

Root-N Consistent Semiparametric Estimators of a Dynamic Panel-Sample-Selection Model

George-Levi Gayle

Tepper School of Business, Carnegie Mellon University

Christelle Viauoux

Department of Economics, University of Cincinnati

First Version: September 2004. Current Version : July 2006

Abstract

This paper considers the problem of identification and estimation in panel-data sample-selection models with a binary selection rule when the latent equations contain possibly predetermined variables, lags of the dependent variables, and unobserved individual effects. The selection equation contains lags of the dependent variables from both the latent and the selection equations as well as other possibly predetermined variables relative to the latent equations. We derive a set of conditional moment restrictions that are then exploited to construct a three-step sieve estimator for the parameters of the main equation including a nonparametric estimator of the sample-selection term. In the second step the unknown parameters of the selection equation are consistently estimated using a transformation approach in the spirit of Berkson's minimum chi-square sieve method and a first-step kernel estimator for the selection probability. This second-step estimator is of interest in its own right. It can be used to semiparametrically estimate a panel-data binary response model with a nonparametric individual specific effect without making any other distributional assumptions. We show that both estimators (second and third stage) are \sqrt{n} -consistent and asymptotically normal.

1 Introduction

In this paper, we study a panel-data sample-selection model of the form

$$(1) \quad y_{it}^* = \rho y_{it-1}^* + x_{it}\beta + \alpha_i + \varepsilon_{it}^*,$$

$$(2) \quad y_{it} = d_{it}y_{it}^*,$$

and

$$(3) \quad d_{it} = \mathbb{I}\{\phi_0 d_{it-1} + \phi_1 y_{it-1} + z_{it}\gamma + \eta_i - u_{it} > 0\},$$

where $\mathbb{I}\{\cdot\}$ denotes the usual indicator function, i indexes individuals ($i = 1, \dots, n$), t indexes time ($t = 1, \dots, T$), y_{it}^* is a latent outcome variable, $(y_{it}, d_{it}, x_{it}, z_{it})$ are observed random variables, (α_i, η_i) are unobserved time-invariant individual-specific effects, and $(\varepsilon_{it}^*, u_{it})$ are unobserved random individual time-specific effects assumed to be independent across individuals. The first equation, also called the outcome equation, has an autoregressive structure with x_{it} containing both strictly exogenous variables x_{it}^e and predetermined variables x_{it}^p with respect to ε_{it}^* . The second equation, the selection equation, summarizes the process of observations entering into the sample. It has an autoregressive structure, but also depends explicitly on the lagged outcome of the first equation and on variables z_{it} that may be predetermined with respect to ε_{it}^* and u_{it} .

Although sample-selection models have been studied extensively in the econometrics literature (see Heckman, 1974, 1976, Das, Newey and Vella, 2003, Kyriazidou, 1997, 2001, among others), the model described in (1)–(3) is new in two respects. First, the selection process may depend on variables that are predetermined with respect to the error structure of both the outcome and the selection equations. Second, the outcome equation contains a lagged dependent variable along with other predetermined variables. Therefore our model can be derived from a dynamic utility maximization problem with time-inseparable preferences.

The economic literature contains many examples where panel-sample-selection models that include lags or other predetermined variables would be of interest. Models of this form arise for example in the study of the intertemporal behavior of economic agents (see Hotz, Kydland and Sedlacek, 1988, Altug and Miller, 1998, to name a few). Several applications have shown that the current realization of outcomes or the current decision to participate in the sample is affected by both current and past variables (see Kydland and Prescott, 1982, Altug and Miller, 1998, among others). One illustration is the analysis of a company's investment behavior (Bond and Meghir, 1994). In this case it is reasonable to expect that whether a company decides to invest today will depend on how much it invested in the last period, hence creating a feedback from the continuous outcome (the amount invested last period) to the selection process (the decision to invest today). Variables explaining investment include variables in the agent's information sets that would be correlated with past shocks and hence past values of the dependent variable. These variables are therefore predetermined with respect to the system. Another example of interest is the study of models of life-cycle behavior (Gayle and Miller, 2004) where the dynamic utility maximization of time nonseparable preferences gives rise to feedback effects from the outcome equation to the selection equation. For example, the

decision to participate in the labor force may depend on the wage earned in the past. As the past explains the present, so does the expectation about the future. For example, parents' decision to work in the current period is affected by their expectations about future birth events. Therefore, explanatory variables such as the number of children are predetermined (that is, correlated with lagged values of the error term of the outcome equation, but uncorrelated with its present and future values).¹

Panel-data dynamic models of sample selection have been studied by Kyriazidou (2001). Her model, however differs from ours in that the selection equation does not have any autoregressive structure, nor does it depend on the lag of the dependent variable from the outcome equation. Consequently, it cannot be directly derived from a dynamic utility maximization problem. Nevertheless, the extensions we consider in this paper are not trivial since Kyriazidou's (2001) estimation and identification strategy critically depends on the assumption that the selection equation only contains strictly exogenous regressors. Hence, we will have to pursue a different identification and estimation strategy from Kyriazidou (2001).

The method that we adopt in this paper allows us to identify the structural parameters of the outcome and of the selection equation without placing any parametric assumption on the distribution of the error terms. We will still be able to obtain a \sqrt{n} -consistent estimator for the structural parameters that is asymptotically normally distributed. The \sqrt{n} -consistency of our estimator is an important improvement over that of Kyriazidou (2001). However this generalization and improvement comes at a cost of two additional restrictions both of which enable us to obtain identification of our model. The first major restriction is to impose the individual specific effect in the selection equation to be a nonparametric function of strictly exogenous individual specific variables. The second restriction requires that the distribution of the predetermined and lagged dependent variables in the selection equation conditional on the instruments in the outcome equation forms a complete family of distribution.

The contribution of this paper is twofold. First it develops a new semiparametric estimator for a more general class of panel sample selection models than is found in the current literature. Second, it develops a new semiparametric estimator for dynamic panel binary choice model. This estimation technique does not rely on specifying the distribution of the errors nor does it rely on the distribution of the initial condition.

The remainder of this paper is organized as follows. In Section 2, the identification of the sample-selection model is discussed. In Section 3, our proposed estimator is defined. Consistency and asymptotic normality of our estimator are established in Section 4. A Monte Carlo study of our estimator is conducted in Section 5 to compare the small-sample performance of our estimator with that of Kyriazidou (2001). An appendix

¹A predetermined variable is a regressor that is Granger caused by past values of the dependent variable of the system.

contains proofs of the technical results while proofs that are critical to the flow of the paper are left in the text.

2 Identification

Five major issues need to be addressed in order to achieve identification of the model specified in equations (1)–(3). These issues include the presence in equation (1) of the unobserved time-invariant-individual-specific component, α_i , the presence of the lagged dependent variable, y_{it-1}^* , the presence of predetermined variables in x_{it} , the sample-selection mechanism specified by Equations (2) and (3) and the identification of the binary selection model with lagged dependent variables. In the absence of sample selection, the first three issues have been resolved in the panel-data literature (see Anderson and Hsiao, 1982, Arellano and Bond, 1991, Ahn and Schmidt, 1995, to name a few). The presence of sample selection, however, complicates the identification issue. Another issue is the identification of the selection equation as it contains lagged dependent variables and an individual specific effect.

A number of papers have looked at this problem before (Honore, 1993; Arellano and Carrasco, 2003, Hahn, 1997 and Honore and Kyriazidou, 2000). However, these papers have either studied the case where the distribution of the error term is assumed to be parametric or the case where the distribution of the initial condition is parametrically specified. Moreover, they take one of the following approaches for the individual specific effect: either they assume that it is random with a parametric distribution or that it is fixed and can be eliminated in some ways. In this paper we do not assume that the distribution of the error is parametric. We take an alternative approach on the individual specific effects. The following assumptions summarize that approach.

Assumption 2.1 ($\forall i = 1, \dots, n, t = 1, \dots, T$):

1. $\eta_i = \eta(z_i^1)$ where z_i^1 are the strictly exogenous time-invariant components of z_{it} , $\eta(z_i^1)$ is an unknown function of z_i^1 .
2. $(\varepsilon_{it}^*, u_{it})$ are jointly independent of z_i^1 and z_{it} .

This restriction is still mild since the distribution of the error terms $(\varepsilon_{it}^*, u_{it})$ is left unspecified. This is a nonparametric extension of the approach in MaCurdy (1981), which was formalized in Altug and Miller (1998). A version of this restriction can be found in Chamberlain (1986), Nijman and Verbeek (1992), and Zabel (1992), where $\eta(\cdot)$ is specified as a linear function of strictly exogenous variables and the distribution of the error terms are assumed to be normal. Newey (1994) relaxed the linear functional assumption on $\eta(\cdot)$ while retaining the normality assumption on the errors. In what

follows, we use an approach similar to Chen (1998) by relaxing the assumptions on both the functional form and the distribution of the error term. One viable alternative would be to use the fixed-effect specification of Honore and Lewbel (2002). Their fixed-effect model relies, however, on the existence of a special regressor conditionally independent of the individual-specific effect and of the error term in the model, whereas our original formulation does not. We now state additional conditions under which the selection equation is semiparametrically identified. Let

$$P_{it0} \equiv E[d_{it}|d_{it-1}, y_{it-1}, z_{it}, z_i^1],$$

be the true conditional expectation. Let $w_{it} = [d_{it-1}, y_{it-1}, z_{it}']'$, $\omega_{it} = [w_{it}', z_i^1]'$ and $\tilde{\gamma} = [\phi_0, \phi_1, \gamma']'$.

Assumption 2.2 ($\forall i = 1, \dots, n, t = 1, \dots, T$):

1. $\|\tilde{\gamma}\| = 1$.
2. *The random vector w_{it} contains at least one continuous regressor.*
3. *$E[\Delta w_{it} \Delta w_{it}']$ is invertible.*
4. $E[\eta(z_i^1)] = 0$.
5. *Let $F_u(\cdot)$ be the distribution function of u_{it} , conditional on ω_{it} . $F_u(\cdot)$ is strictly increasing, differentiable and non constant on its support.*

Assumptions 2.2.1 and 2.2.2 are standard in the semiparametric literature on the identification of binary choice models (See Manski, 1987 among others). A well known alternative to Assumption 2.2.2 is that a component of w_{it}' , say w_{itk} (of associated parameter $\tilde{\gamma}_k$) has a probability distribution conditional on the remaining components that are absolutely continuous with respect to the Lebesgue measure, and then assume that $|\tilde{\gamma}_k| = 1$. We do not need to make this assumption here since under Assumption 2.2.5 we are able to estimate the signs of all coefficients. We could however assume that $\tilde{\gamma}_k = 1$ (or $\tilde{\gamma}_k = -1$) which along with Assumption 2.2.2 would allow us to estimate the remaining parameters relative to $\tilde{\gamma}_k$. Assumption 2.2.3 is the traditional full rank condition used for identification in the linear panel data literature. Assumption 2.2.4 is a version of the traditional zero mean assumption in fixed effect models, it serves here as the location normalization. It could be relaxed but then all the nonparametric functions would be identified up to an additive constant. Assumption 2.2.5 is the critical assumption which allows us to identify and estimate our model. Since the identification strategy is very important in understanding the estimation, we will prove it in the text.

Proposition 1 *Under Assumptions 2.1 and 2.2 ($\tilde{\gamma}_0, F_{0u}(\cdot), \eta_0(z_i^1)$) are identified.*

Proof. Let $\pi_0 = (\tilde{\gamma}_0, F_{0u}(\cdot), \eta_0(z_i^1))$ denote the true parameters of our model and let $\pi_1 = (\tilde{\gamma}_1, F_{1u}(\cdot), \eta_1(z_i^1))$ be another set of parameters that are observationally equivalent to π_0 . Then by equation (3) and Assumption 2.1 we obtain that

$$(4) \quad F_{0u}(w'_{it}\tilde{\gamma}_0 + \eta_0(z_i^1)) = F_{1u}(w'_{it}\tilde{\gamma}_1 + \eta_1(z_i^1)),$$

because by observational equivalence both sides are equal to P_{it0} . Assumption 2.2.5 implies that

$$(5) \quad w'_{it}\tilde{\gamma}_0 + \eta_0(z_i^1) = F_{0u}^{-1}(F_{1u}(w'_{it}\tilde{\gamma}_1 + \eta_1(z_i^1))).$$

Since $F_u(\cdot)$ is strictly increasing it is differentiable almost everywhere. Differentiating (5) with respect to the continuous regressor w_{itk} gives

$$(6) \quad \begin{aligned} \tilde{\gamma}_{0k} &= \frac{\partial_{w_{itk}} F_{0u}^{-1}(F_{1u}(w'_{it}\tilde{\gamma}_1 + \eta_1(z_i^1)))}{\partial_{w_{itk}} F_{1u}^{-1}(F_{1u}(w'_{it}\tilde{\gamma}_1 + \eta_1(z_i^1)))} \tilde{\gamma}_{1k}, \\ &= \frac{\partial_{w_{itk}} F_{0u}^{-1}(P_{it0})}{\partial_{w_{itk}} F_{1u}^{-1}(P_{it0})} \tilde{\gamma}_{1k}. \end{aligned}$$

Let $\xi_k = \frac{\tilde{\gamma}_{0k}}{\tilde{\gamma}_{1k}}$, note that $\xi_k > 0$ which follows directly from Assumption 2.2.5. Hence we can rewrite equation (6) as

$$(7) \quad \partial_{w_{itk}} F_{0u}^{-1}(P_{it0}) = \xi_k \partial_{w_{itk}} F_{1u}^{-1}(P_{it0}).$$

Integrating equation (7) over the range $[0, P_{it0}]$ gives us

$$(8) \quad F_{0u}^{-1}(P_{it0}) = \xi_k F_{1u}^{-1}(P_{it0}) + K,$$

where $K = F_{0u}^{-1}(0) - \xi_k F_{1u}^{-1}(0)$. Note that by equation (3), Assumptions 2.1 and 2.2.5, we obtain

$$(9) \quad F_{0u}^{-1}(P_{it0}) = w'_{it}\tilde{\gamma}_0 + \eta_0(z_i^1),$$

$$(10) \quad F_{1u}^{-1}(P_{it0}) = w'_{it}\tilde{\gamma}_1 + \eta_1(z_i^1).$$

Taking the first difference of equations (8), (9) and (10) gives us

$$(11) \quad \Delta F_{0u}^{-1}(P_{it0}) = \xi_k \Delta F_{1u}^{-1}(P_{it0}),$$

$$(12) \quad \Delta F_{0u}^{-1}(P_{it0}) = \Delta w'_{it}\tilde{\gamma}_0,$$

$$(13) \quad \Delta F_{1u}^{-1}(P_{it0}) = \Delta w'_{it}\tilde{\gamma}_1.$$

From equations (11), (12) and (13) we obtain that

$$(14) \quad \Delta w'_{it}\tilde{\gamma}_0 = \xi_k \Delta w'_{it}\tilde{\gamma}_1.$$

Premultiplying equation (14) by Δw_{it} and taking the expectation gives

$$(15) \quad E[\Delta w_{it} \Delta w'_{it}] \tilde{\gamma}_0 = \xi_k E[\Delta w_{it} \Delta w'_{it}] \tilde{\gamma}_1.$$

Assumption 2.2.3 and equation (15) imply that

$$(16) \quad \tilde{\gamma}_0 = \xi_k \tilde{\gamma}_1.$$

Assumption 2.2.1 that $\|\tilde{\gamma}_0\| = \|\tilde{\gamma}_1\| = 1$ implies that $|\xi_k| = 1$. But since $\xi_k > 0$, this means that $\xi_k = 1$, thus

$$(17) \quad \tilde{\gamma}_0 = \tilde{\gamma}_1,$$

$$(18) \quad F_{0u}^{-1}(P_{it0}) = F_{1u}^{-1}(P_{it0}) + K.$$

Equations (9), (10) and (18) imply that

$$(19) \quad w'_{it} \tilde{\gamma}_0 + \eta_0(z_i^1) = w'_{it} \tilde{\gamma}_1 + \eta_1(z_i^1) + K.$$

Equation (17) further gives

$$(20) \quad \eta_0(z_i^1) = \eta_1(z_i^1) + K.$$

Finally, Assumption 2.2.4 implies that $K = 0$, hence the model is identified. ■

The literature on sample selection normally takes two different approaches to the identification of the structural parameters, the first is the standard Heckman's correction and the second is to find a way to eliminate the selection bias indirectly (see Ahn and Powell, 1993 for example). Under the first approach one could either assume the parametric form of the joint error distribution and then obtain the correction term or one could nonparametrically identify the correction term. We will use the nonparametric approach in this paper. Below we state some regularity conditions that allow us to proceed.

Assumption 2.3:

1. ε_{it}^* is independent of y_{i0}^* for all t and for each i .
2. ε_{it}^* is independent of α_i for all t and for each i .
3. The $(\varepsilon_{it}^*, u_{it})$'s are mutually independent for all t and for each i , with $E(\varepsilon_{it}^*) = 0$ and $\text{var}(\varepsilon_{it}^*) = \sigma_\varepsilon^2$.
4. x_{is}^p is independent of ε_{it}^* for all $s \leq t$ and for each i .

5. x_{is}^e is independent of ε_{it}^* for all $s, t = 1, \dots, T$ and for each i .
6. z_{is} is independent of ε_{it}^* for all $s \leq t$ and for each i .
7. y_{i0}^* is i.i.d. with density $f_{y_0}(\cdot)$ for each i .
8. $d_{i0} \in \{0, 1\}$ i.i.d. with $\Pr[d_{i0} = 1] = P_o$ for each i .

Assumption 2.3 is a strengthening of the standard type of assumptions used in the linear panel-data literature to achieve identification. In the standard linear panel-data literature, these are typically conditional mean independence assumptions, while here we need full conditional independence. This is not unusual in the semiparametric literature, however, and is similar to the assumptions used in Kyriazidou (2001).

