{\partial G_r(\omega_{\cdot,r}, \lambda_r, \theta_1)}{\partial(\omega_{\cdot,r}, \lambda_r)} &= \begin{bmatrix} -\sum_{i=1}^n \frac{f(w_i; \bar{v}_{\cdot,r}; \theta_1) f(w_i; \bar{v}_{\cdot,r}; \theta_1)^T}{\left(\sum_{s=1}^{q_n} \omega_{sr} f(w_i; \bar{v}_{\cdot,r}; \theta_1)\right)^2} & I \\ \text{diag}(\lambda_r) & \text{diag}(\omega_{\cdot,r}) \end{bmatrix} \\ \frac{\partial G_r(\omega_{\cdot,r}, \lambda_r, \theta_1)}{\partial \theta_1} &= \begin{bmatrix} -\sum_{i=1}^n \frac{\nabla_{\theta_1} f(w_i; \bar{v}_{\cdot,r}; \theta_1) \sum_{s=1}^{q_n} \omega_{sr} f(w_i; \bar{v}_{\cdot,r}; \theta_1) - f(w_i; \bar{v}_{\cdot,r}; \theta_1) \sum_{s=1}^{q_n} \omega_{sr} \nabla_{\theta_1} f(w_i; \bar{v}_{\cdot,r}; \theta_1)}{\left(\sum_{s=1}^{q_n} \omega_{sr} f(w_i; \bar{v}_{\cdot,r}; \theta_1)\right)^2} \\ 0 \end{bmatrix} \end{aligned}$$

Summing over  $r$ , we obtain the derivative of the full profile likelihood function,  $\frac{d\mathcal{L}^*(\theta_1)}{d\theta_1} = \sum_r \frac{d\mathcal{L}_r^*(\theta_1)}{d\theta_1}$

Finally, note that from the KKT conditions,  $G_r(\omega_{\cdot,r}, \lambda_r, \theta_1) = 0$  imply that  $\lambda_r^* =$



$-\left(1 + \sum_{i=1}^n \frac{f(w_i; \bar{v}_{sr}; \theta_1)}{\sum_{s'=1}^{qn} \omega_{s'r}^* f(w_i; \bar{v}_{s'r}; \theta_1)}\right)$ . Therefore, it is possible to calculate this gradient even if the values of  $\lambda_r^*(\theta_1)$  available from the maximization procedure used.

## E.2 Details on DGP

This section gives further details on the DGP used for Monte Carlo simulations discussed in Section 5.2. The values of the finite parameters used in the DGP are given in the table below.

Table 4: Finite parameter values

$\alpha_1(1) = 0$	$\beta_{z1,1}(1) = -0.5$	$\beta_{z1,2}(1) = -0.58$	$F_{u1}(1) = 1$	$F_{k1}(1) = 0.3$
$\alpha_2(1) = 0.1$	$\beta_{z2,1}(1) = -0.8$	$\beta_{z2,2}(1) = -0.83$	$F_{u2}(1) = 1.05$	$F_{k2}(1) = 0.35$
$\alpha_3(1) = 0.2$	$\beta_{z3,1}(1) = 0.12$	$\beta_{z3,2}(1) = -0.83$	$F_{u3}(1) = 1.01$	$F_{k3}(1) = 0.33$
$\sigma^2(1) = 0.5$				
$\alpha_1(2) = -0.1$	$\beta_{z1,1}(2) = 0.13$	$\beta_{z1,2}(2) = 0.71$	$F_{u1}(2) = 0.4$	$F_{k1}(2) = 1$
$\alpha_2(2) = -0.22$	$\beta_{z2,1}(2) = 0.89$	$\beta_{z2,2}(2) = -0.36$	$F_{u2}(2) = 0.36$	$F_{k2}(2) = 1.05$
$\alpha_3(2) = -0.33$	$\beta_{z3,1}(2) = 0.32$	$\beta_{z3,2}(2) = -0.36$	$F_{u3}(2) = 0.44$	$F_{k3}(2) = 1.02$
$\sigma^2(2) = 0.7$				
$\sigma_u^2 = 1.5$	$\rho = 2.0$	$\gamma = 0.5$		

The distribution of  $X_k^*$  is a finite mixture of three truncated normal distributions. The means of the component distributions are  $(-1.2, 0, 1.5)$  and the variances of the component distributions are  $(0.2, 0.1, 0.3)$ , and the weights of the mixing distribution are  $(0.4, 0.3, 0.3)$ . Each component distribution is truncated at the third standard deviation of its distribution.

The distribution of the covariates  $X = (X_1, X_2)$  is as follows:  $X_1$  has a standard normal distribution and  $X_2$  as a Bernoulli distribution with equal weights. We assume that  $X_1$  and  $X_2$  are independent from each other and from  $X^*$ .