

NBER WORKING PAPER SERIES

IDENTIFICATION AND ESTIMATION OF 'IRREGULAR' CORRELATED RANDOM
COEFFICIENT MODELS

Bryan S. Graham
James Powell

Working Paper 14469
<http://www.nber.org/papers/w14469>

NATIONAL BUREAU OF ECONOMIC RESEARCH
1050 Massachusetts Avenue
Cambridge, MA 02138
November 2008

We would like to thank seminar participants at UC - Berkeley, UCLA, USC, Harvard, Yale, NYU, Princeton, Rutgers, Syracuse, Penn State, members of the Berkeley Econometrics Reading Group and participants in the Conference in Economics and Statistics in honor of Theodore W. Anderson's 90th Birthday (Stanford University), the Copenhagen Microeconometrics Summer Workshop and the JAE Conference on Distributional Dynamics (CEMFI, Madrid) for comments and feedback. Discussions with Manuel Arellano, Stéphane Bonhomme, Gary Chamberlain, Iván Fernández-Val, Jinyong Hahn, Jerry Hausman, Bo Honoré, Michael Jansson, Roger Klein, Arthur Lewbel, Ulrich Müller, John Strauss, and Edward Vytlačil were helpful in numerous ways. This revision has also benefited from the detailed comments of a co-editor as well as three anonymous referees. Max Kasy and Alex Poirier provided research assistance. Financial support from the National Science Foundation (SES #0921928) is gratefully acknowledged. All the usual disclaimers apply. The views expressed herein are those of the author(s) and do not necessarily reflect the views of the National Bureau of Economic Research.

NBER working papers are circulated for discussion and comment purposes. They have not been peer-reviewed or been subject to the review by the NBER Board of Directors that accompanies official NBER publications.

© 2008 by Bryan S. Graham and James Powell. All rights reserved. Short sections of text, not to exceed two paragraphs, may be quoted without explicit permission provided that full credit, including © notice, is given to the source.

Identification and Estimation of 'Irregular' Correlated Random Coefficient Models
Bryan S. Graham and James Powell
NBER Working Paper No. 14469
November 2008. "Tgxkugf "Qexqdt"4232
JEL No. C14,C23,I1,O1,O15

ABSTRACT

In this paper we study identification and estimation of a correlated random coefficients (CRC) panel data model. The outcome of interest varies linearly with a vector of endogenous regressors. The coefficients on these regressors are heterogenous across units and may covary with them. We consider the average partial effect (APE) of a small change in the regressor vector on the outcome (cf., Chamberlain, 1984; Wooldridge, 2005a). Chamberlain (1992) calculates the semiparametric efficiency bound for the APE in our model and proposes a \sqrt{N} consistent estimator. Nonsingularity of the APE's information bound, and hence the appropriateness of Chamberlain's (1992) estimator, requires (i) the time dimension of the panel (T) to strictly exceed the number of random coefficients (p) and (ii) strong conditions on the time series properties of the regressor vector. We demonstrate irregular identification of the APE when $T = p$ and for more persistent regressor processes. Our approach exploits the different identifying information in the subpopulations of 'stayers' — or units whose regressor values change little across periods — and 'movers' — or units whose regressor values change substantially across periods. We propose a feasible estimator based on our identification result and characterize its large sample properties. While irregularity precludes our estimator from attaining parametric rates of convergence, its limiting distribution is normal and inference is straightforward to conduct. Standard software may be used to compute point estimates and standard errors. We use our methods to estimate the average elasticity of calorie consumption with respect to total outlay for a sample of poor Nicaraguan households.

Bryan S. Graham
New York University
19 W 4th Street, 6FL
New York, NY 10012
and NBER
bsg1@nyu.edu

James Powell
UC, Berkeley
Department of Economics
508-1 Evans Hall #3880
Berkeley, CA 94720-3880
powell@econ.berkeley.edu

That the availability of multiple observations of the same sampling unit (e.g., individual, firm, etc.) over time can help to control for the presence of unobserved heterogeneity is both intuitive and plausible. The inclusion of unit-specific intercepts in linear regression models is among the most widespread methods of ‘controlling for’ omitted variables in empirical work (e.g., Card, 1996). The appropriateness of this modelling strategy, however, hinges on any time-invariant correlated heterogeneity entering the outcome equation additively. Unfortunately, additivity, while statistically convenient, is difficult to motivate economically (cf., Imbens, 2007).² Browning and Carro (2007) present a number of empirical panel data examples where non-additive forms of unobserved heterogeneity appear to be empirically relevant.

In this paper we study the use of panel data for identifying and estimating what is arguably the simplest statistical model admitting nonseparable heterogeneity: the *correlated random coefficients* (CRC) model. Let $\mathbf{Y} = (Y_1, \dots, Y_T)'$ be a $T \times 1$ vector of outcomes and $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_T)'$ a $T \times p$ matrix of regressors with $\mathbf{X}_t \in \mathbb{X}_t \subset \mathbb{R}^p$ and $\mathbf{X} \in \mathbb{X}^T$ where $\mathbb{X}^T = \times_{t \in \{1, \dots, T\}} \mathbb{X}_t$. We assume that \mathbf{X}_t is strictly exogenous. This rules out feedback from the period t outcome Y_t to the period $s \geq t$ regressor \mathbf{X}_s . One implication of this assumption is that lags of the dependent variable may not be included in \mathbf{X}_t . Our model is a static one.