Let

$$\zeta_{it} \equiv \{y_{it-1}, y_{it-2}, d_{it-2}, x_{it}, z_{it}, z_{it-1}, z_i^1, \alpha_i, d_{it}d_{it-1} = 1\}.$$

Taking the expectation of y_{it}^* in Equation (1) conditional on ζ_{it} gives for $i = 1, \dots, n$, $t = 2, \dots, T$,

$$(21) \quad E(y_{it}^* \mid \zeta_{it}) = \rho y_{it-1} + x_{it}\beta + \alpha_i + \bar{\lambda}(\bar{v}_{it}, \bar{v}_{it-1}, z_i^1),$$

where

$$(22) \quad \bar{v}_{it} = \phi_0 d_{it-1} + \phi_1 y_{it-1} + z_{it}\gamma$$

and

$$(23) \quad \bar{\lambda}(\bar{v}_{it}, \bar{v}_{it-1}, z_i^1) \equiv E(\varepsilon_{it}^* \mid \zeta_{it}).$$

Equation (21) follows by noting that

$$(24) \quad E(\varepsilon_{it}^* \mid \zeta_{it}) = E\{E(\varepsilon_{it}^* \mid \zeta_{it}, \eta_i) \mid \zeta_{it}\}.$$

The inner expectation $E(\varepsilon_{it}^* \mid \zeta_{it}, \eta_i)$ can be expressed as

$$\begin{aligned} & E(\varepsilon_{it}^* \mid \zeta_{it}, \eta_i) \\ &= \frac{E[\varepsilon_{it}^* \mid y_{it-1}, y_{it-2}, d_{it-2}, x_{it}, z_{it}, z_{it-1}, z_i^1, \alpha_i, \eta_i, u_{it} < \bar{v}_{it} + \eta_i, u_{it-1} < \bar{v}_{it-1} + \eta_i]}{\int_{-\infty}^{\bar{v}_{it} + \eta_i} \int_{-\infty}^{\bar{v}_{it-1} + \eta_i} \int_{-\infty}^{\infty} \varepsilon^* f(\varepsilon^*, u_2, u_1 \mid y_{it-1}, y_{it-2}, d_{it-2}, x_{it}, z_{it}, z_{it-1}, z_i^1, \alpha_i, \eta_i) d\varepsilon^* du_2 du_1} \\ &= \frac{\int_{-\infty}^{\bar{v}_{it} + \eta_i} \int_{-\infty}^{\bar{v}_{it-1} + \eta_i} \int_{-\infty}^{\infty} f(\varepsilon^*, u_2, u_1 \mid y_{it-1}, y_{it-2}, d_{it-2}, x_{it}, z_{it}, z_{it-1}, z_i^1, \alpha_i, \eta_i) d\varepsilon^* du_2 du_1}{\int_{-\infty}^{\bar{v}_{it} + \eta_i} \int_{-\infty}^{\bar{v}_{it-1} + \eta_i} \int_{-\infty}^{\infty} f(\varepsilon^*, u_2, u_1 \mid y_{it-1}, y_{it-2}, d_{it-2}, x_{it}, z_{it}, z_{it-1}, z_i^1, \alpha_i, \eta_i) d\varepsilon^* du_2 du_1}. \end{aligned}$$

Furthermore, by Assumptions 2.1.2, 2.3.2, 2.3.4, and 2.3.5 to 2.3.8, the conditional density $f(\varepsilon^*, u_2, u_1 \mid y_{it-1}, y_{it-2}, d_{it-2}, x_{it}, z_{it}, z_{it-1}, z_i^1, \alpha_i, \eta_i)$ equals the joint density $f(\varepsilon^*, u_2, u_1)$, which implies that

$$\begin{aligned}
E(\varepsilon_{it}^* \mid \zeta_{it}, \eta_i) &= \frac{\int_{-\infty}^{\bar{v}_{it} + \eta_i} \int_{-\infty}^{\bar{v}_{it-1} + \eta_i + \infty} \int_{-\infty}^{\infty} \varepsilon^* f(\varepsilon^*, u_2, u_1) d\varepsilon^* du_2 du_1}{\int_{-\infty}^{\bar{v}_{it} + \eta_i} \int_{-\infty}^{\bar{v}_{it-1} + \eta_i + \infty} \int_{-\infty}^{\infty} f(\varepsilon^*, u_2, u_1) d\varepsilon^* du_2 du_1}, \\
(25) \qquad \qquad \qquad &\equiv \lambda(\bar{v}_{it} + \eta_i, \bar{v}_{it-1} + \eta_i).
\end{aligned}$$

Denote by $f_\eta(\cdot)$ the density of η_i , then we can integrate out the unobserved component η_i over $f_\eta(\cdot)$ and

$$\begin{aligned}
E(\varepsilon_{it}^* \mid \zeta_{it}) &= \int \lambda(u + \bar{v}_{it}, u + \bar{v}_{it-1}) f_\eta(u) du, \\
(26) \qquad \qquad \qquad &\equiv \bar{\lambda}(\bar{v}_{it}, \bar{v}_{it-1}, z_i^1).
\end{aligned}$$

Note that if (23) were equal to zero, then (21) could be estimated as a standard dynamic, linear panel-data model. However, (23) in general is not equal to zero and the model to be estimated is of the form

$$(27) \qquad y_{it} = \rho y_{it-1} + x_{it} \beta + \alpha_i + \bar{\lambda}(\bar{v}_{it}, \bar{v}_{it-1}, z_i^1) + \varepsilon_{it},$$

where

$$(28) \qquad \varepsilon_{it} = \varepsilon_{it}^* - \bar{\lambda}(\bar{v}_{it}, \bar{v}_{it-1}, z_i^1).$$

The above correction is now similar to that considered by Heckman (1976) except that it is a multi-index specification instead of a single-index specification. In this paper, we allow $\bar{\lambda}(\cdot)$ to have an unknown functional form, which is similar to the formulation used by Das, Newey and Vella (2003) in the cross-section contexts. Equation (27) is an additive semiparametric regression equation similar to that considered by Robinson (1988), except that $\bar{\lambda}(\cdot)$ depends on the unknown parameters of the selection equation. Suppose for the moment that these parameters were known; then, in order to obtain consistent estimates of the remaining parameters in (27), we would need to correct for the presence of y_{it-1} and of the predetermined variables in x_{it} . Following Arellano and Bover (1995), we would look within the system for instruments that will lead to consistent estimates of the parameters of interest.

Unlike in the linear case, we will put restrictions on the latent variables and derive orthogonality and identification conditions implied on observed variables. These conditions are summarized in the proposition below.

Proposition 2 For $i = 1, \dots, n$, $t = 4, \dots, T$, let

$F_{it} \equiv \{y_{i0}, \dots, y_{it-3}, d_{i0}, \dots, d_{it-3}, x_{i1}^p, \dots, x_{it-2}^p, x_{i1}^e, \dots, x_{iT}^e, z_{i1}, \dots, z_{it-2}, z_i^1, d_{it}d_{it-1} = 1\}$.
Under Assumptions 2.1 and 2.3, the following holds:

$$E[(\Delta y_{it} - \rho_0 \Delta y_{it-1} - \Delta x_{it} \beta_0 - \Delta \bar{\lambda}_0(\bar{v}_{it}, \bar{v}_{it-1}, \bar{v}_{it-2}, z_i^1) \mid F_{it})] = 0.$$

Identification of the parameters $(\rho_0, \beta_0, \Delta \bar{\lambda}_0)$ conditional on (ϕ_0, ϕ_1, γ) is not as straightforward as in the standard linear model, because Δy_{it-1} and \bar{v}_{it} are endogenous with respect to $\Delta \varepsilon_{it}$. This endogeneity arises as both Δy_{it-1} and \bar{v}_{it} are functions of y_{it-1}^* . Moreover, y_{it-1}^* is a function of ε_{it-1}^* , of which $\Delta \varepsilon_{it}$ is a function. Also in the presence of predetermined variables, x_{it}^p and z_{it}^p may be functions of ε_{it-1}^* . If \bar{v}_{it} was not a function y_{it-1}^* and hence Δy_{it-1} and Δx_{it}^p were the only endogenous variables in the model, we could use the standard instrument variable conditions for identification. However, since $\Delta \bar{\lambda}_0$ is a nonparametric function of endogenous variables it becomes a problem of semiparametric identification with endogeneity (see Darolles, Florens and Renault, 2002 and Newey and Powell, 2003 for a discussion of the problem). In order to obtain identification we will impose the following conditions:

Assumption 2.4 ($\forall i = 1, \dots, n$, $t = 3, \dots, T$):

1. The distribution of $(\Delta y_{it-1}, x_{it}^p, \bar{v}_{it})$ conditional on

$$\tilde{F}_{it} \equiv \left\{ \begin{array}{l} y_{i0}, \dots, y_{it-3}, x_{i1}^p, \dots, x_{it-2}^p, x_{i1}^e, \dots, x_{iT}^e, z_{i1}^p, \dots, z_{it-2}^p, \\ x_{i1}^p, \dots, x_{iT}^p, z_{i1}^p, \dots, z_{iT}^p, z_i^1, d_{i0} \times \dots \times d_{it} = 1 \end{array} \right\}$$

forms a complete family of distribution in the sense of Newey and Powell (2003).

2. $\Delta \bar{\lambda}_0(\cdot) \in \Lambda_{c_2}^{p_2}$ with $p_2 > 1$; $E[\Delta \bar{\lambda}_0(\bar{v}_{it}, \bar{v}_{it-1}, \bar{v}_{it-2}, z_i^1) \mid \tilde{F}_{it}] \notin \text{linear span}(\Delta y_{it-1}, \Delta x_{it})$ and $E[\tilde{F}_{it} \tilde{F}_{it}']$ is finite positive definite.²

Proposition 3 Under Assumptions 2.1, 2.3, 2.4 and for \bar{v}_{it} , \bar{v}_{it-1} and \bar{v}_{it-2} known, $(\rho_0, \beta_0, \bar{\lambda}_0(\cdot))$ is identified.

Assumption 2.4.1 states that the conditional distribution of the endogenous variables in the model is complete conditional on the instruments. This assumption places restrictions on the joint distribution of the errors terms and the predetermined variables.

² Λ_c^p is a Hölder ball, this controls the smoothness of the functional space to which $\Delta \bar{\lambda}_0$ may belong. We will formally define a Hölder ball in the next section .

Consider the simple case where x_{it} and z_{it} contains only strictly exogenous variables with respect both $(\varepsilon_{it}^*, u_{it})$ so that the only predeterminedness comes from y_{it-1}^* . Let $f_{\varepsilon u}(y_{it-1}^*, u)$ be the density of $(\varepsilon_{it}^*, u_{it})$. Then a sufficient condition for Assumption 2.4.1 to hold is that the cumulative distribution function associated with the density

$$(29) \quad f(y_{it-1}^* | \tilde{F}_{it}) \equiv \frac{\int_{-\infty}^{\phi_0 + \phi_1 y_{it-1}^* + \gamma z_{it} + \eta_i} f_{\varepsilon u}(y_{it-1}^*, u) du}{\int_{-\infty}^{+\infty} \left[\int_{-\infty}^{\phi_0 + \phi_1 y_{it-1}^* + \gamma z_{it} + \eta_i} f_{\varepsilon u}(y_{it-1}^*, u) du \right] dy_{it-1}^*},$$

forms a complete family of distribution. Assumption 2.4.2 is the same as those imposed by Cosslett (1991) and Newey (1999) in the cross-section sample-selection context and Robinson (1988) in the selection version of the additive semiparametric regression. The proof of Proposition 3 follows from checking appropriate conditions of Theorem 4.3 of Newey and Powell (2003). Note that we could relax the completeness assumption to a bounded completeness assumption. This would restrict the nonparametric function to be bounded and at the same time would broaden the class of distributions (see Blundell, Chen and Kristensen, 2004 for details).

3 Estimation

We consider a three-step estimator, where the first step is a nonparametric estimator of the individual conditional probability of being observed in the sample each period, the second and third step are semiparametric minimum-distance estimators. In a sense, the estimators are analogous to Heckman's (1976) two-step procedure for the cross-sectional Gaussian disturbance model. The difference is that the selection equation is estimated by a distribution-free method that depends on a preliminary nonparametric estimator rather than by a Probit, and a nonparametric approximation of the selection correction function, $\Delta \bar{\lambda}_0(\cdot)$, is used instead of the inverse Mills ratio. In this regard, our estimator is similar to Newey's (1999) two-step procedure for the cross-sectional case, except that we have a preliminary nonparametric estimate of the conditional selection probability. Our estimator is closest to Chen's (1998) three-step procedure for the static panel-sample-selection model. However, he did not introduce lagged-dependent and other predetermined variables, he used a least-square regression in the second and third steps, whereas we have conditional moment restrictions, which lead to different estimators.

We will use sieve-extremum estimation methods to estimate our model. There is a growing literature on this topic with important theoretical contributions by Shen and

Wong (1994), Shen (1997), Chen and Shen (1998), Chen, Linton, Van Keilegom (2003), Ai and Chen (2003), among many others. See also Chen (2005) for an extensive survey on the sieve estimation literature. The most important consideration in choosing a sieve's space for an approximation is how well it approximates a given class of functions. Restricting our attention to functions that belong to a Hölder space, we first introduce a measure of the approximation error that will play an important role in the large-sample properties of our estimator.

Suppose that $\mathcal{X} = \mathcal{X}_1 \times \mathcal{X}_2 \times \cdots \times \mathcal{X}_d$ is the Cartesian product of compact intervals $\mathcal{X}_1, \dots, \mathcal{X}_d$. Let $0 < \tau \leq 1$. A real-valued function \tilde{h} on \mathcal{X} is said to satisfy a Hölder condition with exponent τ if there is a positive number c such that $|\tilde{h}(x) - \tilde{h}(y)| \leq c \|x - y\|_E^\tau$ for all $x, y \in \mathcal{X}$, where $\|x\|_E = \left(\sum_{l=1}^d x_l^2\right)^{\frac{1}{2}}$ is the Euclidean norm of $x = (x_1, \dots, x_d) \in \mathcal{X}$. Given a d-tuple $\delta = (\delta_1, \dots, \delta_d)$ of nonnegative integers, set $[\delta] = \delta_1 + \cdots + \delta_d$ and let \mathcal{D}^δ denote the differential operator defined by

$$(30) \quad \mathcal{D}^\delta = \frac{\partial^{[\delta]}}{\partial x_1^{\delta_1} \cdots \partial x_d^{\delta_d}}.$$

Let φ be a nonnegative integer and set $p = \varphi + \tau$. A real-valued function h on \mathcal{X} is said to be p -smooth if it is φ times continuously differentiable on \mathcal{X} and \mathcal{D}^δ satisfies a Hölder condition with exponent τ for all δ with $[\delta] = \varphi$. More generally, a real valued function $h(x, \theta)$ is said to be Hölder continuous in $\theta \in \Theta$ if there exists a constant $\kappa \in (0, 1]$ and a measurable function $c(x)$ with $E(c(x)^2)$ bounded and such that $|h(x, \theta_1) - h(x, \theta_2)| \leq c(x) \|\theta_1 - \theta_2\|_s^\kappa$ for all $x \in \mathcal{X}$, $\theta_1, \theta_2 \in \Theta$ where $\|\cdot\|_s$ is a norm such as the sup or the L_2 -norm. Denote the class of all p -smooth real-valued functions on \mathcal{X} by $\Lambda^p(\mathcal{X})$ and the space of all φ -fold continuously differentiable real-valued functions on \mathcal{X} by $C^\varphi(\mathcal{X})$. Define a Hölder ball with smoothness $p = \varphi + \tau$ as

$$(31) \quad \Lambda_c^p(\mathcal{X}) = \left\{ \tilde{h} \in C^\varphi(\mathcal{X}) : \sup_{x \in \mathcal{X}} |\tilde{h}(x)| \leq c, \sup_{[\delta]=\varphi} \sup_{x, y \in \mathcal{X}, x \neq y} \frac{|\mathcal{D}^\delta \tilde{h}(x) - \mathcal{D}^\delta \tilde{h}(y)|}{\|x - y\|_E^\tau} \leq c \right\}.$$

We restrict all our nonparametric functions to belong to a Hölder ball because these functions are well approximated by linear sieves,³ which we choose in this paper.

Let us denote by θ a generic real-valued function with bounded domain $\mathcal{X} \subset \mathcal{R}^d$, let $\|\theta\|_\infty \equiv \sup_{x \in \mathcal{X}} |\theta(x)|$ be the L_∞ norm, and $\|\theta\|_{2,leb} \equiv \left\{ \int_{\mathcal{X}} [\theta(x)]^2 dx / vol(\mathcal{X}) \right\}^{\frac{1}{2}}$ be the scaled L_2 norm relative to the Lebesgue measure on \mathcal{X} . The sieve approximation errors to $\theta_0 \in \Lambda_c^p(\mathcal{X})$ in $L_\infty(\mathcal{X}, leb)$ -norm and $L_2(\mathcal{X}, leb)$ -norm are defined as

$$e_{\infty n} \equiv \inf_{h \in \Theta_n} \|h - \theta_0\|_\infty \quad \text{and} \quad e_{2n} \equiv \inf_{h \in \Theta_n} \|h - \theta_0\|_{2,leb}.$$

³A sieve is called a linear sieve if it is a linear span of finitely many known basis functions.

3.1 Selection Equation Estimation

Let $\omega_{it} \equiv (d_{it-1}, y_{it-1}, z'_{it}, z_i^1)'$ be the vector of observed variables that affect the probability of selection. Let (ϕ_{00}, γ_0) denote the true value of (ϕ_0, γ) . In order to obtain a \sqrt{n} -consistent estimator of the finite-dimensional parameters in the outcome equation, a sufficient condition is to require that any estimator of $(\phi_0, \phi_1, \gamma)'$ be asymptotically equivalent to the sample average that depends only on ω_{it} . In particular, let there exist a function $\Psi(\omega)$ such that $\sqrt{n}((\hat{\phi}_0, \hat{\phi}_1, \hat{\gamma})' - (\phi_{00}, \phi_{10}, \gamma_0)') = \sum_{i=1}^n \Psi(\omega_i) / \sqrt{n} + o_p(1)$, where $\omega_i = (\omega_{i1}, \dots, \omega_{iT})'$, $E(\Psi(\omega_i)) = 0$ then $E(\Psi(\omega_i)\Psi(\omega_i)')$ exists and is nonsingular. This is the same requirement as in Newey (1999) in the cross-section context. However, while in the cross-section context there are a number of distribution-free estimators that have this property, in the context of our specific selection equation, there are none to the best of our knowledge. The closest distribution-free estimator that could be used in our context is in Honore and Lewbel (2002), but as pointed out in Section 2 above, this would require us to fundamentally change our identification assumptions. What we will do instead is to develop an example of how one such an estimator can be constructed.

The estimator of the selection parameters $(\phi_0, \phi_1, \gamma)'$ is derived from the relationship between the conditional probability of selection and the selection parameters themselves over the full sample of observations. Conditional on ω_{it} , the probability of selection is $P(\omega_{it}) \equiv E[d_{it} \mid \omega_{it}]$. Equation (3) implies that $\forall i = 1, \dots, n, t = 1, \dots, T$,

$$(32) \quad P_0(\omega_{it}) = \Pr [u_{it} < \phi_0 d_{it-1} + \phi_1 y_{it-1} + z_{it}\gamma + \eta(z_i^1) \mid \omega_{it}].$$

Therefore, an alternative representation of the conditional probability of selection is

$$(33) \quad P_0(\omega_{it}) = F_u(\phi_0 d_{it-1} + \phi_1 y_{it-1} + z_{it}\gamma + \eta(z_i^1)).$$

Under Assumption 2.2.5, $F_u(\cdot)$ is a strictly monotone increasing function, therefore its inverse exists and we obtain the following relationship between the conditional selection probability $P(\omega_{it})$ and the parameters in the selection equation $(\phi_0, \phi_1, \gamma)'$:

$$(34) \quad F_u^{-1}(P_0(\omega_{it})) = \phi_0 d_{it-1} + \phi_1 y_{it-1} + z_{it}\gamma + \eta(z_i^1).$$

Suppose there exists a consistent estimator of $P_0(\omega_{it})$, say $\hat{P}(\omega_{it})$. By taking a mean value expansion of $F_u^{-1}(\hat{P}(\omega_{it}))$ around the true selection probability, $P_0(\omega_{it})$, we obtain

$$(35) \quad F_u^{-1}(\hat{P}(\omega_{it})) = \phi_0 d_{it-1} + \phi_1 y_{it-1} + z_{it}\gamma + \eta(z_i^1) + v_{it} + \xi_{it},$$

where $P_*(\omega_{it})$ lies between $\hat{P}(\omega_{it})$ and $P_0(\omega_{it})$. Let $f_u(\cdot)$ be the density of $F_u(\cdot)$, then $v_{it} = \frac{1}{f_u[F_u^{-1}(P_0(\omega_{it}))]}(\hat{P}(\omega_{it}) - P_0(\omega_{it}))$ and $\xi_{it} = \left(\frac{1}{f_u[P_*(\omega_{it})]} - \frac{1}{f_u[F_u^{-1}(P_0(\omega_{it}))]} \right) (\hat{P}(\omega_{it}) - P_0(\omega_{it}))$. The terms v_{it} and ξ_{it} are weighted discrepancy measures between the true

conditional selection probability and the estimated one. If $\widehat{P}(\omega_{it})$ is a consistent estimator, then v_{it} and ξ_{it} are asymptotically zero. Equation (35) then approximately defines a heteroscedastic transformation model of which we can take first differences to eliminate the unknown nuisance function, $\eta(z_i^1)$, which gives

$$(36) \quad \Delta F_u^{-1}(\widehat{P}(\omega_{it}), \widehat{P}(\omega_{it-1})) \equiv \phi_0 \Delta d_{it-1} + \phi_1 \Delta y_{it-1} + \Delta z_{it} \gamma + \Delta v_{it} + \Delta \xi_{it}.$$

As discussed in section 2 (see Manski, 1987 or Ichimura, 1993 for details), we normalize the parameter of the continuous lagged outcome ϕ_1 to -1 ,⁴ so that (36) can be expressed as

$$\Delta y_{it-1} \equiv \phi_0 \Delta d_{it-1} + \Delta z_{it} \gamma - \Delta F_u^{-1}(\widehat{P}(\omega_{it}), \widehat{P}(\omega_{it-1})) + \Delta v_{it} + \Delta \xi_{it}.$$

Finally, to estimate the selection parameters, one need to find a consistent nonparametric estimator of the conditional selection probability. Among many types available in the nonparametric literature, we will choose the kernel density estimator

$$(37) \quad \widehat{P}(\omega_{it}) = \frac{\sum_{j=1}^n d_{jt} K_h(\omega_{jt} - \omega_{it})}{\sum_{k=1}^n K_h(\omega_{kt} - \omega_{it})},$$

where h is a positive smoothing parameter that goes to zero as the sample size increases, $K_h(u) = \frac{1}{h^{d_\omega}} K(\frac{u}{h})$ for a given kernel K (with compact support S_ω) and d_ω is the number of continuous variables in ω_{it} .

To ease the exposition, the following additional notations are required. Define $\omega_i = \{\omega_{is}\}_{s=1}^T$, $\omega_{i-1} = \{\omega_{is-1}\}_{s=1}^T$, $\omega_{i-2} = \{\omega_{is-2}\}_{s=2}^T$, $\Delta \omega_i = \{\omega_{is}\}_{s=2}^T$, $\Delta \omega_{i-1} = \{\omega_{is}\}_{s=1}^{T-1}$. We define the same way the variables $y_i, y_{i-1}, y_{i-2}, d_i, d_{i-1}, d_{i-2}, z_i, z_{i-1}, z_{i-2}, \bar{v}_i, \bar{v}_{i-1}, \bar{v}_{i-2}, \Delta y_i, \Delta y_{i-1}, \Delta y_{i-2}, \Delta d_i, \Delta d_{i-1}, \Delta d_{i-2}, \Delta z_i, \Delta z_{i-1}, \Delta z_{i-2}, \Delta \bar{v}_i, \Delta \bar{v}_{i-1}, \Delta \bar{v}_{i-2}$.

Let $P(\omega_i) \equiv \{P(\omega_{is})\}_{s=2}^T$ and $P(\omega_{i-1}) \equiv \{P(\omega_{is})\}_{s=1}^{T-1}$. Let $\theta_1 \equiv (\phi_0, \gamma, \Delta F^{-1}(\cdot))' \in \Theta_1$ denote the vector of parameters to be estimated where the infinite-dimensional parameter space $\Theta_1 \equiv \Lambda_1 \times \mathcal{H}_1$ can be decomposed into a finite-dimensional space, Λ_1 , and an infinite-dimensional space \mathcal{H}_1 with $(\phi_0, \gamma) \in \Lambda_1$ and $\Delta F^{-1}(\cdot) \in \mathcal{H}_1$.

Let

$$(38) \quad \ell(\theta_1, \omega_i, \widehat{P}(\omega_i), \widehat{P}(\omega_{i-1})) = -\frac{1}{2} \left([\Delta \widehat{v}_i + \Delta \widehat{\xi}_i]' [\Delta \widehat{v}_i + \Delta \widehat{\xi}_i] \right)',$$

where $\Delta \widehat{v}_i + \Delta \widehat{\xi}_i = \Delta y_{i-1} - \phi_0 \Delta d_{i-1} - \Delta z_i \gamma + \Delta F^{-1}(\widehat{P}(\omega_i), \widehat{P}(\omega_{i-1}))$. We then define our Least Squared (LS) estimator of θ_1 as

$$(39) \quad \sup_{\theta_1 \in \Theta_1} Q_{1n}(\theta_1, \widehat{P}(\omega_i), \widehat{P}(\omega_{i-1})) = \sup_{\theta_1 \in \Theta_1} \frac{1}{n} \sum_{i=1}^n \ell(\theta_1, \omega_i, \widehat{P}(\omega_i), \widehat{P}(\omega_{i-1})).$$

⁴Note that we could alternatively normalize any other continuous variable in z_{it} .

However, since \mathcal{H}_1 is infinite dimensional, maximizing over Θ_1 may not be well defined; even if the maximizer exists, it will generally be too difficult to compute. Instead, the maximization will be restricted to a sequence of approximating spaces, Θ_{1n} , such that $\bigcup_n \Theta_{1n}$ is dense in Θ_1 . These types of estimators are called sieve LS in the econometric literature (see Chen, 2005, for a comprehensive survey of this literature). Following this literature, the selection equation estimator can be redefined as

$$(40) \quad \hat{\theta}_1 = \arg \max_{\theta_1 \in \Theta_{1n}} Q_{1n}(\theta_1, \hat{P}(\omega_i), \hat{P}(\omega_{i-1})),$$

where $\Theta_{1n} \equiv \Lambda_1 \times \mathcal{H}_{1n}$ such that $\bigcup_n \mathcal{H}_{1n}$ is dense in \mathcal{H}_1 .

We will restrict our analysis to linear sieve spaces that are compact, nondecreasing (i.e., $\mathcal{H}_{1n} \subseteq \mathcal{H}_{1n+1} \subseteq \dots \subseteq \mathcal{H}_1$). We provide below the example of two univariate linear sieves, which bivariate extensions will be used to estimate $\Delta F^{-1}(P(\omega_i), P(\omega_{i-1}))$. Since $P(\omega_i)$ lies in the interval $[0, 1]$, let $Pol(J_n)$ denote the space of polynomials on $[0, 1]$ of degree J_n or less; that is,

$$(41) \quad Pol(J_n) = \left\{ \sum_{k=0}^{J_n} a_k x^k, x \in [0, 1] : a_k \in \mathcal{R} \right\}.$$

Let $TriPol(J_n)$ denote the space of trigonometric polynomials on $[0, 1]$ of degree J_n or less; that is

$$(42) \quad TriPol(J_n) = \left\{ \frac{a_0}{2} + \sum_{k=1}^{J_n} a_k \cos(2k\pi x) + b_k \sin(2k\pi x), x \in [0, 1] : a_k, b_k \in \mathcal{R} \right\}.$$

Since we will assume that $\Delta F_0^{-1}(P(\omega_i), P(\omega_{i-1}))$ belongs to a Hölder space, say $\Lambda_{c_1}^{p_1}([0, 1]^2)$, it will be well approximated by the bivariate versions of both $Pol(J_n)$ and $TriPol(J_n)$. We will consider the tensor product linear sieve space \mathcal{H}_{1n} , which is constructed as a tensor product space of the univariate linear approximating spaces $\mathcal{H}_{1n1}, \mathcal{H}_{1n2}$. Let $\dim(\mathcal{H}_{1n}) = k_{1n}$ and $[p]$ be the biggest integer satisfying $[p] < p_1$. Then, the approximation error rates for polynomials or orthogonal wavelets for ΔF_0^{-1} are of order $O(k_{1n}^{-\frac{p_1}{2}})$ (see Timan, 1963).

3.2 Outcome Equation Estimation

For notational ease, the following additional notations will be required. Let $E[g(\theta_2, Z_i) | F_i] \equiv \{E[g_s(\theta_2, Z_{is}) | F_{is}]\}_{s=4}^T$ where $\theta_2 \equiv (\rho, \beta, \Delta \bar{\lambda}(\cdot))$, $Z_{is} \equiv (d_{is}, d_{is-1}, \Delta y_{is}, \Delta y_{is-1}, \Delta x_{is}, \bar{v}_{is}, \bar{v}_{is-1}, \bar{v}_{is-2}, z_i^1)$, $Z_i \equiv \{Z_{is}\}_{s=4}^T$, $F_i \equiv \{F_{is}\}_{s=4}^T$ and

$$(43) \quad g_s(\theta_2, Z_{is}) \equiv d_{is} d_{is-1} (\Delta y_{is} - \rho \Delta y_{is-1} - \Delta x_{is} \beta - \Delta \bar{\lambda}(\bar{v}_{is}, \bar{v}_{is-1}, \bar{v}_{is-2}, z_i^1)).$$

From Proposition 2, we have a set of conditional moment restrictions, which can be expressed as

$$(44) \quad E[g(\theta_{20}, Z_i) \mid F_i] = 0,$$

where θ_{20} is the true value of θ_2 .

Let $\widehat{m}(\theta_2, F_i) \equiv \{\widehat{m}_s(\theta_2, F_{is})\}_{s=4}^T$, where $\widehat{m}_s(\theta_2, F_{is})$ is a consistent nonparametric estimator of the s^{th} element of the vector $E[g(\theta_2, Z_i) \mid F_i]$. Following Ai and Chen (2003), our third-step estimator can be expressed as

$$(45) \quad \inf_{\theta_2 \in \Theta_{2n}} \widehat{Q}_{2n}(\theta_2) = \inf_{\theta_2 \in \Theta_{2n}} \frac{1}{n} \sum_{i=1}^n \widehat{m}(\theta_2, F_i)' [\widehat{\Sigma}_o(F_i)]^{-1} \widehat{m}(\theta_2, F_i),$$

where $\Theta_{2n} = \Lambda_2 \times \mathcal{H}_{2n}$, $\widehat{\Sigma}_o(F_i) \longrightarrow \Sigma_o(F_i)$ in probability and $\Sigma_o(F_i)$ is a positive definite weighting matrix of the same dimension as $E[g(\theta_2, Z_i) \mid F_i]$. We restrict our analysis to linear sieve spaces that are compact, nondecreasing (i.e., $\mathcal{H}_{2n} \subseteq \mathcal{H}_{2n+1} \subseteq \dots \subseteq \mathcal{H}_2$).

Let $b_{oj}(F_{is}), j = 1, 2, \dots, k_{s,n}$ be a sequence of known basis functions that approximate any real-valued L_2 -functions of F_{is} well as $k_{s,n} \longrightarrow \infty$. Denote by $B^{k_{s,n}}(F_{is}) = (b_{o1}(F_{is}), \dots, b_{ok_{s,n}}(F_{is}))'$ and $B_s = (B^{k_{s,n}}(F_{1s}), \dots, B^{k_{s,n}}(F_{ns}))'$. Following Ai and Chen (2003),⁵ a series LS estimator of the conditional expectation $E[g_s(\theta_2, Z_{is}) \mid F_{is}]$ is

$$(46) \quad \widehat{m}_s(\theta_2, F_{is}) = \frac{1}{n} \sum_{j=1}^n g_s(\theta_2, Z_{js}) B^{k_{s,n}}(F_{js})' (B_s' B_s)^{-1} B^{k_{s,n}}(F_{is})$$

As with the standard GMM type estimator, in order to implement the above estimator, one needs to be able to estimate $\Sigma_o(F_i)$. One method is to use a nonparametric conditional-variance estimator of the moment restrictions calculated from a preliminary consistent estimator of θ_2 (see Robinson, 1987, for example). We use a LS estimator of the conditional variance, $\Sigma_o(F_i)$ as $\widehat{\Sigma}_o(F_i) \equiv \widehat{\Sigma}_o(F_i, \widetilde{\theta}_2) \equiv \{\sigma_{0st}(F_{is}, \widetilde{\theta}_2)\}_{s,t=2,\dots,T}$, where

$$(47) \quad \sigma_{0st}(F_{is}, \widetilde{\theta}_2) \equiv \begin{cases} \frac{1}{n} \sum_{j=1}^n g_s(\widetilde{\theta}_2, Z_{js}) g_t(\widetilde{\theta}_2, Z_{jt}) B^{k_{s,n}}(F_{js})' (B_s' B_s)^{-1} B^{k_{s,n}}(F_{is}) & \text{for } s=t \\ 0 & \text{otherwise} \end{cases}$$

and $\widetilde{\theta}_2$ is a preliminary consistent estimator, normally obtained by minimizing

$$(48) \quad \frac{1}{n} \sum_{i=1}^n \widehat{m}(\theta_2, F_i)' \widehat{m}(\theta_2, F_i).$$

⁵Note that any consistent nonparametric estimator of $E[g(\theta_2, Z_i) \mid F_i]$ could be used. However, we follow Ai and Chen (2003) because of the possible large dimension of F_i .

The line in (47) follows from the fact the $F_{is} \subset F_{is+1}$ in conjunction with Proposition 1. At this stage, we will use the multivariate version of $Pol(J_n)$ or $TriPol(J_n)$ for estimating $\Delta\bar{\lambda}(\cdot)$ and corresponding tensor product sieves for the estimation of $\widehat{m}(\theta_2, F_i)$ and $\widehat{\Sigma}_o(F_i)$ (see Chen, 2005).

4 Large-Sample Properties

Recall that $\omega_i \equiv \{\omega_{is}\}_{s=1}^T$. Let $d = \dim(\omega_i)$. For ease of exposition, we redefine the first-step kernel estimator as a component of

$$(49) \quad \widehat{p}(\omega_i) \equiv (\widehat{p}_1(\omega_i)', \widehat{p}_2(\omega_i)')' \equiv \sum_{j=1}^n \widetilde{d}_j K_h(\omega_j - \omega_i),$$

where $\widetilde{d}_j = [1, d_j]'$ so that $\widehat{p}(\omega_i)$ is the kernel estimate of $p_0(\omega_i) \equiv (p'_{01}, p'_{02})' \equiv f_\omega(\omega_i)E[\widetilde{d}_i | \omega_i]$ and $f_\omega(\cdot)$ denotes the marginal density of ω . In particular, $\widehat{P}(\omega_i) = \frac{\widehat{p}_2(\omega_i)}{\widehat{p}_1(\omega_i)}$. This notational change for the conditional expectation⁶ eases the exposition in the results that follow. The conditions summarized in Assumption 4.1 below, ensure that \widehat{p} is close to p_0 for n large enough (see Newey and McFadden, 1994).

Assumption 4.1 ($\forall i = 1, \dots, n$):

1. *There is a version of $p_0(\omega_i)$ that is continuously differentiable of order q ($> r$) with bounded derivatives on an open set,⁷ and $p_{01}(\omega_i) = f_\omega(\omega_i)$ is bounded away from 0 on S_ω , the compact support of ω_i .*
2. *$K(u)$ is differentiable of order q , $K(u)$ is zero outside a bounded set, $\int K(u)du = 1$, and there is a positive integer L , such that for all $j < L$, $\int K(u)(\bigotimes_{l=1}^j u)du = 0$.*
3. *There is $p \geq 4$ such that $E\left(\left\|\widetilde{d}_i\right\|^p\right) < \infty$ and $E\left(\left\|\widetilde{d}_i\right\|^p | \omega_i\right) f_\omega(\omega_i)$ is bounded.*
4. *The bandwidth h satisfies $h(n) \rightarrow 0$ and $n^{1-\frac{2}{p}}h(n)^d \ln n \rightarrow \infty$.*

Assumption 4.2 ($\forall i = 1, \dots, n$):

1. $\Theta_{1n} \subseteq \Theta_{1n+1} \subseteq \dots \subseteq \Theta_1$ compact for all $n \geq 1$ and for any $\theta_1 \in \Theta_1$ there exists $\pi_n \theta_1 \in \Theta_{1n}$ such that $\|\theta_1 - \pi_n \theta_1\| = o(1)$ as n gets large.

⁶Now $E[d | \omega] = p_{02}(\omega)/p_{01}(\omega) = f_\omega(\omega)E[\widetilde{d} | \omega]/f_\omega(\omega)$.

⁷Note that r is the dimension of the continuous components of ω .

2. $\sigma_1^2(P_0, \omega_i) = E\{(\Delta v_i)^2 \mid P_0, \omega_i\} < \infty$.
3. $\Delta F^{-1}(\cdot) \in H_1 = \Lambda_{c_1}^{p_1}[s_l, s_u]^2$, with $p_1 > d/2$ and $[s_l, s_u] \subset [0, 1]$.

Theorem 1: *Under Assumptions 2.1, 2.2, 4.1 and 4.2, if $k_{1n} = O(n^{1/(2p_1+d)})$ then $\|\widehat{\theta}_1 - \theta_{10}\| = O_p(n^{-p_1/(2p_1+d)})$.*

Assumption 4.1 is standard in the nonparametric literature (see Newey, 1994, and Newey and McFadden, 1994, for a proof and discussions of this result).

Assumption 4.2.1 is the regularity condition on the sieve space (see Chen and Shen, 1998). Assumption 4.2.2 is standard, bounding the second conditional moments. Assumption 4.2.3 imposes a smoothness condition on the class of functions. Theorem 1 is a consistency and rate of convergence result, which is an application of Theorem 1 in Chen and Shen (1998).⁸

To study the asymptotic distribution of $\widehat{\theta}_1$, we use a linear approximation of the criterion difference by the corresponding derivatives and the degree of smoothness of $g(\cdot)$, where $g(\theta_1)$ is a real functional of θ_1 . We will show that the sieve estimator has a normal distribution and is efficient when the empirical criterion satisfies certain stochastic equicontinuity conditions. The degree of smoothness of $g(\cdot)$ can compensate for the slowness of the convergence rate of the estimate.

Let $P_i \equiv (P(\omega_i), P(\omega_{i-1}))$ and $P_{i0} \equiv (P_0(\omega_i), P_0(\omega_{i-1}))$. Let $\omega \equiv (d_{-1}, y_{-1}, z', z^1)' \equiv \{\omega_i\}_{i=1}^n$, $P \equiv \{P_i\}_{i=1}^n$ and $P_0 \equiv \{P_{i0}\}_{i=1}^n$ and all other relevant variables in the same manner.

Let

$$(50) \quad D_{u_1^*}(\omega)' \equiv \begin{bmatrix} \Delta d_{-1} \\ \Delta z \end{bmatrix} - u_1^*(P_0),$$

where $u_1^*(P)$ solves the following programming problem

$$(51) \quad \inf_{u_1: E \left[\left\| \begin{bmatrix} \Delta d_{-1} \\ \Delta z \end{bmatrix} - u_1(P_{i0}) \right\|_e^2 \right] > 0} E \left[\left(\begin{bmatrix} \Delta d_{-1} \\ \Delta z \end{bmatrix} - u_1(P_{i0}) \right) \left(\begin{bmatrix} \Delta d_{-1} \\ \Delta z \end{bmatrix} - u_1(P_0) \right)' \right].$$

Assumption 4.3 ($\forall i = 1, \dots, n$):

1. $(\phi_{00}, \gamma_0) \in \text{int}(\Theta_1)$.

⁸These results are summarized in Chen (2005).

2. $E[D_{u_1^*}(\omega_i)'D_{u_1^*}(\omega_i)]$ is positive definite; $\Sigma_1(\omega_i) = E[(\Delta v_i + \Delta \xi_i)^2 | \omega_i]$ is positive definite.
3. Each element of $u_1^*(P_i)$ belongs to the Hölder space $\Lambda^{p_{11}}$ with $p_{11} > \frac{d}{2}$.