Available is a random sample $\{(\mathbf{Y}_i, \mathbf{X}_i)\}_{i=1}^N$ from a distribution F_0 . The t^{th} period outcome is given by

$$Y_t = \mathbf{X}_t' b_t(A, U_t), \quad (1)$$

where A is time-invariant unobserved unit-level heterogeneity and U_t a time-varying disturbance. Both A and U_t may be vector-valued. The $p \times 1$ vector of functions $b_t(A, U_t)$, which we allow to vary over time, map A and U_t into unit-by-period-specific slope coefficients. By ‘random’ coefficients we mean that $b_t(A, U_t)$ varies across units. By ‘correlated’, we mean that the entire path of regressor values, \mathbf{X} , may have predictive power for $b_t(A, U_t)$. This implies that an agent’s incremental return to an additional unit of \mathbf{X}_t may vary with \mathbf{X}_t . In this sense \mathbf{X}_t may be endogenous.

Equation (1) is structural in the sense that the *unit-specific* function

$$Y_t(\mathbf{x}_t) = \mathbf{x}_t' b_t(A, U_t) \quad (2)$$

traces out a unit’s period t potential outcome across different hypothetical values of $\mathbf{x}_t \in \mathbb{X}_t$.³ Let $\mathbf{X}_t = (1, X_t)'$; setting $b_{1t}(A, U_t) = \beta_1 + A + U_t$ (with A and U_t scalar and mean zero) and $b_{kt}(A, U_t) = \beta_k$ for $k = 2, \dots, p$ yields the textbook linear panel data model:

$$Y_t(\mathbf{x}_t) = \mathbf{x}_t' \beta + A + U_t. \quad (3)$$

Equation (2), while preserving linearity in \mathbf{X}_t , is more flexible than (3) in that it allows for time-

²Chamberlain (1984) presents several well-formulated economic models that *do* imply linear specifications with unit-specific intercepts.

³Throughout we use capital letters to denote random variables, lower case letters specific realizations of them, and blackboard bold letters to denote their support (e.g., X , x and \mathbb{X}).

varying random coefficients on all of the regressors (not just the intercept). Furthermore these coefficients may nonlinearly depend on A and/or U_t .

Our goal is to characterize the effect of an exogenous change in \mathbf{X}_t on the probability distribution of Y_t . By ‘exogenous change’ we mean an external manipulation of \mathbf{X}_t in the sense described by Blundell and Powell (2003) or Imbens and Newey (2009). We begin by studying identification and estimation of the *average partial effect* (APE) of \mathbf{X}_t on Y_t (cf., Chamberlain, 1984; Blundell and Powell, 2003; Wooldridge, 2005a). Under (1) the average partial effect is given by

$$\beta_{0t} \stackrel{def}{=} \mathbb{E} \left[\frac{\partial Y_t(\mathbf{x}_t)}{\partial \mathbf{x}_t} \right] = \mathbb{E} [b_t(A, U_t)]. \quad (4)$$

Identification and estimation of (4) is nontrivial because, in our setup, \mathbf{X}_t may vary systematically with A and/or U_t . To see the consequences of such dependence observe that the derivative of the mean regression function of Y_t given \mathbf{X} does not identify a structural parameter. Differentiating through the integral we have

$$\frac{\partial \mathbb{E} [Y_t | \mathbf{X} = \mathbf{x}]}{\partial \mathbf{x}_t} = \beta_{0t}(\mathbf{x}) + \mathbb{E} [Y_t(\mathbf{X}_t) \mathbb{S}_{X_t}(A, U_t | \mathbf{X}) | \mathbf{X} = \mathbf{x}], \quad (5)$$

with $\beta_{0t}(\mathbf{x}) = \mathbb{E} [b_t(A, U_t) | \mathbf{X} = \mathbf{x}]$ and $\mathbb{S}_{X_t}(A, U_t | \mathbf{X}) = \nabla_{X_t} \log f(A, U_t | \mathbf{X})$. The second term is what Chamberlain (1982) calls heterogeneity bias. If the (log) density of the unobserved heterogeneity varies sharply with \mathbf{x}_t – corresponding to ‘selection bias’ or ‘endogeneity’ in a unit’s choice of \mathbf{x}_t – then the second term in (5) can be quite large.

Chamberlain (1982) studies identification of $\beta_0 \equiv \beta_{00}$ using panel data (cf., Mundlak, 1961, 1978b). In a second paper, Chamberlain (1992, pp. 579 - 585) calculates the semiparametric variance bound for β_0 and proposes an efficient method-of-moments estimator.⁴ His approach is based on a generalized within-group transformation; naturally extending the idea that panel data allow the research to control for time-invariant heterogeneity by ‘differencing it away’.⁵ Under regularity conditions, which ensure nonsingularity of β_0 ’s information bound, Chamberlain’s estimator converges at the standard \sqrt{N} rate.