Theorem 2: Under Assumptions 2.1, 2.2, 4.1, 4.2 and 4.3, if $k_{1n} = O(n^{1/(2p_1+d)})$, then $\sqrt{n}[(\widehat{\phi}_0, \widehat{\gamma})' - (\phi_{00}, \gamma_0)'] \implies N(0, \Omega_1)$, where

$$\Omega_1 = E[D_{u_1^*}(\omega_i)'D_{u_1^*}(\omega_i)]^{-1}E[D_{u_1^*}(\omega_i)'\Sigma_1(\omega_i)D_{u_1^*}(\omega_i)]E[D_{u_1^*}(\omega_i)'D_{u_1^*}(\omega_i)]^{-1}.$$

Assumption 4.3 is a standard regularity condition in sieve estimation (see Shen, 1997, Chen and Shen, 1998, and Chen, 2005, for discussions of these conditions). Note that as with the standard Berkson minimum chi-square estimation in standard binary-choice models, the asymptotic variance of the estimator depends only on the variance of Δv_i and not on that of $\Delta \xi_i$ (see Amemiya, 1994, p. 277 for an example). Let k_{2n} denote the dimension of the approximating sieve space, \mathcal{H}_{2n} . Let

$$\Sigma_o(F_i) \equiv \begin{pmatrix} \sigma_2^2(F_{i2}) & \dots & o \\ \cdot & \cdot & \cdot \\ 0 & & \sigma_T^2(F_{iT}) \end{pmatrix}$$

where $\sigma_2^2(F_{is}) = E[g_s(\theta_{20}, Z_{is})^2 | F_{is}]$.

Assumption 4.4 ($\forall i = 1, \dots, n$):

1. α_i and z_i^1 are i.i.d over individuals.
2. The support of $\{x_{it}, z_{it}, y_{io}^*, \alpha_i, \varepsilon_{it}^*\}_{t=1}^T$ is compact with nonempty support.
3. $(\rho, \beta, \phi_0, \phi_1, \gamma) \in \Lambda_1 \times \Lambda_2$, with Λ_1, Λ_2 are compact with nonempty interior and $|\rho| < 1, |\phi_0| < 1$.
4. The density of $\{x_{it}, z_{it}, y_{io}^*, d_{io}, \alpha_i, \varepsilon_{it}^*, u_{it}\}_{t=1}^T$ is bounded and bounded away from zero.
5. Either $B^{k_{s,n}}(F_{is})$ is a tensor product of Fourier series with $k_{s,n}k_{2n} \ln(n)/\sqrt{n} = o(1)$ or a tensor product power series with $k_{s,n}^2 k_{2n} \ln(n)/\sqrt{n} = o(1)$, where k_{2n} is defined in (70).
6. $\dim(g(\cdot)) k_{s,n} \geq 1 + \dim(x_{it}) + k_{2n}$; $k_{s,n}^{-\frac{p_3}{\dim(F)}} = o(n^{-\frac{1}{4}})$ and $k_{2n}^{-\frac{p_2}{3+\dim(z^1)}} = o(n^{-\frac{1}{4}})$.

7. The smallest and largest eigenvalues of $E\{B^{k_s,n}(F_{is})B^{k_s,n}(F_{is})'\}$ are bounded and bounded away from zero for all $k_{s,n}$.
8. (i) $\sigma_2^2(F_i, \theta_2) > 0$ and $\Sigma_o(F_i, \theta_2)$ is finite positive definite uniformly over $F_i \in \mathfrak{F}$, $\theta_2 \in N_{on}$; (ii) each element of $g(\theta_2, Z_i)g(\theta_2, Z_i)'$ satisfies an envelope condition and is Hölder continuous in $\theta_2 \in N_o$; (iii) Each element of $\Sigma_o(F_i, \theta_2)$ is in $\Lambda_{\Sigma}^{p_\Sigma}$ with $p_\Sigma > \frac{d_{\mathcal{F}}}{2}$ for all $\theta_2 \in N_{on}$.
9. $\Delta\bar{\lambda}_0(\cdot) \in H_2 = \Lambda_{c_2}^{p_2}(R^{3+\dim(z^1)})$ with $p_2 > 1$ and $E[(1+(\bar{v}_i, \bar{v}_{i-1}, \bar{v}_{i-2}, z_i^1)'(\bar{v}_i, \bar{v}_{i-1}, \bar{v}_{i-2}, z_i^1))^a | F_i]$ is bounded for some $a > p_2$.
10. $E[dd_{-1}\Delta y_i | F_i = \bar{f}]$, $E[dd_{-1}\Delta x_i | F_i = \bar{f}]$, $E[dd_1q^{k_{2n}}(\bar{v}_i, \bar{v}_{i-1}, \bar{v}_{i-2}, z_i^1) | F_i = \bar{f}] \in \Lambda_{c_3}^{p_3}(\mathfrak{F})$, $p_3 > (\dim(F_i))/2$.
11. (i) $E(|\Delta y_i|^4) < \infty$; (ii) $E(\|\Delta y_i\|_E | F_i) < \infty$, $E(\|\Delta y_{i-1}\|_E | F_i) < \infty$, $E(\|\Delta x_i\|_E | F_i) < \infty$ and $E(\sup_{\theta_1 \in \Theta_1} \|\Delta\bar{\lambda}_0(\bar{v}_i(\theta_1), \bar{v}_{i-1}(\theta_1), \bar{v}_{i-2}(\theta_1), z_i^1)\|_E | F_i) < \infty$.

Let v denote $(\bar{v}_{is}, \bar{v}_{is-1}, \bar{v}_{is-2}, z_i^1)$ for $i = 1, \dots, n, s = 3, \dots, T$. Let $F \equiv \{F_i\}_{i=1}^n$ and denote by $\frac{dm(F, \theta_{20})}{d\theta_2}[\theta_{21} - \theta_{22}]$ the first pathwise derivative of $m(\cdot)$ at the direction $[\theta_{21} - \theta_{22}]$ evaluated at θ_{20} .

Theorem 3: For any $\theta_2, \theta_{21}, \theta_{22} \in \Theta_2$, define

$$\|\theta_{21} - \theta_{22}\| = \sqrt{E \left\{ \left\{ \frac{dm(F, \theta_{20})}{d\theta_2}[\theta_{21} - \theta_{22}] \right\}' \Sigma(F, \theta_2)^{-1} \frac{dm(F, \theta_{20})}{d\theta_2}[\theta_{21} - \theta_{22}] \right\}}. \text{ Under Ass-}$$

sumptions 2.1-2.4 and 4.1-4.4, $\|\hat{\theta}_2 - \theta_{20}\| = o_p(n^{-\frac{1}{4}})$.

Assumptions 4.4.1-4.4.10 are standard in the sieve minimum-distance literature (see Ai and Chen, 2003, for a complete discussion on the importance of these regularity conditions). Assumption 4.4.11 is an additional boundedness condition needed to ensure that our moment conditions are continuous in the preliminary estimates of the selection equation.

Let $u_2^* = (u_{21}^*, \dots, u_{2,1+\dim(x)}^*)$ be the solution to

$$(52) \quad \min_{u_{21}^* \in \bar{\mathcal{U}}} E\{[\Delta y_{-1} - E\{u_{21}^*(\bar{v}, \bar{v}_{-1}, \bar{v}_{-2}, z^1) | F\}]^2\}$$

and

$$(53) \quad \min_{u_{2,1+j}^* \in \bar{\mathcal{U}}} E\{[\Delta x_j - E\{u_{21}^*(\bar{v}, \bar{v}_{-1}, \bar{v}_{-2}, z^1) | F\}]^2\} \quad (j = 1, \dots, \dim(x))$$

$$\begin{aligned}
&\text{Let } D_{u_2^*}(F) = [\Delta y_{-1}, \Delta x] - E\{u_2^*(\bar{v}, \bar{v}_{-1}, \bar{v}_{-2}, z^1) \mid F\}, \\
&H_1 = E \left[D_{u_2^*}(F)' \Sigma_0(F, \theta_{20})^{-1} d d_{-1} \left\{ \frac{\partial \bar{\lambda}_0(\bar{v}, \bar{v}_{-1}, z^1)}{\partial \bar{v}} (d_{-1}, z)' + \frac{\partial \bar{\lambda}_0(\bar{v}, \bar{v}_{-1}, z^1)}{\partial \bar{v}} (d_{-2}, z_{-1})' \right\} \right], \\
&\Omega_2 = E\{D_{u_2^*}(F)' \Sigma_0(F, \theta_{20})^{-1} D_{u_2^*}(F)\} \text{ and} \\
&H_2 = E\{D_{u_2^*}(F)' \Sigma_0(F, \theta_{20})^{-1} D_{u_2^*}(F)\}.
\end{aligned}$$

Assumption 4.5:

1. $(\rho_0, \beta_0) \in A_2$.
2. (i) For $j=1, \dots, \dim(\beta) + 1$,
 $E[u_{2j}^*(\bar{v}, \bar{v}_{-1}, \bar{v}_{-2}, z^1) \mid F = \bar{f}] \in \Lambda_{c_4}^{p_4}(\mathfrak{S})$, for $p_4 > (\dim(F))/2$; (ii) H_1 is bounded.

Theorem 4: Under Assumptions 2.1-2.4, 4.1-4.5, $\sqrt{n}[(\hat{\rho}, \hat{\beta})' - (\rho_0, \beta_0)'] \implies N(0, V)$, where $V = H_2^{-1}[\Omega_2 + H_1 \Omega_1 H_1'] H_2^{-1}$.

Assumptions 4.5.1 and 4.5.2(i) are standard in the sieve minimum-distance literature (see Ai and Chen, 2003). Assumption 4.5.2 (ii) is an additional condition needed to correct the asymptotic variance for the preliminary estimates of the selection equation.

It is important to note that as the number of variables in F increases, the smoothness requirement on $E[\Delta y \mid F = \bar{f}]$ and $E[q^{k_{2n}}(\bar{v}, \bar{v}_{-1}, \bar{v}_{-2}, z^1) \mid F = \bar{f}]$ where $q^{k_{2n}}$ is defined as in (70), increases as well. This suggests that if we use less instruments in our estimation, it may have better small-sample properties. This is not new in this literature as it was pointed out by Ahn and Schmidt (1995) in linear panel-data models. To this end it may be practical to use less variables in the conditioning set when estimating $E[g(\theta_2, Z) \mid F]$.

5 Monte Carlo Study

In this section, we present simulation results illustrating the performance of the estimation procedure described in the preceding sections. We consider two models: the first illustrates how our complete model (including both exogenous and predetermined variables in the selection equation) performs in a limited Monte Carlo study estimating both the selection and outcome equations. The second compares the performance of our estimator of the outcome equation to that of Kyriazidou (2001). However, because her estimator cannot handle predetermined variables in the selection equation, we do not include any such regressors in the second model.

5.1 Model 1 : Predetermined Variable in the Selection Equation

Data for the Monte Carlo experiment of this model are generated for $t = 0, 1, 2, 3$ according to the following specification.

For $i = 1, \dots, n; t = 1, 2, 3$,

$$y_{it}^* = \rho y_{it-1}^* + x_{it}\beta + \alpha_i + \varepsilon_{it}^*,$$

$$y_{it} = d_{it}y_{it}^*,$$

$$d_{it} = \mathbb{I}\{\phi_0 d_{it-1} + \phi_1 y_{it-1} + \gamma_1 z_{1it} + \gamma_2 z_{2it} + \eta_i + u_{it} \geq 0\}.$$

In the selection equation, the predetermined variables are generated so that $z_{1it} = 6 + 0.5y_{it-1}^* + v_{it}$, with v_{it} independently distributed as $U(0, 1)$, whereas the strictly exogenous variables z_{2it} are distributed as $N(0, 1)$ and are the same strictly exogenous latent variables as those of the outcome equation, x_{it} . The individual-specific effect $\eta_i = \frac{1}{3}(z_{2i1} + \dots + z_{2i3}) + 2\xi - 0.75$, where ξ is an independently distributed $U(0, 1)$, whereas the time-varying error term, u_{it} , is distributed as $N(0, 0.5)$. In order to satisfy the scale normalization needed for identification of parameters in the selection equation, we normalized $\gamma_2 = 1$ in the estimation. Finally for simulation purposes, we assume that $\phi_0 = \phi_1 = 0.5, \gamma_1 = 1$. In the main equation, $\alpha_i = \frac{1}{3}(z_{2i1} + \dots + z_{2i3})$ and $\varepsilon_{it}^* = 0.8\xi^* + 0.6(u_{it} - E(u_{it}))$, where ξ^* is an independent standard normal vector. The initial observation, y_{i0} , is generated as $y_{i0} = d_{i0}y_{i0}^*$, where y_{i0}^* is $N(-2, 1)$ and d_{i0} is generated according to a binomial distribution. We investigate the small-sample properties of the estimators of ρ and β , whereas their true values are respectively assumed to be $\rho = 0.5, \beta = 3$. Three sample sizes, n , are considered: 500, 1000, 2000.

We first estimate our selection equation. In all cases that follow, the approximation of unknown functions with power or Fourier series, respectively, use the results of the small-sample experiment of our selection equation estimator, displayed in Table 1 for 100 replications. For the design under investigation, we note that, in general, estimation using Fourier series leads to better mean and variance estimates than does the power series. We use the normal multivariate kernel with the bandwidth as the Silverman's rule of thumb. We trimmed the upper and low 2.5% of the data.

We nonparametrically estimate $\Delta F^{-1}(\hat{P}(\omega_{it}), \hat{P}(\omega_{it-1}))$ for $t = 2, 3$ and all i by using the the bivariate Fourier and power series basis functions. We use an order of ten for $n=2000$ (i.e. $k_{1,2000} = 10$), an order of seven for $n=1000$ (i.e. $k_{1,1000} = 7$), and an order of five for $n=500$ (i.e. $k_{1,500} = 5$). The results are reported in Table 1 below. As can be seen from the results, our estimator performs well in this limited Monte Carlo.

TABLE 1
FINITE-SAMPLE PERFORMANCE OF THE SECOND-STEP ESTIMATOR

Sample Size		500		1000		2000	
		Power Series	Fourier Series	Power Series	Fourier Series	Power Series	Fourier Series
ϕ_0	Mean	0.456	0.473	0.466	0.483	0.491	0.492
	Std Dev.	0.0122	0.107	0.092	0.079	0.031	0.025
	RMSE	0.012	0.010	0.007	0.006	2e-3	1e-3
ϕ_1	Mean	0.473	0.482	0.483	0.505	0.492	0.498
	Std Dev.	0.118	0.093	0.067	0.053	0.032	0.019
	RMSE	0.0113	0.0112	0.005	0.004	8e-3	1e-4
γ_1	Mean	0.943	0.963	0.973	0.978	1.003	1.001
	Std Dev.	0.145	0.134	0.076	0.091	0.008	0.006
	RMSE	0.011	9.4e-3	5.6e-3	0.003	8e-3	7e-4
γ_2	Mean	-	-	-	-	-	-
	Std Dev.	-	-	-	-	-	-
	RMSE	-	-	-	-	-	-

In order to estimate the structural parameters ρ and β of our model we use two versions of our estimator. The first version of our estimator is what we will call the sieve instrumental variable (SIV) estimator. This estimator is obtained by noting that the conditional moment restrictions can be written as an increasing sequence of unconditional moment restrictions. Using the identity matrix as the weighting matrix, our conditional moment restrictions imply the following unconditional moment restrictions

$$E[b_{0j}(\mathcal{F}_{it})g_t(\theta_{20}, Z_{it})] = 0$$

for $t = 2, \dots, T$ and $j = 1, \dots, k_{t,n}$. The model can now be estimated using any standard instrumental variable or two stage least square method with $\{b_{0j}(\mathcal{F}_{it})\}_{j=1}^{k_{t,n}}$ as the instruments (see Ai and Chen, 2003 and Chen, 2005 for details). The second version of our estimator is the standard sieve minimum distance (SMD) described in the text.

In this version of the Monte Carlo simulation, we have two moment restrictions for

$$g_3(\theta_2, Z_{i3}) = d_{i3}d_{i2}[\Delta y_{i3} - \rho\Delta y_{i2} - \Delta x_{i3}\beta - \Delta\bar{\lambda}(\bar{v}_{i3}, \bar{v}_{i2}, \bar{z}_{2i})],$$

and

$$g_2(\theta_2, Z_{i2}) = d_{i2}d_{i1}[\Delta y_{i2} - \rho\Delta y_{i1} - \Delta x_{i2}\beta - \Delta\bar{\lambda}(\bar{v}_{i2}, \bar{v}_{i1}, \bar{z}_{2i})],$$

where

$$\bar{v}_{i3} = \phi_0 d_{i2} + \phi_1 y_{i2} + \gamma_1 z_{1i3} + \gamma_2 z_{2i3},$$

$$\bar{v}_{i2} = \phi_0 d_{i1} + \phi_1 y_{i1} + \gamma_1 z_{1i2} + \gamma_2 z_{2i2},$$

$$\bar{v}_{i1} = \phi_0 d_{i0} + \phi_1 y_{i0} + \gamma_1 z_{1i1} + \gamma_2 z_{2i1},$$

and

$$\bar{z}_{2i} = \frac{1}{3}(z_{2i1} + z_{2i2} + z_{2i3}).$$

The conditioning sets becomes

$$\mathcal{F}_{i2} = \{y_{i0}, d_{i0}, x_{i1}, x_{i2}, z_{1i1}, \bar{z}_{2i}, d_{i2}d_{i1} = 1\}$$

and

$$\mathcal{F}_{i3} = \{y_{i0}, y_{i1}, d_{i0}, d_{i1}, x_{i1}, x_{i2}, x_{i3}, z_{1i1}, \bar{z}_{2i}, d_{i3}d_{i1} = 1\}.$$

The series LS estimator of the conditional expectation, $E[g(\theta_2, Z_i) | F_i]$ used in the SMD estimator is

$$(54) \quad \hat{m}(\theta_2, F_i) = \{\hat{m}_t(\theta_2, F_{it})\}_{t=2,3},$$

where

$$(55) \quad \hat{m}_t(\theta_2, F_{it}) = \sum_{j=1}^n g_t(\theta_2, Z_{jt}) B^{k_t, n}(F_{jt})' (B_t' B_t)^{-1} B^{k_t, n}(F_{it})$$

and the LS estimator of the conditional variance $\Sigma_o(F_i)$ as $\hat{\Sigma}_o(F_i) \equiv \{\sigma_{0ts}(F_{it})\}_{t,s=2,3}$, where

$$(56) \quad \sigma_{0st}(F_{is}) \equiv \begin{cases} \sum_{j=1}^n g_s(\tilde{\theta}_{2SIV}, Z_{js}) g_t(\tilde{\theta}_{2SIV}, Z_{jt}) B^{k_t, n}(F_{js})' (B_t' B_t)^{-1} B^{k_s, n}(F_{is}) & t = s \\ 0 & \text{otherwise} \end{cases}$$

where $\tilde{\theta}_{2SIV}$ is the estimator from the SIV estimation above.

As mentioned above, we use both the multivariate Fourier and power series for our estimation. Furthermore, we increase the number of approximating basis functions with the sample size. Table 2 reports the results for the Mean, Standard Deviation, and Root Mean Squared Error of the estimates of ρ and β for 100 replications. For the design under investigation, we note that, both the power and Fourier series do a very good job of approximation. As is expected the SMD has much smaller standard deviation than the SIV. In all cases, however, the RMSE of the proposed estimators decreases as sample size increases at rate at least equal to \sqrt{n} .

TABLE 2
FINITE-SAMPLE PERFORMANCE OF THE THIRD-STEP ESTIMATOR
($\rho = 0.5$ and $\beta = 3$)

Sample Size			500			1000			2000
			Power Series	Fourier Series	Power Series	Fourier Series	Power Series	Fourier Series	
SIV	ρ	Mean	0.493	0.496	0.514	0.499	0.496	0.499	
		Std Dev.	0.069	0.038	0.066	0.006	0.066	0.005	
		RMSE	0.002	0.002	7.4e-4	1.1e-4	4.4e-4	1.6e-5	
	β	Mean	3.026	2.999	3.095	3.000	3.002	3.000	
		Std Dev.	0.710	0.019	0.459	0.002	0.020	0.001	
		RMSE	0.013	5.8e-5	2.7e-3	5.7e-6	1.5e-5	2.4e-7	
SMD	ρ	Mean	0.485	0.485	0.485	0.485	0.478	0.483	
		Std Dev.	0.059	0.057	0.045	0.044	0.033	0.037	
		RMSE	0.027	0.026	0.015	0.015	0.009	0.009	
	β	Mean	2.940	2.956	2.974	2.991	2.968	2.976	
		Std Dev.	0.173	0.198	0.136	0.127	0.094	0.101	
		RMSE	0.082	0.090	0.043	0.040	0.022	0.023	

5.2 Model 2: No Predetermined Variable in the Selection Equation

In order to see how our estimator performs in a small sample study, we compared our SMD estimator to the three estimators proposed in Kyriazidou (2001). Since Kyriazidou's (2001) estimator does not apply to the case in which there are variables in the selection that are predetermined with respect to the outcome we use the following framework for this study.

For $i = 1, \dots, n; t = 1, 2, 3$

$$y_{it}^* = \rho y_{it-1}^* + \alpha_i + \varepsilon_{it}^*,$$

$$y_{it} = d_{it} y_{it}^*,$$

$$d_{it} = \mathbb{I}\{\gamma_1 z_{1it} + \eta_i + u_{it} > 0\},$$

where z_1 is strictly exogenous and normally distributed $N(0, 1)$. Individual-specific effects as well as the time-varying error terms and initial observations follow the same structure and distributional assumptions as the one used in our model above. Namely $\eta_i = \frac{1}{3}(z_{1i1} + z_{1i2} + z_{1i3}) + 2\xi - 0.75$, where ξ is independently distributed $U(0, 1)$, u_{it} is distributed as $N(0, 0.5)$. In the main equation, $\alpha_i = \frac{1}{3}(z_{2i1} + z_{2i2} + z_{2i3})$ and

$\varepsilon_{it}^* = 0.8\xi^* + 0.6(u_{it} - E(u_{it}))$, where ξ^* is an independent standard normal vector. The initial observation, y_{i0} , is generated as $y_{i0} = d_{i0}y_{i0}^*$, where y_{i0}^* is $N(-2, 1)$ and d_{i0} is generated according to a binomial distribution. In both cases, we report the small-sample properties of the structural estimator, $\hat{\rho}$, for the true selection parameter $\gamma_1 = 1$. Again we consider three sample sizes of 500, 1000, and 2000.

In this version of the Monte Carlo study, we have two moment restrictions associated to

$$g_3(\theta_2, Z_{i3}) = d_{i3}d_{i2}[\Delta y_{i3} - \rho\Delta y_{i2} - \Delta\bar{\lambda}(\bar{v}_{i3}, \bar{v}_{i2}, \bar{z}_{2i})]$$

and

$$g_2(\theta_2, Z_{i2}) = d_{i2}d_{i1}[\Delta y_{i2} - \rho\Delta y_{i1} - \Delta\bar{\lambda}(\bar{v}_{i2}, \bar{v}_{i1}, \bar{z}_{2i})],$$

where

$$\bar{v}_{i3} = \phi_0 d_{i2} + \phi_1 y_{i2} + \gamma_1 z_{1i3} + \gamma_2 z_{2i3},$$

$$\bar{v}_{i2} = \phi_0 d_{i1} + \phi_1 y_{i1} + \gamma_1 z_{1i2} + \gamma_2 z_{2i2},$$

$$\bar{v}_{i1} = \phi_0 d_{i0} + \phi_1 y_{i0} + \gamma_1 z_{1i1} + \gamma_2 z_{2i1},$$

and

$$\bar{z}_{2i} = \frac{1}{3}(z_{2i1} + z_{2i2} + z_{2i3}).$$

The conditioning sets becomes

$$\mathcal{F}_{i2} = \{y_{i0}, d_{i0}, z_{1i1}, z_{1i2}, z_{1i3}, \bar{z}_{2i}, d_{i2}d_{i1} = 1\}$$

and

$$\mathcal{F}_{i3} = \{y_{i0}, y_{i1}, d_{i0}, d_{i1}, z_{1i1}, z_{1i2}, z_{1i3}, \bar{z}_{2i}, d_{i3}d_{i1} = 1\}.$$

Below we specify the kernel weighted moment restrictions used in Kyriazidou (2001):

$$(57) \quad \frac{1}{n} \sum_{i=1}^n d_{i0}d_{i1}d_{i2}y_{i0} (\Delta y_{i2} - \rho\Delta y_{i1}) \bar{w}_{i2},$$

$$(58) \quad \frac{1}{n} \sum_{i=1}^n d_{i0}d_{i1}d_{i2}d_{i3}y_{i0} (\Delta y_{i3} - \rho\Delta y_{i2}) \bar{w}_{i3},$$

$$(59) \quad \frac{1}{n} \sum_{i=1}^n d_{i1}d_{i2}d_{i3}y_{i1} (\Delta y_{i3} - \rho\Delta y_{i2}) \bar{w}_{i3},$$

$$(60) \quad \frac{1}{n} \sum_{i=1}^n d_{i0}d_{i1}d_{i2}d_{i3} (y_{i3} - \rho y_{i2}) (\Delta y_{i2} - \rho\Delta y_{i1}) \bar{w}_{i2},$$

$$(61) \quad \frac{1}{n} \sum_{i=1}^n d_{i0} d_{i1} d_{i2} [(y_{i2} - \rho y_{i1})^2 - (y_{i1} - \rho \Delta y_{i0})^2] \bar{\omega}_{i2},$$

and

$$(62) \quad \frac{1}{n} \sum_{i=1}^n d_{i1} d_{i2} d_{i3} [(y_{i3} - \rho y_{i2})^2 - (y_{i2} - \rho \Delta y_{i1})^2] \bar{\omega}_{i3},$$

where $\bar{\omega}_{ij}$ are the kernel weights constructed using a standard normal kernel with bandwidth parameter $h = n^{-\frac{1}{5}}$. We report here only estimates using the optimal weighting matrix. IV refers to the estimation using only moment restrictions (57), (58) and (59). GMM1 refers to the estimation where moment restriction (60) is added to those used in the IV estimation. Finally, GMM2 refers to the estimation using all the moment restrictions, (57)-(62).

TABLE 3
FINITE-SAMPLE PERFORMANCE OF SIEVE ESTIMATORS

Sample Size		500		1000		2000		
$\rho_0 = .5$		Power	Fourier	Power	Fourier	Power	Fourier	
		Series	Series	Series	Series	Series	Series	
	SMD	Mean	0.493	0.467	0.487	0.468	0.486	0.460
		Std Dev.	0.070	0.099	0.057	0.087	0.036	0.050
		RMSE	0.031	0.047	0.018	0.029	0.009	0.014

TABLE 4
FINITE-SAMPLE PERFORMANCE OF WEIGHTED KERNEL ESTIMATORS

Sample Size		500	1000	2000	
$\rho_0 = .5$	IV	Mean	0.5020	0.5017	0.5019
		Std Dev.	0.0844	0.0616	0.0447
		RMSE	0.0755	0.0390	0.0200
GMM1		Mean	0.5032	0.5020	0.5022
		Std Dev.	0.0858	0.0644	0.0457
		RMSE	0.0767	0.0407	0.0204
GMM2		Mean	0.5041	0.5003	0.5004
		Std Dev	0.0880	0.0653	0.0462
		RMSE	0.0787	0.0413	0.0206

As expected, our estimator performs as well if not better than Kyriazidou's (2001) estimators in terms of a smaller RMSE.

6 Conclusion

In this paper we consider the problem of identification and estimation in panel-data sample-selection models with a binary selection rule when the latent equations contain possible predetermined variables, lags of the dependent variables, and unobserved individual effects. The selection equation contains lags of the dependent variables from the latent equations and other possible predetermined variables relative to the latent equations. We derive a set of conditional moment restrictions that are then exploited to construct a three-step sieve extremum estimator for the parameters of the main equation including a nonparametric estimator of the sample-selection term. In the second step, the unknown parameters of the selection equation are consistently estimated using a transformation approach in the spirit of Berkson's minimum chi-square sieve method and a first-step kernel estimator for the selection probability. This second step estimator is of interest in its own right: it can be used to semiparametrically estimate a panel-data binary-response model with a nonparametric individual specific effect without making any other distributional assumptions. We show that both our second- and third-step estimators are \sqrt{n} -consistent and asymptotically normal. This has not been previously established for this class of dynamic sample-selection models. Our framework is also more general than the ones previously studied in that ours can estimate equations derived from an intertemporal utility maximization problem, whereas the alternative estimators in the literature cannot. The major limitation of our model is that it imposes an individual specific effect structure similar to Altug and Miller (1998) in the selection equation; however, it is still general enough since we do not make any parametric assumption about the functional form of either the mean component or the distribution function. Even this restriction can be relaxed further at the expense of making an exclusion-independence assumption. Our estimators perform well in a limited Monte Carlo study and does better than Kyriazidou's estimator in small sample.

7 Appendix

In this appendix, the letters c and C will denote diverse constants, not necessarily the same, and diverse occurrences.

Proof of Proposition 1. First, note that for $i = 1, \dots, n$, $t = 4, \dots, T$,

$$E[(\Delta y_{it} - \rho \Delta y_{it-1} - \Delta x_{it} \beta - \Delta \bar{\lambda}(\bar{v}_{it}, \bar{v}_{it-1}, \bar{v}_{it-2}, z_i^1)) \mid F_{it}] = E[\Delta \varepsilon_{it} \mid F_{it}].$$

From (28),

$$(63) \quad E[\Delta \varepsilon_{it} \mid F_{it}] = E[\varepsilon_{it}^* - \bar{\lambda}(\bar{v}_{it}, \bar{v}_{it-1}, z_i^1) \mid F_{it}] - E[\varepsilon_{it-1}^* - \bar{\lambda}(\bar{v}_{it-1}, \bar{v}_{it-2}, z_i^1) \mid F_{it}].$$

Furthermore,

$$(64) \quad E[\varepsilon_{it}^* - \bar{\lambda}(\bar{v}_{it}, \bar{v}_{it-1}, z_i^1) \mid F_{it}] = E[E[\varepsilon_{it}^* - \bar{\lambda}(\bar{v}_{it}, \bar{v}_{it-1}, z_i^1) \mid \zeta_{it}, F_{it}] \mid F_{it}]$$

$$(65) \quad = E[E[\varepsilon_{it}^* - \bar{\lambda}(\bar{v}_{it}, \bar{v}_{it-1}, z_i^1) \mid \zeta_{it}] \mid F_{it}].$$

The second equality follows from the independence of F_{it}/ζ_{it} from ε_{it}^* (see Assumptions 2.1–2.2). Finally from (23),

$$(66) \quad E[\varepsilon_{it}^* - \bar{\lambda}(\bar{v}_{it}, \bar{v}_{it-1}, z_i^1) \mid \zeta_{it}] = 0,$$

which implies that

$$(67) \quad E[\varepsilon_{it}^* - \bar{\lambda}(\bar{v}_{it}, \bar{v}_{it-1}, z_i^1) \mid F_{it}] = 0.$$

A similar reasoning on the second term of the right-hand side of (63) leads to

$$(68) \quad E[\varepsilon_{it-1}^* - \bar{\lambda}(\bar{v}_{it-1}, \bar{v}_{it-2}, z_i^1) \mid F_{it}] = 0.$$

Finally, Equations (67) and (68) prove that (63) is equal to zero. ■

Proof of Theorem 1. First, by Assumption 2.3, $\{x_{it}, z_{it}, \alpha_i, \eta_i, \varepsilon_{it}^*, u_{it}\}_{t=1}^T$ is an i.i.d sample drawn from a distribution that satisfies (1)-(3). We observe a random sample $\{(y_{i0}, d_{i0})\}_{i=1}^n$ from a distribution that satisfies (2) where d_{i0} and y_{i0}^* take values on a subset of the real line. Assumptions 2.4 and 4.4.1 imply that y_{it}^* are i.i.d over individuals. From equation (3) and by Assumptions 2.1-2.2, 4.4.1 and 4.4.2 we obtain that $d_i \equiv \{d_{it}\}_{t=1}^T$ are i.i.d over individuals. Therefore, by definition of $\omega_{it} \equiv (d_{it-1}, y_{it-1}, z'_{it}, z_i^1)'$, $\omega_i = \{\omega_{it}\}_{t=1}^T$ are i.i.d..

Let $\theta_1(\Delta d_{-1}, \Delta z, P) \equiv \phi_0 \Delta d_{-1} + \Delta z \gamma - \Delta F^{-1}(P)$, and $\hat{\theta}_1(\Delta d_{-1}, \Delta z, P) \equiv \hat{\phi}_0 \Delta d_{-1} + \Delta z \hat{\gamma} - \widehat{\Delta F^{-1}}(P)$. Let $d^2(\theta_1, \theta_{10}) = \|\theta_1 - \theta_{10}\|^2 = E\{[\phi_0 - \phi_{00}] \Delta d_{-1} + \Delta z [\gamma - \gamma_0] - [\Delta F^{-1}(P_0) - \Delta F_0^{-1}(P_0)]\}^2$. Consistency of θ_1 is established using the following decomposition:

$$\begin{aligned} \left\| \hat{\theta}_1(\hat{P}) - \theta_{10} \right\| &\leq \left\| \hat{\theta}_1(P_0) - \theta_{10} \right\| + \left\| \hat{\theta}_1(\hat{P}) - \hat{\theta}_1(P_0) \right\| \\ &\leq A + B. \end{aligned}$$

Rate of Convergence of A.

The rate of convergence of A is obtained by application of Theorem 3-2 in Chen (2005). For clarity, these conditions are recalled here:

(C05)1.1: $Q_1(\theta_1, P_0) \equiv E[\ell(\theta_1, \omega, P_0)]$ is uniquely maximized on Θ_1 at $\theta_{10} \in \Theta_1$.
(C05)1.2 : $\Theta_{1n} \subseteq \Theta_{1n+1} \subseteq \dots \subseteq \Theta_1$ for all $n \geq 1$ and for any $\theta_1 \in \Theta_1$ there exists $\pi_n \theta_1 \in \Theta_{1n}$ such that $\|\theta_1 - \pi_n \theta_1\| = o(1)$ as n gets large. (C05)1.3: The criterion $Q_1(\theta_1, P_0)$ is continuous in $\theta_1 \in \Theta_1$ with respect to $d(\cdot)$. (C05)1.4: The sieve spaces Θ_{1n} are compact under $d(\cdot)$. (C05)1.5: $\text{plim}_{n \rightarrow \infty} \sup_{\theta_1 \in \Theta_{1n}} |Q_{1n}(\theta_1, P_0) - Q_1(\theta_1, P_0)| = 0$.
(C05)1.6 : $\{\omega_i\}_{i=1}^n$ are i.i.d. (C05)1.7: There is $C_1 > 0$ such that for all small $\varepsilon > 0$, $\sup_{\theta_1 \in \Theta_{1n}} V[\ell(\theta_1, \omega, P_0) - \ell(\theta_{10}, \omega, P_0)] \leq C_1 \varepsilon^2$. (C05)1.8: For any $\delta > 0$, there exists a constant $s \in (0, 2)$ such that

$$\sup_{\theta_1 \in \Theta_{1n}} |\ell(\theta_1, \omega, P_0) - \ell(\theta_{10}, \omega, P_0)| \leq \delta^s U(\omega),$$

where $E([U(\omega)]^\gamma) \leq C_2$ for some $\gamma \geq 2$. (C05)1.9 : Let

$$\mathcal{F}_n = \{\ell(\theta_1, \omega, P_0) - \ell(\theta_{10}, \omega, P_0) : d(\theta_1, \theta_{10}) \leq \delta, \theta_1 \in \Theta_{1n}\}$$

and for some constant b , there exists $\delta_n \in (0, 1)$ such that

$$\delta_n = \inf \left\{ \delta \in (0, 1) : \frac{1}{\sqrt{n} \delta^2} \int_0^\delta \sqrt{H_{\square}(s, \mathcal{F}_n, \|\cdot\|_2)} ds \leq \text{const.} \right\}.$$

(C05)1.1 is a direct consequence of Proposition 1. (C05)1.2 is assumed in 4.2.1. (C05)1.3 follows from the properties of scalar products. Assumption 4.2.1 guaranties that the sieve spaces Θ_{1n} are compact, hence (C05)1.4 is satisfied. (C05)1.5 follows from the continuity (C05)1.3, from simple convergence of Q_{1n} and (C05)1.4. (C05)1.6 is implied by Assumption 2.3. The remaining of the proof consists in checking conditions (C05)1.7, (C05)1.8 and (C05)1.9 controlling for the rate of convergence. To simplify notations, we denote by $\Delta F^{-1} = \Delta F_u^{-1}(P_0)$ and $\Delta F_0^{-1} = \Delta F_{u_0}^{-1}(P_0)$. Note that

$$\begin{aligned} & \ell(\theta_1, \omega, P_0) - \ell(\theta_{10}, \omega, P_0) \\ &= -\frac{1}{2} \left\{ \begin{array}{l} [\Delta y_{-1} - \phi_0 \Delta d_{-1} - \Delta z \gamma + \Delta F^{-1}]' [\Delta y_{-1} - \phi_0 \Delta d_{-1} - \Delta z \gamma + \Delta F^{-1}] \\ - [\Delta y_{-1} - \phi_{00} \Delta d_{-1} - \Delta z \gamma_0 + \Delta F_0^{-1}]' [\Delta y_{-1} - \phi_{00} \Delta d_{-1} - \Delta z \gamma_0 + \Delta F_0^{-1}] \end{array} \right\} \\ &= -\frac{1}{2} [(\phi_{00} - \phi_0) \Delta d_{-1} + \Delta z (\gamma_0 - \gamma) + \Delta F^{-1} - \Delta F_0^{-1}]' \\ & \times [2 \Delta y_{-1} - (\phi_{00} + \phi_0) \Delta d_{-1} + \Delta z (\gamma_0 + \gamma) + \Delta F^{-1} + \Delta F_0^{-1}] \\ &= [(\phi_0 - \phi_{00}) \Delta d_{-1} + \Delta z (\gamma - \gamma_0) + \Delta F_0^{-1} - \Delta F^{-1}]' \\ & \times \left[(\Delta v + \Delta \xi) + \frac{(\phi_0 - \phi_{00})}{2} \Delta d_{-1} + \Delta z \frac{(\gamma - \gamma_0)}{2} + \frac{\Delta F_0^{-1} - \Delta F^{-1}}{2} \right], \end{aligned}$$

and by Assumption 4.2.2,

$$E[\ell(\theta_1, \omega, P_0) - \ell(\theta_{10}, \omega, P_0)]^2 \leq 2C \|\theta_1 - \theta_{10}\|^2 + E \left[\begin{array}{l} \left(\begin{array}{l} ((\phi_0 - \phi_{00}) \Delta d_{-1} + \Delta z (\gamma - \gamma_0) + \Delta F_0^{-1} - \Delta F^{-1})' \\ \times ((\phi_0 - \phi_{00}) \Delta d_{-1} + \Delta z (\gamma - \gamma_0) + \Delta F_0^{-1} - \Delta F^{-1}) \end{array} \right)' \\ \times \left(\begin{array}{l} ((\phi_0 - \phi_{00}) \Delta d_{-1} + \Delta z (\gamma - \gamma_0) + \Delta F_0^{-1} - \Delta F^{-1})' \\ \times ((\phi_0 - \phi_{00}) \Delta d_{-1} + \Delta z (\gamma - \gamma_0) + \Delta F_0^{-1} - \Delta F^{-1}) \end{array} \right)' \end{array} \right].$$

Let D denote the second term on the right hand side of this inequality.

$$D \leq \sup_{\theta_1 \in \Theta_{1n}} \left[\begin{aligned} & \left((\phi_0 - \phi_{00}) \Delta d_{-1} + \Delta z (\gamma - \gamma_0) + \Delta F_0^{-1} - \Delta F^{-1} \right)' \\ & \times \left((\phi_0 - \phi_{00}) \Delta d_{-1} + \Delta z (\gamma - \gamma_0) + \Delta F_0^{-1} - \Delta F^{-1} \right) \end{aligned} \right] \\ \times E \left[\begin{aligned} & \left((\phi_0 - \phi_{00}) \Delta d_{-1} + \Delta z (\gamma - \gamma_0) + \Delta F_0^{-1} - \Delta F^{-1} \right)' \\ & \times \left((\phi_0 - \phi_{00}) \Delta d_{-1} + \Delta z (\gamma - \gamma_0) + \Delta F_0^{-1} - \Delta F^{-1} \right) \end{aligned} \right].$$

Using Theorem 1 of Gabushin (1967) and Lemma 2 in Chen and Shen (1998), by Assumption 4.2.3 we have for any $p_1 > 0$, $\|\theta_1 - \theta_{10}\|_\infty \leq C \|\theta_1 - \theta_{10}\|^{2p_1/(2p_1+d)}$, where $d = \dim(\omega)$. Therefore, $D \leq C_2 \|\theta_1 - \theta_{10}\|^{2(1+2p_1/(2p_1+d))} \leq C_2 \|\theta_1 - \theta_{10}\|^2$ and (C05)1.7 is satisfied.

Moreover,

$$\begin{aligned} |\ell(\theta_1, \omega, P_0) - \ell(\theta_{10}, \omega, P_0)| \\ \leq \|\theta_1 - \theta_{10}\|_\infty \left[|\Delta v + \Delta \xi| + \frac{1}{2} (\|\theta_1\|_\infty + \|\theta_{10}\|_\infty) \right] \text{ a.s.} \\ \leq C \|\theta_1 - \theta_{10}\|^{2p_1/(2p_1+d)} \left[|\Delta v + \Delta \xi| + \frac{1}{2} (\|\theta_1\|_\infty + \|\theta_{10}\|_\infty) \right] \text{ a.s.} \end{aligned}$$

Using Lemma 2 in Chen and Shen (1998), we have that (C05)1.8 is also satisfied for $s = 2p_1/(2p_1 + d)$, $\gamma = 2$ and $U(\omega) = |\Delta v + \Delta \xi| + \frac{1}{2} (\|\theta_1\|_\infty + \|\theta_{10}\|_\infty)$.

Finally, in order for $\hat{\theta}_1$ to converge to θ_{10} at a fast rate, not only does the approximation error $d(\theta_0, \pi_n \theta_0)$ have to approach zero suitably fast, the sieve space must not be too complex. More precisely, we denote by $H_{[\cdot]}(s, \mathcal{F}_n, \|\cdot\|_2)$ the logarithm of the minimum number of closed intervals denoted by $N_{[\cdot]}(s, \mathcal{F}_n, \|\cdot\|_2)$ and of the form $\{f : g \leq f \leq h\}$ that cover \mathcal{F}_n for g and h given such that $\|h - g\|_r \leq s$. In \mathbb{R}^2 for example, these intervals are rectangles whose lower left and upper right summits are respectively g and h . Determination of the final convergence rate is obtained by setting $\delta_n \equiv d(\theta_{10}, \pi_n \theta_{10})$. To calculate δ_n from (C05)1.9, an upper bound of $H_{[\cdot]}(s, \mathcal{F}_n, \|\cdot\|_2)$ suffices.

Note that $\|\theta_{10} - \pi_n \theta_{10}\| \leq \text{const.} \|\Delta F_0^{-1} - \pi_n \Delta F_0^{-1}\|_\infty$. By Lemma 2.1 in Os-
siander (1987), $H_{[\cdot]}(s, \mathcal{F}_n, \|\cdot\|_2) \leq \log N_{[\cdot]}(s, \mathcal{H}_{1n}, \|\cdot\|_2)$. Let $C = \sqrt{E[U(\omega)]^2}$, then $H_{[\cdot]}(s, \mathcal{F}_n, \|\cdot\|_2) \leq \log N_{[\cdot]}(\frac{s}{C}, \mathcal{H}_{1n}, \|\cdot\|_\infty)$. By Lorentz (1966), $\|\Delta F_0^{-1}(P_0) - \pi_n \Delta F_0^{-1}(P_0)\| = O(k_n^{-p_1})$ for $\mathcal{H}_1 \in \Lambda_{c_1}^{p_1}$. By Lemma 2.5 of Van de Geer (2000), $\log N_{[\cdot]}(\frac{\omega}{C}, \mathcal{H}_{1n}, \|\cdot\|_\infty) \leq \text{const.} k_n \log(1 + \frac{4c_1}{\omega})$ and δ_n solves

$$\begin{aligned} \frac{1}{\sqrt{n} \delta_n^2} \int_0^{\delta_n} \sqrt{H_{[\cdot]}(s, \mathcal{F}_n, \|\cdot\|_2)} ds &\leq \frac{1}{\sqrt{n} \delta_n^2} \sqrt{k_n} \int_0^{\delta_n} \log \left(1 + \frac{4c_1}{s} \right) ds \\ &\leq \frac{1}{\sqrt{n} \delta_n^2} \sqrt{k_n} \delta_n \leq \text{const.} \end{aligned}$$

The solution is $\delta_n \asymp \sqrt{\frac{k_n}{n}}$ where \asymp means "bounded above and below". Finally, by Theorem 3.2 in Chen (2005), $\|\theta_1(P_0) - \theta_{10}\| = O_p \left(\max \left(\sqrt{\frac{k_n}{n}}, O(k_n^{-p_1}) \right) \right)$ and for

$$k_n = O\left(n^{\frac{1}{2p_1+d}}\right), \|\theta_1(P_0) - \theta_{10}\| = O_p\left(n^{-\frac{p_1}{2p_1+d}}\right).$$

Rate of convergence of B.

By Lemma 8-10 in Newey and McFadden (1994) relying on Assumptions 4.1.1-4.1.3, we have $\sqrt{n}\|\widehat{p}(\omega) - p_0(\omega)\|^2 \rightarrow 0$. By Assumptions 4.1.1-4.1.4, Newey (1994) shows that

$$\|\widehat{p} - p_0\|_S = O_p\left((\ln n)^{\frac{1}{2}}(nh^{d+2q})^{-\frac{1}{2}} + h^L\right).$$

Hence, we obtain that $\|\widehat{P} - P_0\|_S = O_p\left((\ln n)^{\frac{1}{2}}(nh^{d+2q})^{-\frac{1}{2}} + h^L\right)$, where $\|\cdot\|_S$ is the Sobolev norm defined as $\|P\|_S \equiv \max_{l \leq q} \sup_{\omega \in \mathcal{W}} \left\| \frac{\partial P(\omega)}{\partial \omega^l} \right\|$.

Finally, using the fact that π_n are uniformly Lipschitzian, we have $\|\widehat{\theta}_1(\widehat{P}) - \widehat{\theta}_1(P_0)\| = \|\pi_n \triangle F^{-1}(\widehat{P}) - \pi_n \triangle F^{-1}(P_0)\| \leq C \|\widehat{P} - P_0\|_S^{p_1}$.

Hence, if $\|P - \widehat{P}\|_S = O(k_n^{-1})$, then $\|\pi_n \triangle F^{-1}(\widehat{P}) - \pi_n \triangle F^{-1}(P)\| = O(k_n^{-p_1})$ and $\|\widehat{\theta}_1(\widehat{P}) - \widehat{\theta}_1(P_0)\| = O_p\left(n^{-\frac{p_1}{2p_1+d}}\right)$. Recall that $\|\widehat{P} - P_0\|_S = O_p\left((\ln n)^{\frac{1}{2}}(nh^{d+2q})^{-\frac{1}{2}} + h^L\right)$ and note that $(\ln n)^{\frac{1}{2}}(nh^{d+2q})^{-\frac{1}{2}} + h^L \leq (\ln n)^{\frac{1}{2}}\left[(nh^{d+2q})^{-\frac{1}{2}} + h^L\right]$.

We choose the optimal bandwidth h^* such that $(nh^{d+2q})^{-\frac{1}{2}} = h^L$, that is $h^* = n^{-\frac{1}{2L+d+2q}}$. Finally, $(n(h^*)^{d+2q})^{-\frac{1}{2}} + (h^*)^L = n^{-\frac{L}{2L+d+2q}}$, so that $\|P - \widehat{P}\| = O_p\left((\ln n)^{\frac{1}{2}}n^{-\frac{L}{2L+d+2q}}\right)$ and $\|\widehat{\theta}_1(\widehat{P}) - \widehat{\theta}_1(P_0)\| = O_p\left((\ln n)^{\frac{p_1}{2}}n^{-\frac{L}{2L+d+2q}}\right)$.

Finally, we can choose q and L large enough such that $\frac{L}{2L+d+2q} > \frac{p_1}{2p_1+d}$ and $\|\widehat{\theta}_1(\widehat{P}) - \theta_0\| = O_p\left(n^{-\frac{p_1}{2p_1+d}}\right)$. ■

Proof of Theorem 2. All computations being computed at the true value of the selection probability P_0 , we will define $l(\theta_1, \omega, P_0) \equiv l(\theta_1, \omega)$ for the remaining of the proof. For all $\theta_1 \in \Theta_1$ and all ω , there exists $l'_{\theta_0}[\theta - \theta_0, \omega]$ such that the remainder in the linear approximation is

$$r(\theta_1 - \theta_{10}, \omega) = l(\theta_1, \omega) - l(\theta_{10}, \omega) - l'_{\theta_0}[\theta - \theta_0, \omega],$$

where $l'_{\theta_0}[\theta - \theta_0, \omega] = \lim_{t \rightarrow 0} [l(\theta_1(\theta_{10}, t), \omega) - l(\theta_{10}, \omega)]/t$ is the pathwise derivative of l at θ_0 and $\theta_1(\theta_{10}, t) \in \Theta_1$ is a path in t connecting θ_{10} to θ_1 such that $\theta_1(\theta_{10}, 0) = \theta_{10}$ and $\theta_1(\theta_{10}, 1) = \theta_1$. Suppose that L_2 norm defined above induces an inner product $\langle \cdot, \cdot \rangle$ on the completion of the space spanned by $\Theta_1 - \theta_{10}$ denoted V . Let ε_n denote any sequence satisfying $\varepsilon_n = o\left(n^{-\frac{1}{2}}\right)$ and let $\nu_n(g(\omega)) = \frac{1}{n} \sum_{i=1}^n (g(\omega_i) - E_0(g(\omega_i)))$ denote the empirical

process indexed by the function $g(\cdot)$. Let $K(\theta_{10}, \theta_1) \equiv n^{-1} \sum_{i=1}^n E_0 [l(\theta_{10}, \omega_i) - l(\theta_1, \omega_i)]$.

The proof consists in checking the following conditions (see Chen, 2005).

(C05)2.1: Suppose the functional of interest f has the following smoothness properties, (i) there is a $\alpha > 0$ such that $|f(\theta_1) - f(\theta_{10}) - f'_{\theta_{10}}[\theta_1 - \theta_{10}]| = O(\|\theta_1 - \theta_{10}\|^\alpha)$ uniformly in $\theta_1 \in \Theta_{1n}$ with $\|\theta_1 - \theta_{10}\| = o(1)$; (ii) $\sup_{\{\theta_1 \in \Theta_{1n} : \|\theta_1 - \theta_{10}\| > 0\}} \frac{|f'_{\theta_{10}}[\theta_1 - \theta_{10}]|}{\|\theta_1 - \theta_{10}\|} < \infty$; (iii) there is $\pi_n v^* \in \Theta_{1n}$ such that $\|\pi_n v^* - v^*\| \|\widehat{\theta}_1 - \theta_{10}\| = o_p(n^{-\frac{1}{2}})$. (C05)2.2:

$$\sup_{\{\theta \in \Theta_n : \|\theta_1 - \theta_{10}\| \leq \delta_n\}} \nu_n(l(\theta_1, \omega) - l((\theta_1 \pm \varepsilon_n \pi_n v^*), \omega) - l'_{\theta_{10}}([\pm \varepsilon_n \pi_n v^*, \omega]) = O_p(\varepsilon_n^2).$$

$$(C05)2.3: K(\theta_{10}, \widehat{\theta}_1) - K(\theta_{10}, \widehat{\theta}_1 \pm \varepsilon_n \pi_n v^*) = \pm \varepsilon_n \langle \widehat{\theta}_1 - \theta_{10}, \pi_n v^* \rangle + o(n^{-1}). \quad (C05)2.4:$$

(i) $\nu_n(l'_{\theta_{10}}[\pi_n v^* - v^*, \omega]) = o_p(n^{-\frac{1}{2}})$; (ii) $E[l'_{\theta_{10}}[\pi_n v^*, \omega]] = o(n^{-\frac{1}{2}})$. (C05)2.5: $n^{\frac{1}{2}} \nu_n(l'_{\theta_{10}}[v^*]) \rightarrow N(0, \sigma_{v^*}^2)$ where $\sigma_{v^*}^2 \equiv \text{Var}_0(l'_{\theta_{10}}[v^*]) > 0$ for i.i.d data.

By the Riesz representation theorem, there exists $v^* \in V$ such that for any $\theta_1 \in \Theta_1$

$$f'_{\theta_{10}}[\theta_1 - \theta_{10}] = \langle \theta_1 - \theta_{10}, v^* \rangle$$

if and only if $\|g'_{\theta_{10}}\| < \infty$.

Let $f(\theta_1) = \lambda' \begin{bmatrix} \phi_0 \\ \gamma \end{bmatrix}$ where λ is a unit vector in $\mathbb{R}^{d_{\phi_0} + \gamma}$. Then, (C05)2.1 is satisfied for any arbitrary large α since

$$|f(\theta_1) - f(\theta_{10}) - f'_{\theta_{10}}[\theta_1 - \theta_{10}]| = \lambda' \begin{bmatrix} \phi_0 \\ \gamma \end{bmatrix} - \lambda' \begin{bmatrix} \phi_{00} \\ \gamma_0 \end{bmatrix} - \begin{bmatrix} \phi_0 - \phi_{00} \\ \gamma - \gamma_0 \end{bmatrix}' \lambda = 0.$$

Moreover,

$$\begin{aligned} \|f'_{\theta_{10}}\|^2 &= \sup_{\{\theta_1 \in \Theta_{1n} : \|\theta_1 - \theta_{10}\| > 0\}} \frac{\left\{ \begin{pmatrix} \phi_0 - \phi_{00} \\ \gamma - \gamma_0 \end{pmatrix}' \lambda \right\}^2}{\|\theta_1 - \theta_{10}\|^2} \\ &= \frac{\{b' \lambda\}^2}{\sup_{\left\{ (b, u_1) : \left\| b' \begin{bmatrix} \Delta d_{-1} \\ \Delta z \end{bmatrix} - u_1(P_i) \right\| > 0 \right\}} b' E \left[\left(\begin{bmatrix} \Delta d_{-1} \\ \Delta z \end{bmatrix} - u_1(P) \right) \left(\begin{bmatrix} \Delta d_{-1} \\ \Delta z \end{bmatrix} - u_1(P) \right)' \right] b} \\ &= \lambda' E [D_{u_1^*}(\omega)' D_{u_1^*}(\omega)]^{-1} \lambda = \lambda' \Sigma_*^{-1} \lambda, \end{aligned}$$

where $D_{u_1^*}(\omega) \equiv \begin{bmatrix} \Delta d_{-1} \\ \Delta z \end{bmatrix} - u_1(P)$ solves

$$\left(u_1: \left\| \begin{bmatrix} \Delta d_{-1} \\ \Delta z \end{bmatrix} - u_1(P) \right\| > 0 \right) E \left[\left(\begin{bmatrix} \Delta d_{-1} \\ \Delta z \end{bmatrix} - u_1(P) \right) \left(\begin{bmatrix} \Delta d_{-1} \\ \Delta z \end{bmatrix} - u_1(P) \right)' \right].$$

Using the definition of the norm above, we have

$$\left(\nu_{(\phi_0, \gamma)}^*, \nu_F^* \right) = \left(\Sigma_*^{-1} \lambda, -u_1^* \Sigma_*^{-1} \lambda \right).$$

Assumption 4.3 ensures that the denominator is different from zero, hence that $\|f'_{\theta_{10}}\|$ is bounded and (C05)2.1 (ii) is satisfied. Since $\|\widehat{\theta}_1 - \theta_{10}\| = O\left(n^{-\frac{p_1}{2p_1+d}}\right)$, therefore, by Assumption 4.3.3, (C05)2.1(iii) is satisfied. In order to check (C05)2.2, note that

$$l(\theta_1, \omega) - l(\theta_1 \pm \varepsilon_n \pi_n v^*, \omega) - l'_{\theta_{10}}[\pm \varepsilon_n \pi_n v^*, \omega] = r(\theta_1 - \theta_{10}, \omega) - r(\theta_1 \pm \varepsilon_n \pi_n v^* - \theta_{10}, \omega),$$

with

$$r(\theta_1 - \theta_{10}, \omega) = \left((\phi_0 - \phi_{00}) \Delta d_{-1} + \Delta z (\gamma - \gamma_0) - \Delta F^{-1}(P_0) + \Delta F_0^{-1}(P_0) \right)^2$$

and

$$r(\theta_1 \pm \varepsilon_n \pi_n v^* - \theta_{10}, \omega) = \left(\begin{array}{c} (\phi_0 - \phi_{00}) \Delta d_{-1} + \Delta z (\gamma - \gamma_0) \\ - (\Delta F^{-1} - \Delta F_0^{-1}) \pm \varepsilon_n \pi_n v^* \end{array} \right)^2.$$

Therefore,

$$l(\theta_1, \omega) - l(\theta_1 \pm \varepsilon_n \pi_n v^*, \omega) - l'_{\theta_{10}}[\pm \varepsilon_n \pi_n v^*, \omega] = -\frac{1}{2} [\mp 2(\theta_1 - \theta_{10}) \varepsilon_n \pi_n v^* - (\varepsilon_n \pi_n v^*)^2].$$

Finally, let $S_n = \{(\theta_1 - \theta_{10}) v^* : \|\theta_1 - \theta_{10}\| \leq \delta_n, \theta_1 \in \Theta_1\}$. It follows from Kolmogorov and Tihomirov (1961) and by Lemma 4 of Shen and Wong (1994) that the convergence rate of the empirical process

$$\sup_{\{\theta_1 \in \Theta_1: \|\theta_1 - \theta_{10}\| \leq \delta_n\}} n^{-\frac{1}{2}} v_n((\theta_1 - \theta_{10}) v^*)$$

is of order $O_p\left(n^{-\frac{2p_1}{2p_1+d}}\right)$, and (C05)2.2 holds for $p_1 > \frac{d}{2}$.

$$\begin{aligned} \text{Note that } l(\theta_1, \omega) - l(\theta_{10}, \omega) &= [(\phi_0 - \phi_{00}) \Delta d_{-1} + \Delta z (\gamma - \gamma_0) + \Delta F_0^{-1} - \Delta F^{-1}]' \\ &\times \left[(\Delta v + \Delta \xi) + \frac{(\phi_0 - \phi_{00})}{2} \Delta d_{-1} + \Delta z \frac{(\gamma - \gamma_0)}{2} + \frac{\Delta F_0^{-1} - \Delta F^{-1}}{2} \right]. \end{aligned}$$

Therefore, $E [\ell(\theta_1, \omega) - \ell(\theta_{10}, \omega)] = \frac{1}{2} \|\theta_1 - \theta_{10}\|^2 = K(\theta_{10}, \theta_1)$ and

$$\begin{aligned}
& K(\theta_{10}, \widehat{\theta}_1) - K(\theta_{10}, \widehat{\theta}_1 \pm \varepsilon_n \pi_n v^*) \\
&= \frac{1}{2} \langle \theta_{10} - \widehat{\theta}_1, \theta_{10} - \widehat{\theta}_1 \rangle - \frac{1}{2} \langle \theta_{10} - \widehat{\theta}_1 \pm \varepsilon_n \pi_n v^*, \theta_{10} - \widehat{\theta}_1 \pm \varepsilon_n \pi_n v^* \rangle \\
&= \pm \frac{1}{2} \left(2\varepsilon_n \langle \theta_{10} - \widehat{\theta}_1, \pi_n v^* \rangle \right) - \varepsilon_n^2 \langle \pi_n v^*, \pi_n v^* \rangle \\
&= \pm \varepsilon_n \langle \theta_{10} - \widehat{\theta}_1, \pi_n v^* \rangle + o_p(n^{-1}),
\end{aligned}$$

and (C05)2.3 is satisfied. Furthermore,

$$l'_{\theta_{10}}[\theta_1 - \theta_{10}, \omega] = (\Delta v + \Delta \xi) ((\phi_0 - \phi_{00}) \Delta d_{-1} + \Delta z (\gamma - \gamma_0) + \Delta F_0^{-1} - \Delta F^{-1}).$$

Hence the expectation of this term is equal to zero and (C05)2.4 (ii) is satisfied.

Moreover,

$$\frac{1}{n} \sum_{i=1}^n l'_{\theta_{10}}[\pi_n v^* - v^*, \omega_i] = \frac{1}{n} \sum_{i=1}^n (\Delta v_i + \Delta \xi_i) \{\pi_n v_F^* - v_F^*\}.$$

Using Chebyshev's inequality, for any real number α ,

$$P \left(\left| \frac{1}{n} \sum_{i=1}^n l'_{\theta_{10}}[\pi_n v^* - v^*, \omega_i] \right| > n^\alpha \right) \leq \frac{\sum_{i=1}^n \{\pi_n v_F^* - v_F^*\}^2 V(\Delta v_i + \Delta \xi_i | \omega)}{n^{2\alpha+2}}.$$

Therefore, by Assumption 4.2.2, (C05)2.4(i) is satisfied for $\alpha = -\frac{1}{4}$.

We know that

$$l'_{\theta_{10}}[\pi_n v^*, \omega] = (\Delta v + \Delta \xi) \left\{ \begin{pmatrix} \Delta d_{-1} \\ \Delta z \end{pmatrix} \nu_{(\phi_0, \gamma)}^* + \pi_n v_F^*(P) \right\}.$$

Hence, (C05)4.2(ii) is trivially satisfied. Using the expression above for $\nu_{(\phi_0, \gamma)}^*$ and $v_F^*(P)$, we obtain that

$$\left(\nu_{(\phi_0, \gamma)}^*, \nu_F^* \right) = (\Sigma_*^{-1} \lambda, -u_1^* \Sigma_*^{-1} \lambda)$$

and

$$\begin{aligned}
l'_{\theta_{10}}[v^*, \omega] &= (\Delta v + \Delta \xi) \left\{ \begin{pmatrix} \Delta d_{-1} \\ \Delta z \end{pmatrix} \Sigma_*^{-1} \lambda - u_1^* \Sigma_*^{-1} \lambda \right\} \\
&= (\Delta v + \Delta \xi) D_{u_1^*}(\omega) \Sigma_*^{-1} \lambda.
\end{aligned}$$

(C05)2.5 is satisfied under Assumption 4.3.1. Finally,

$$\begin{aligned}
V(l'_{\theta_{10}}[v^*, \omega]) &= E(\lambda' \Sigma_*^{-1} D_{u_1^*}(\omega)' V(\Delta v + \Delta \xi | \omega) D_{u_1^*}(\omega) \Sigma_*^{-1} \lambda), \\
&= \lambda' E[D_{u_1^*}(\omega)' D_{u_1^*}(\omega)]^{-1} E(D_{u_1^*}(\omega)' V(\Delta v + \Delta \xi | \omega) D_{u_1^*}(\omega)) \\
&\quad \times E[D_{u_1^*}(\omega)' D_{u_1^*}(\omega)]^{-1} \lambda > 0.
\end{aligned}$$

The result follows from Theorem 1 and a direct application of Theorem 4.2 in Chen (2005). ■

Proof of Theorem 3. As for the selection equation, we decompose the proof in successive stages:

$$\begin{aligned}
\left\| \widehat{\theta}_2 \left(\widehat{\theta}_1, \widehat{\Sigma}_o(\widetilde{\theta}_2) \right) - \theta_{20} \right\| &\leq \left\| \widehat{\theta}_2(\theta_{10}, \Sigma_0(\theta_{20})) - \theta_{20} \right\| \\
&\quad + \left\| \widehat{\theta}_2 \left(\widehat{\theta}_1, \Sigma_o(\theta_{20}) \right) - \widehat{\theta}_2(\theta_{10}, \Sigma_0(\theta_{20})) \right\| \\
&\quad + \left\| \widehat{\theta}_2 \left(\widehat{\theta}_1, \widehat{\Sigma}_o(\widetilde{\theta}_2) \right) - \widehat{\theta}_2 \left(\widehat{\theta}_1, \Sigma_o(\widetilde{\theta}_2) \right) \right\| \\
&\quad + \left\| \widehat{\theta}_2 \left(\widehat{\theta}_1, \Sigma_o(\widetilde{\theta}_2) \right) - \widehat{\theta}_2 \left(\widehat{\theta}_1, \Sigma_o(\theta_{20}) \right) \right\|, \\
&\leq A' + B' + C' + D'.
\end{aligned}$$

We first find the rate of convergence of A' and $\widetilde{\theta}_2$. This is done by verifying Assumptions 3.1-3.9 necessary to apply Theorem 3.1 of Ai and Chen (2003) for the respective weighting matrices $\Sigma_0(\theta_{20})$ and the identity I . For clarity, we recall these assumptions in the context of our model. We then show that the set of primitive Assumptions 4.1-4.11 are sufficient for these conditions to hold. Let us introduce the following norm:

$$(69) \quad \|\theta_2\|_s = \|(\rho, \beta)'\|_E + \sup_{v \in \mathcal{R}^3} \left| \Delta \bar{\lambda}(v) \times (1 + \|v\|_E^2)^{-\frac{a}{2}} \right| \text{ or some } a > p_2.$$

(AC03)3.1: (i) The data $\{Z_i, F_i\}_{i=1}^n$ are i.i.d.; (ii) the support of F , \mathfrak{S} is compact with a nonempty interior; (iii) the density of F is bounded and bounded away from zero; (AC03)3.2: (i) The smallest and largest eigenvalues of $E\{B^{k_{s,n}}(F)B^{k_{s,n}}(F)'\}$ are bounded and bounded away from zero for all $k_{s,n}$; (ii) for any $f(\cdot)$ with $E[f(F)^2] < \infty$, there exists a $B^{k_{s,n}}(F)'\pi$ such that $E\left\{[f(F) - B^{k_{s,n}}(F)'\pi]^2\right\} = o(1)$; (iii) for any $f(\cdot) \in \Lambda_c^\gamma(\mathfrak{S})$ with $\gamma > d_F/2$, there exists $B^{k_{s,n}}(F)'\pi \in \Lambda_c^\gamma(\mathfrak{S})$ such that

$$\sup_{F \in \mathfrak{S}} |f(F) - B^{k_{s,n}}(F)'\pi| = O(k_{s,n}^{-\gamma/d_F})$$

and $k_{s,n}^{-\gamma/d_F} = o(n^{-\frac{1}{4}})$. (AC03)3.3: $\theta_{20} \in \Theta_2$ is the only $\theta_2 \in \Theta_2$ satisfying $E[g(\theta_2, Z)|F] = 0$. (AC03)3.4: (i) $\Sigma_o(F)$ is finite positive definite uniformly over $F \in \mathfrak{S}$; (ii) $\widehat{\Sigma}_o(F) =$

$\Sigma_o(F) + o_p(n^{-\frac{1}{4}})$ uniformly over $F \in \mathfrak{F}$. (AC03)3.5: (i) There is a metric $\|\cdot\|_s$ such that $\Theta_2 \equiv \Lambda_2 \times H_2$ is compact under $\|\cdot\|_s$; (ii) for any $\theta_2 \in \Theta_2$, there exists $\pi_n \theta_2 \in \Theta_{2n} \equiv \Lambda_2 \times H_{2n}$ with $\|\pi_n \theta_2 - \theta_2\|_s = o(1)$; (iii) there is a constant $\mu_1 > 0$ such that for any $\theta_2 \in \Theta_2$, there exists $\pi_n \theta_2 \in \Theta_{2n} \equiv \Lambda_2 \times H_{2n}$ such that $\|\pi_n \theta_2 - \theta_2\|_s = O(k_{2n}^{-\mu_1})$, and $k_{2n}^{-\mu_1} = o(n^{-\frac{1}{4}})$. (AC03)3.6: (i) $E[g(\theta_{20}, Z)^2 | F]$ is bounded; (ii) $g(\theta_2, Z)$ is Hölder continuous in $\alpha \in A$; (iii) each element of $g(\theta_2, Z)$ satisfies an envelope condition over $\theta_2 \in \Theta_{2n}$ that is, there exists a measurable function $c(Z)$ with $E[c(Z)^4] < \infty$ such that $|g(\theta_2, Z)| \leq c(Z)$ for all Z and $\theta_2 \in \Theta_{2n}$; (iv) each element of $m(\cdot, \theta_2) \equiv E(g(\cdot, \theta_2) | F) \in \Lambda_c^{\gamma}(\mathfrak{F})$ with $\gamma > d_F/2$ for all $\theta_2 \in \Theta_{2n}$. (AC03)3.7: (i) $d_g k_{s,n} \geq d_{1+\dim(x)} + k_{2n}$, $k_{2n} \rightarrow \infty$ and $k_{s,n}/n \rightarrow 0$. Let $\xi_{0n} \equiv \sup_{F \in \mathfrak{F}} \|B^{k_{s,n}}(F)\|_E$, which is nondecreasing in $k_{s,n}$. Denote $N(r, \Theta_2, \|\cdot\|_s)$ as the minimal number of radius r covering balls of H_{2n} under $\|\cdot\|_s$. (ii) $k_{2n} \times \ln n \times \xi_{0n}^2 \times n^{-\frac{1}{2}} = o(1)$. (AC03)3.8: $\ln[N(\varepsilon^{1/\kappa}, \Theta_{2n}, \|\cdot\|_s)] \leq \text{const.} \times k_{2n} \times \ln(k_{2n}/\varepsilon)$. (AC03)3.9: (i) Θ_2 is convex in θ_{20} , and $g(\theta_2, Z)$ is pathwise differentiable at θ_{20} ; (ii) for some $C_1, C_2 > 0$

$$C_1 E\{g(F, \theta_2)' \Sigma(F)^{-1} g(F, \theta_2)\} \leq \|\theta_2 - \theta_{20}\|^2 \leq C_2 E\{g(F, \theta_2)' \Sigma(F)^{-1} g(F, \theta_2)\}$$

for all $\theta_2 \in \Theta_{2n}$ with $\|\theta_2 - \theta_{20}\|_s = o(1)$.

As is standard we need to show that the criterion function is continuous in the first step estimator. This is done by adding a boundedness condition in 4.4.7 that ensures that our moment conditions are continuous in the preliminary estimates of the selection equation. Given the definition of Z_{is} and F_{is} , the data $\{Z_i, F_i\}_{i=1}^n$ are i.i.d (see proof of Theorem 1) and (AC03)3.1(i) is satisfied. Assumptions 4.4.2, 4.4.3, 4.4.4 directly imply that the support of $\{F_i\}_{i=1}^n$ is compact with a nonempty interior and (AC03)3.1 (ii) is satisfied. (AC03)3.1(iii) is directly assumed in Assumption 4.4.4. (AC03)3.2(i) is directly assumed by Assumption 4.4.7. Note that (AC03)3.2(i) is satisfied with the linear sieves satisfying Assumption 4.4.5 (see Newey (1997) for details). Note that 3.2 (ii) is implied by Assumptions 4.4.5, while (AC03)3.2(iii) is implied by Assumptions 4.4.5 and 4.4.6 (ii). Propositions 2 and 3 ensure that θ_2 is identified, hence (AC03)3.3 is satisfied. Note that (AC03)3.4(ii) and (iii) are trivially satisfied with an identity weighting matrix or for the true value Σ_0 . Otherwise Assumptions 4.4.8 (ii) and (iii) are sufficient to satisfy this condition. By Assumptions 4.4.3 and 4.4.9, Θ_2 is compact under the norm $\|\theta_2\|_s = \|(\rho, \beta)\|_E + \sup_{v \in \mathcal{R}^{3+\dim(z^1)}} |\Delta \bar{\lambda}(v) \times (1 + \|v\|_E^2)^{-\frac{a}{2}}|$ for some $a > p_2$. From Chen, Hansen and Scheinkman (1997) for any $\theta_2 \in \Theta_2$, $\|\pi_n \theta_2 - \theta_2\|_s = \|\pi_n \Delta \bar{\lambda} - \Delta \bar{\lambda}\|_{\infty, \varpi} = \sup_v |[\pi_n \Delta \bar{\lambda}(v) - \Delta \bar{\lambda}(v)] \varpi(v)| \leq C(k_{2n})^{-\frac{p_2}{3+\dim(z^1)}}$, where $\varpi(v) = (1 + \|v\|_E^2)^{-\frac{a}{2}}$ for some $a > p_2$. Hence (AC03)3.5(i) is satisfied by Assumption 4.4.6. By the same argument (AC03)3.5(iii) is satisfied with $\mu_1 = p_2/(3 + \dim(z^1))$. (AC03)3.6(i) is satisfied by As-

sumptions 4.4.2, 4.4.3 and 4.4.8 while for any $\theta_{21}, \theta_{22} \in \Theta_2$,

$$\begin{aligned}
|g(\theta_{21}, Z) - g(\theta_{22}, Z)| &\leq |dd_{-1}\Delta y_{-1}(\rho_1 - \rho_2)| + |dd_{-1}\Delta x(\beta_1 - \beta_2)| \\
&\quad + |dd_{-1}(\Delta\bar{\lambda}_1(v) - \Delta\bar{\lambda}_2(v))| \\
&\leq \|dd_{-1}\Delta y_{-1}\|_E \|\rho_1 - \rho_2\|_E + \|dd_{-1}\Delta x\|_E \|\beta_1 - \beta_2\|_E \\
&\quad + [\varpi(v)]^{-1} \|dd_{-1}\|_\infty \|\Delta\bar{\lambda}_1 - \Delta\bar{\lambda}_2\|_{\infty, \varpi} \\
&\leq C(Z) \|\theta_{21} - \theta_{22}\|_s,
\end{aligned}$$

where $C(Z) = \|dd_{-1}\Delta y_{-1}\|_E + \|dd_{-1}\Delta x\|_E + [\varpi(v)]^{-1} \|dd_{-1}\|_\infty$ and $E[C^2(Z)|F] < \infty$. Therefore, (AC03)3.6(ii) for $\kappa = 1$ by Assumptions 4.4.2, 4.4.9 and 4.4.11(ii). For any $\theta_2 \in \Theta_{2n} \equiv \Theta \times \mathcal{H}_{2n}$, with \mathcal{H}_{2n} , being the sieve space for the tensor product of Fourier or power series of the form

$$(70) \quad \mathcal{H}_{2n} = \left\{ \begin{array}{l} \Delta\bar{\lambda}(\bar{v}, \bar{v}_{-1}, \bar{v}_{-2}, z^1) = q^{k_{2n}} (\bar{v}, \bar{v}_{-1}, \bar{v}_{-2}, z^1)' \delta \\ \text{for all } \delta \text{ satisfying } \|\Delta\bar{\lambda}\|_{\Lambda^{p_2}} \leq c \end{array} \right\},$$

we have

$$\begin{aligned}
\sup_{\theta_2 \in \Theta_{2n}} |g(Z, \theta_2)| &\leq |dd_{-1}\Delta y| + |dd_{-1}\Delta y_{-1}| + \sup_{\beta \in \Lambda_2} |dd_{-1}\Delta x\beta| \\
&\quad + |dd_{-1}| \sup_{\Delta\bar{\lambda} \in \mathcal{H}_n} |\Delta\bar{\lambda}(v)| \leq C(Z)
\end{aligned}$$

with $C(Z) = |\Delta y| + |\Delta y_{-1}| + \sup_{\beta} |\Delta x\beta| + \sup_{\Delta\bar{\lambda} \in \mathcal{H}_n} |\Delta\bar{\lambda}(v)|$. Hence (AC03)3.6(iii) is satisfied by Assumptions 2.2, 4.4.11(i) and 4.4.2-4.4.4. Note that

$$\begin{aligned}
m(F, \theta_2) &= E[dd_{-1}\Delta y | F] - \rho E[dd_{-1}\Delta y_{-1} | F] \\
&\quad - E[dd_{-1}\Delta x | F]\beta - E[dd_{-1}\Delta\bar{\lambda}(v) | F].
\end{aligned}$$

Therefore (AC03)3.6(iv) is satisfied by Assumption 4.4.10. The first part of (AC03)3.7(i) directly assumed by Assumption 4.4.6(i), (AC03)3.7(ii) is trivially implied by Assumption 4.4.5 and 4.4.6(ii). Consider the case of power series. From Newey (1997), if a) $\text{var}(g(\theta_{20}, Z)|F)$ is bounded; b) The support of F is compact with nonempty interior with probability density bounded away from zero and c) $m(F, \theta) \in \Lambda_c^\gamma(F)$, then $\xi_{0n} \leq k_{s,n}$. Note that a) is satisfied by Assumption 4.4.8, b) was shown above when proving (AC03)3.1(ii) while c) was just verified by (AC03)3.6(ii). We now have

$$k_{2n} \times \ln n \times \xi_{0n}^2 \times n^{-\frac{1}{2}} \leq k_{2n} \times \ln n \times k_{s,n}^2 \times n^{-\frac{1}{2}} = o(1)$$

where the final equality comes from Assumption 4.4.5. In the case of Fourier Series (see Ai and Chen (2003) pp. 1807), $\xi_{0n} = k_{s,n}^{\frac{1}{2}}$, therefore

$$k_{2n} \times \ln n \times \xi_{0n}^2 \times n^{-\frac{1}{2}} = k_{2n} \times \ln n \times k_{s,n} \times n^{-\frac{1}{2}} = o(1)$$

where the last equality comes from Assumption 4.4.5. From Shen and Chen (1998) and Shen and Wong (1994), for $\kappa = 1$ and the linear sieves space in Assumption 4.4.5,

$$\ln[N(\varepsilon, \Theta_{2n}, \|\cdot\|_s)] \leq \text{const.} \times \ln(1/\varepsilon) \leq \text{const.} \times k_{2n} \times \ln(k_{2n}/\varepsilon)$$

Note that $g(\theta_2, Z)$ is linear and continuous in θ_{20} , hence differentiable. As noted in Ai and Chen (2003), when $g(Z, \theta_2)$ is linear in θ_2 then $\|\theta_2 - \theta_{20}\|^2 = E\{m(F, \theta_2)' \Sigma(X)^{-1} m(F, \theta_2)\}$. Assumption 3.9(ii) is satisfied and by Theorem 3.1 in Ai and Chen (2003), we have $\|\theta_2(\Sigma_0(\theta_{20})) - \theta_{20}\| = o_p(n^{-\frac{1}{4}})$, $\|\tilde{\theta}_2(I) - \theta_{20}\| = o_p(n^{-\frac{1}{4}})$. By Assumption 4.4.1-4.4.8 and application of Lemma A1(A and B) in Ai and Chen (2003), we obtain

$$\widehat{\Sigma}_o(F, \theta_2) = \Sigma_o(F, \theta_2) + o_p(n^{-\frac{1}{4}}) \text{ uniformly over } F \in \mathfrak{F}, \theta_2 \in \mathcal{N}_{on}.$$

Finally for $\Theta_{2n} \subseteq \Theta_{2n+1} \subseteq \dots \subseteq \Theta_2$ with $\Theta_2 \in \Lambda_{c_2}^{p_2}$, we obtain

$$\begin{aligned} \left\| \widehat{\theta}_2 \left(\widehat{\theta}_1, \widehat{\Sigma}_o(\tilde{\theta}_2) \right) - \theta_{20} \right\| &= o_p \left(n^{-\frac{1}{4}} \right) + O_p \left(n^{-\frac{p_2 p_1}{2p_1 + d}} \right) + o_p \left(n^{-\frac{p_2}{4}} \right) + o_p \left(n^{-\frac{p_2 p_\Sigma}{4}} \right) \\ &= o_p \left(n^{-\frac{1}{4}} \right), \end{aligned}$$

for $p_2 > 1$, $p_\Sigma \geq 1$ and $p_1 > \frac{d}{2}$, where the first equality follows from the fact that $\|\widehat{\theta}_1 - \theta_{10}\| \leq \|\widehat{\theta}_1 - \theta_{10}\|_2$ for all $\theta_1 \in \Theta_1$. ■

Proof of Theorem 4. This theorem is proved by first verifying Assumptions 4.1–4.6 of Theorem 4.1 in Ai and Chen (2003) given that all the conditions in Theorem 3 are satisfied: (AC03)4.1: (i) $E\{D_{u_2^*}(F)' \Sigma_0(F)^{-1} D_{u_2^*}(F)\}$ is positive definite; (ii) $(\rho_0, \beta_0) \in \text{int}(\Lambda_2)$; (iii) $\Sigma_0(F) \equiv \text{var}[g(\theta_{20}, Z)|F]$ is positive definite for all $F \in \mathcal{F}$. (AC03)4.2: There is $v_n^* = (v_{\theta_2}^*, -\pi_n u_2^* v_{\theta_2}^*) \in \Theta_{2n} - \theta_{20}$ such that $\|v_n^* - v^*\| = O(n^{-\frac{1}{4}})$. Let $\mathcal{N}_{on} \equiv \left\{ \theta_2 \in \Theta_{2n} : \|\theta_2 - \theta_{20}\|_s = o(1), \|\theta_2 - \theta_{20}\| = o(n^{-\frac{1}{4}}) \right\}$ and \mathcal{N}_o is defined the same way with Θ_{2n} replaced by Θ_2 . Denote $\frac{dg(\theta_2, Z)}{d\theta_2} \equiv \frac{dg(Z, \theta_2 + \tau v)}{d\tau}$ a.s. Z and $\frac{dm(F, \theta_2)}{d\theta_2} [v] \equiv E \left\{ \frac{dg(\theta_2, Z)}{d\theta_2} [v] | F \right\}$ a.s. F . (AC03)4.3: (i) For all $\theta_2 \in \mathcal{N}_o$, the pathwise derivative $\frac{dg(Z, \theta_2(t))}{d\theta_2} [v]$ exists a.s. $Z \in Z$. Moreover, each element of $\frac{dg(Z, \theta_2)}{d\theta_2} [v_n^*]$ satisfies an envelope condition and is Hölder continuous in $\theta_2 \in \mathcal{N}_{on}$; (ii) each element of $\frac{dm(F, \theta_2)}{d\theta_2} [v_n^*]$ is in Λ_c^γ with $\gamma > \frac{d_F}{2}$ for all $\theta_2 \in \mathcal{N}_o$; (AC03)4.4: Uniformly over $\theta_2 \in \mathcal{N}_{on}$,

$$E \left[\left\| \frac{dm(F, \theta_2)}{d\theta_2} [v_n^*] - \frac{dm(F, \theta_{20})}{d\theta_2} [v_n^*] \right\|_E^2 \right] = o \left(n^{-\frac{1}{2}} \right).$$

(AC03)4.5: Uniformly over $\theta_2 \in \mathcal{N}_o, \bar{\theta}_2 \in \mathcal{N}_{on}$,

$$E \left[\begin{array}{c} \left\{ \frac{dm(F, \theta_{20})}{d\theta_2} [v^*] \right\}' \Sigma_0(F)^{-1} \\ \times \left\{ \frac{dm(F, \theta_2)}{d\theta_2} [\bar{\theta}_2 - \theta_{20}] - \frac{dm(F, \theta_{20})}{d\theta_2} [\bar{\theta}_2 - \theta_{20}] \right\} \end{array} \right] = o\left(n^{-\frac{1}{2}}\right).$$

(AC03)4.6: For all $\theta_2 \in \mathcal{N}_{on}$, the pathwise second derivative $\left. \frac{d^2 g(Z, \theta_2 + \tau v_n^*)}{d\tau^2} \right|_{\tau=0}$ exists a.s. $Z \in \mathcal{Z}$ and is bounded by a measurable function $c(Z)$ with $E[c^2(Z)] < \infty$.

Assumption (AC03)4.1(i) is implied by Proposition 3 and Assumption 2.4.2. Assumptions (AC03)4.1(ii) and (iii) are assumed directly in Assumptions 4.5.1 and 4.4.8. Assumption (AC03)4.2 is satisfied with Assumption 4.4.10 and Assumption 4.5.2(i). Assumption (AC03)4.3(i) is implied by Assumptions 4.4.1-4.4.4 and 4.5.2(i). Assumption (AC03)4.3(ii) is satisfied with Assumptions 4.5.2(i) and 2.4.2. Conditions (AC03)4.4 and (AC03)4.5 are trivially satisfied and, since the second derivative of $g(\theta_2, Z_i)$ is always zero, (AC03)4.6 is automatically satisfied. This implies that we can write

$$\begin{aligned} \sqrt{n}[(\hat{\rho}, \hat{\beta})' - (\rho_0, \beta_0)'] &= \left\{ -\frac{1}{\sqrt{n}} \sum_{i=1}^n D_{u_2^*}(F_i)' \Sigma_0(F_i)^{-1} g(\theta_{20}(\hat{\phi}_0, \hat{\gamma}), Z_i) \right\} \\ (71) \quad &\times E\{D_{u_2^*}(F_i)' \Sigma_0(F_i)^{-1} D_{u_2^*}(F_i)\}^{-1} + o_p(1). \end{aligned}$$

Taking a first-order Taylor series expansion of $g(\theta_{20}(\hat{\phi}_0, \hat{\gamma}))$ around (ϕ_{00}, γ_0) gives us

$$\begin{aligned} \sqrt{n}[(\hat{\rho}, \hat{\beta})' - (\rho_0, \beta_0)'] &= -\frac{1}{\sqrt{n}} \sum_{i=1}^n D_{u_2^*}(F_i)' \Sigma_0(F_i, \theta_{20})^{-1} \{g(\theta_{20}(\phi_{00}, \gamma_0), Z_i) \\ &+ \frac{\partial g(\theta_{20}(\phi_{00}, \gamma_0), Z_i)}{\partial \theta_{20}} \frac{\partial \theta_{20}(\phi_{00}, \gamma_0)}{\partial (\phi_0, \gamma)'} [(\hat{\phi}_0, \hat{\gamma})' - (\phi_{00}, \gamma_0)']\} \\ &\times E\{D_{u_2^*}(F_i)' \Sigma_0(F_i, \theta_{20})^{-1} D_{u_2^*}(F_i)\}^{-1} \\ &- \frac{1}{\sqrt{n}} \sum_{i=1}^n D_{u_2^*}(F_i)' \Sigma_0(F_i, \theta_{20})^{-1} o\left(\|(\hat{\phi}_0, \hat{\gamma})' - (\phi_{00}, \gamma_0)'\|\right) \\ (72) \quad &\times E\{D_{u_2^*}(F_i)' \Sigma_0(F_i, \theta_{20})^{-1} D_{u_2^*}(F_i)\}^{-1} + o_p(1). \end{aligned}$$

From Theorem 2,

$$(73) \quad \sqrt{n} [(\hat{\phi}_0, \hat{\gamma})' - (\phi_{00}, \gamma_0)'] = \frac{1}{\sqrt{n}} \sum_{i=1}^n (\Delta v_i) D_{u_1^*}(\omega_i) E\{D_{u_1^*}(\omega_i)' D_{u_1^*}(\omega_i)\}^{-1} + o_p(1).$$

Hence, we have

$$(74) \quad \frac{1}{n} \sum_{i=1}^n D_{u_2^*}(F_i)' \Sigma_0(F_i)^{-1} \frac{\partial g(\theta_{20}(\phi_{00}, \gamma_0), Z_i)}{\partial \theta'_{20}} \frac{\partial \theta_{20}(\phi_{00}, \gamma_0)}{\partial (\phi_0, \gamma)'} \\ \longrightarrow E \left[D_{u_2^*}(F_i)' \Sigma_0(F_i)^{-1} \frac{\partial g(\theta_{20}(\phi_{00}, \gamma_0), Z_i)}{\partial \theta'_{20}} \frac{\partial \theta_{20}(\phi_{00}, \gamma_0)}{\partial (\phi_0, \gamma)'} \right].$$

In our framework,

$$(75) \quad \frac{\partial g(\theta_{20}(\phi_{00}, \gamma_0), Z_i)}{\partial \theta'_{20}} \frac{\partial \theta_{20}(\phi_{00}, \gamma_0)}{\partial (\phi_0, \gamma)'} = d_i d_{i-1} \left\{ \frac{\partial \bar{\lambda}_0(\bar{v}_i, \bar{v}_{i-1}, z_i^1)}{\partial \bar{v}_i} (d_{i-1}, z_i)' \right. \\ \left. + \frac{\partial \bar{\lambda}_0(\bar{v}_i, \bar{v}_{i-1}, z_i^1)}{\partial \bar{v}_i} (d_{i-2}, z_{i-1})' \right\}.$$

Therefore by substituting Equations (73), (74) and (75) into Equation (72) and applying a standard CLT for i.i.d. data and the Generalized Slutsky's Theorem, we obtain Theorem 4. ■

Acknowledgement *We are grateful to the Editor and to three anonymous referees, whose comments significantly improved the content of this paper. We also thank Arup Bose, Mehmet Caner, David DeJong, Vilmos Komornik, Robert Miller, Whitney Newey, and Jean-Francois Richard for useful comments on the early draft of this paper. We would also like to thank seminar participants at Carnegie-Mellon, Cincinnati, Pittsburgh, SUNY-Stony Brook, Wisconsin-Madison, the 2003 North American Summer Meetings of the Econometric Society, the 2003 European Summer Meetings of the Econometric Society, and Semiparametrics in Rio 2004. Gayle's research is partially supported by the Andrew Mellon Research Fellowship. Viauoux's research is partially supported by Charles Phelps Taft research grants.*

References

- [1] Ahn, H., and J. L. Powell (1993): "Semiparametric Estimation of Censored Selection Models with a Nonparametric Selection Mechanism," *Journal of Econometrics*, 58, 3–29.
- [2] Ahn, H., and P. Schmidt (1995): "Efficient Estimation of Models for Dynamic Panel Data," *Journal of Econometrics*, 68, 5–27.
- [3] Ai, C., and X. Chen (2003): "Efficient Estimation of Models with Conditional Moment Restrictions Containing Unknown Functions," *Econometrica*, 6, 1795–1843.

- [4] Altug, S. and R. A. Miller (1998), “The Effect of Work Experience on Female Wages and Labour Supply,” *Review of Economic Studies*, 45–85.
- [5] Amemiya, T. (1994): *Advanced Econometrics*. Cambridge, Harvard University Press.
- [6] Anderson, T. W., and C. Hsiao (1982): “Formulation and Estimation of Dynamic Models Using Panel Data,” *Journal of Econometrics*, 18, 47–82.
- [7] Arellano, M., and S. R. Bond (1991): “Some Tests of Specifications for Panel Data: Monte Carlo Evidence and an Application to Employment Equations,” *Review of Economic Studies*, 58, 277–297.
- [8] Arellano, M., and O. Bover (1995): “Another Look at the Instrumental Variable Estimation of Error-Component Models,” *Journal of Econometrics*, 68, 29–51.
- [9] Arellano, M., and R. Carrasco (2003): “Binary Choice Panel Data Models with Predetermined Variables,” *Journal of Econometrics*, 115, 125–157.
- [10] Blundell, R., Chen, X., Kristensen, D. (2004): “Semi-Nonparametric IV Estimation of Shape-Invariant Engel Curves,” CEMMAP Working Paper no CWP15/03, Institute for Fiscal Studies, University College of London.
- [11] Bond, S., and C. Meghir (1994): “Dynamic Investment Models and the Firm’s Financial Policy,” *Review of Economic Studies*, 61, 197–222.
- [12] Chamberlain, G. (1986): “Asymptotic Efficiency in Semiparametric Models with Censoring,” *Journal of Econometrics*, 32, 189–218.
- [13] Chen, S. (1998): “Root N Consistent Estimation of a Panel Data Sample Selection Model,” Unpublished Manuscript, The Hong Kong University of Science and Technology.
- [14] Chen, X. (2005): “Large Sample Sieve Estimation of Semi-Nonparametric Models,” in *Handbook of Econometrics*, Vol. 6, eds. J. Heckman and E. Leamer. Elsevier Science. Forthcoming.
- [15] Chen, X., Hansen L. P., Sheinkman, J. (1997): “Shape-Preserving Estimation of Diffusions,” *Unpublished Manuscript*, University of Chicago.
- [16] Chen, X., Linton, O., Van Keilegom, I. (2003): “Estimation of Semiparametric Models when the Criterion Function is not Smooth,” *Econometrica*, 71, 1591–1608.
- [17] Chen, X., and X., Shen (1998): “Sieve Extremum Estimates for Weakly Dependent Data,” *Econometrica*, 66, 289–314.

- [18] Cosslett, S. R. (1991): “Distribution-Free Estimator of a Regression Model with Sample Selectivity,” in *Nonparametric and Semiparametric Methods in Econometrics*, eds. W. A. Barnett, J. L. Powell and G. Tauchen. Cambridge, Cambridge University Press, 175–197.
- [19] Darolles, S., Florens, J. P., Renault, E. (2003): “Nonparametric Instrumental Regression, IDEI Working Paper 228, University of Toulouse.
- [20] Das, M., W. K. Newey and F. Vella (2003): “Nonparametric Estimation of Sample Selection Models,” *Review of Economic Studies*, 70, 33–58.
- [21] Gabushin, V. N. (1967): “Inequalities for norms of functions and their derivatives in L_p Metric,” *Matematicheskie Zametki*, 1, 291–298.
- [22] Gayle, G. L., and R. A. Miller (2004): “Life-cycle Fertility Behavior and Human Capital Accumulation,” Working Paper no. 2004-E, Carnegie Mellon University.
- [23] Hahn (1997): “Bias Corrected Instrumental Variables Estimation for Dynamic Panel Models with Fixed Effects,” Working Paper 2005-024, Boston University.
- [24] Heckman, J. (1974): “Shadow Prices, Market Wages and Labor Supply,” *Econometrica*, 42, 679–694.
- [25] Heckman, J. (1976): “Common Structures of Statistical Models of Truncations Sample Selection and Limited Dependent Variables, and a Simple Estimator for Such Models,” *Annals of Economic and Social Measurement*, 15, 475–492.
- [26] Heckman, J. (1976): “Sample Selection Bias as a Specification Error,” *Econometrica*, 47, 153–161.
- [27] Honore, B. (1993): “Orthogonality conditions for Tobit Models with Fixed Effects and Lagged Dependent Variables,” *Journal of Econometrics*, 29, 35–61.
- [28] Honore, B. and E. Kyriazidou (2000): “Panel Data Discrete Choice Models with Lagged Dependent Variables,” *Econometrica*, 68, 839–874.
- [29] Honore, B. E., and A. Lewbel (2002): “Semiparametric Binary Choice Panel Data Models Without Strictly Exogenous Regressors,” *Econometrica*, 70, 2053–2063.
- [30] Hotz, V. J., F. Kydland, and G. Sedlacek (1988): “Intertemporal Preferences and Labor Supply,” *Econometrica*, 56, 335–360.
- [31] Ichimura, H. (1993): “Semiparametric Least Squares (SLS) and weighted SLS Estimation of Single Index Models,” *Journal of Econometrics*, 58, 71–120.
- [32] Kydland, F. E. and E. C. Prescott (1982): “Time to Build and Aggregate Fluctuations,” *Econometrica*, 50, 1345–1370.
- [33] Kyriazidou, E. (1997): “Estimation of Panel Data Sample Selection Model,” *Econometrica*, 65, 1335–1364.

- [34] Kyriazidou, E. (2001): "Estimation of Dynamic Panel Data Sample Selection Model," *Review of Economic Studies*, 68, 543–572.
- [35] Kolmogorov, A. N., Tihomirov, V. M. (1961): " ε -entropy and ε -capacity sets in function spaces," in American Mathematical Society Translations, 2, 227-304.
- [36] Lorentz, G (1966): *Approximation of functions*. New York, Holt.
- [37] Manski, C.F. (1987): "Semiparametric Analysis of Random Effects Linear Models from Binary Panel Data," *Econometrica*, 55, 357–362.
- [38] MaCurdy, T. E. (1981): "An Empirical Model of Labor Supply in a Life-Cycle Setting," *Journal of Political Economy*, 89, 1059-1085.
- [39] Newey, W. K. (1994): "The Asymptotic Variance of Semiparametric Estimators," *Econometrica*, 16, 1–32.
- [40] Newey, W. K. (1994): "Kernel Estimation of Partial Means and A General Variance Estimator," *Econometric Theory*, 10, 233-253.
- [41] Newey, W. K. (1997): "Convergence Rates and Asymptotic Normality for Series Estimators," *Journal of Econometrics*, 79, 147-168.
- [42] Newey, W. K. (1999): "Two-Step Series Estimation of Sample Selection Models," Working Paper, Department of Economics, MIT.
- [43] Newey, W. K. and D. McFadden (1994): "Large Sample Estimation and Hypothesis Testing," in *Handbook of Econometrics*, eds. R. F. Engle and D. L. McFadden. Amsterdam: North-Holland. 2111–2245.
- [44] Newey, W. K. and J. L. Powell (2003): Instrumental Variable Estimation of Non-parametric Models, *Econometrica*, 71, 5, 1565-1578.
- [45] Nijman, T. and M. Verbeek (1992): "Nonresponse in Panel Data: The Impact on Estimates of a Life Cycle Consumption Function," *Journal of Applied Econometrics*, 7, 243–257.
- [46] Ossiander, M. (1987): "A central limit theorem under metric entropy with L_2 bracketing," *The Annals of Probability*, 15, 897-919.
- [47] Robinson, P. (1988): "Root N Consistent Semiparametric Regression," *Econometrica*, 56, 931–954.
- [48] Shen, X. (1997): "On Methods of Sieves and Penalization," *The Annals of Statistics*, 25, 2555–2591.
- [49] Shen, X. and W. H. Wong (1994): "Convergence Rate of Sieve Estimates," *Annals of Statistics*, 22, 580–615.
- [50] Timan, A.F. (1963): *Theory of Approximation of Functions of a Real Variable*. McMillan, New York.

- [51] Van de Geer, S. (2000): *Empirical Processes in M-estimation*. Cambridge, Cambridge University Press.
- [52] Zabel, J. (1992): “Estimating Fixed and Random Effects Model with Selectivity,” *Economics Letters*, 40, 269–272.