Nonsingularity of $\mathcal{I}(\beta_0)$, the information for β_0 , generally requires the time dimension of the panel to exceed the number of random coefficients ($T > p$). Depending on the time series properties of the regressors, T may need to substantially exceed p . In extreme cases $\mathcal{I}(\beta_0)$ may be zero for all values of T . In such settings Chamberlain’s method breaks down. We show that, under mild conditions, β_0 nevertheless remains identified. Our method of identification is necessarily ‘irregular’: the information bound is singular and hence no regular \sqrt{N} consistent estimator exists (Chamberlain, 1986). We develop a feasible analog estimator for β_0 and characterize its large sample properties. Although its rate of convergence is slower than the standard parametric one,

⁴Despite its innovative nature, and contemporary relevance given the resurgence of interest in models with heterogeneous marginal effects, Chamberlain’s work on the CRC model is not widely known. The CRC specification is not discussed in Chamberlain’s own *Handbook of Econometrics* chapter (Chamberlain, 1984), while the panel data portion of Chamberlain (1992) is only briefly reviewed in the more recent survey by Arellano and Honoré (2001).

⁵Bonhomme (2010) further generalizes this idea, introducing a notion of ‘functional differencing’.

its limiting distributions is normal. Inference is straightforward.

Our work shares features with other studies of irregularly identified semiparametric models (e.g., Chamberlain, 1986; Manski, 1987; Heckman, 1990; Horowitz, 1992; Abrevaya, 2000; Honoré and Kyriazidou, 1997; Kyriazidou, 1997; Andrews and Schafgans, 1998; Khan and Tamer, 2009). A general feature of irregular identification is its dependence on the special properties of small subpopulations. These special properties are, in turn, generated by specific features of the semiparametric model. Consequently these types of identification arguments tend to highlight the importance, sometimes uncomfortably so, of maintained modelling assumptions (cf., Chamberlain, 1986, pp. 205 - 207; Khan and Tamer, 2009).

Our approach exploits the different properties, borrowing a terminology introduced by Chamberlain (1982), of ‘movers’ and ‘stayers’. Loosely speaking these two subpopulations respectively correspond to those units whose regressors values, \mathbf{X}_t , change and do not change across periods (a precise definition in terms of singularity of a unit-specific design matrix is given below). We identify aggregate time effects using the variation in Y_t in the ‘stayers’ subpopulation. A common trends assumption allows us to extrapolate these estimated effects to the entire population. Having identified the aggregate time effects using stayers, we then identify the APE by the limit of a trimmed mean of a particular unit-specific vector of regression coefficients.

Connection to other work on panel data In order to connect our work to the wider panel data literature it is useful to consider the more general outcome response function:

$$Y_t(\mathbf{x}_t) = m(\mathbf{x}_t, A, U_t).$$

Identification of the APE in the above model may be achieved by one of two main classes of restrictions. The *correlated random effects* approach invokes assumptions on the joint distribution of $(\mathbf{U}, A) | \mathbf{X}$; with $\mathbf{U} = (U_1, \dots, U_T)'$. Mundlak (1978a,b) and Chamberlain (1980, 1984) develop this approach for the case where $m(\mathbf{X}_t, A, U_t)$ and $F(\mathbf{U}, A | \mathbf{X})$ are parametrically specified. Newey (1994a) considers a semiparametric specification for $F(\mathbf{U}, A | \mathbf{X})$ (cf., Arellano and Carrasco, 2003). Recently, Altonji and Matzkin (2005) and Bester and Hansen (2009) have extended this idea to the case where $m(\mathbf{X}_t, A, U_t)$ is either semi- or non-parametric along with $F(\mathbf{U}, A | \mathbf{X})$.

The *fixed effects* approach imposes restrictions on $m(\mathbf{X}_t, A, U_t)$ and $F(\mathbf{U} | \mathbf{X}, A)$, while leaving $F(A | \mathbf{X})$, the distribution of the time-invariant heterogeneity, the so-called ‘fixed effects’, unrestricted. Chamberlain (1980, 1984, 1992), Manski (1987), Honoré (1992), Abrevaya (2000), and Bonhomme (2010) are examples of this approach. Depending on the form of $m(\mathbf{X}_t, A, U_t)$, the fixed effect approach may not allow for a complete characterization of the effect of exogenous changes in \mathbf{X}_t on the probability distribution of Y_t . Instead only certain features of this relationship may be identified (e.g., ratios of the average partial effect of two regressors).

Our methods are of the ‘fixed effect’ variety. In addition to assuming the CRC structure for $Y_t(\mathbf{x}_t)$ we impose a marginal stationarity restriction on $F(U_t | \mathbf{X}, A)$, a restriction also used by Manski (1987), Honoré (1992) and Abrevaya (2000), however, other than some weak smoothness

conditions, we leave $F(A|\mathbf{X})$ unrestricted.

Wooldridge (2005b) and Arellano and Bonhomme (2009) also analyze the CRC panel data model. Wooldridge focuses on providing conditions under which the usual linear fixed effects (FE) estimator is consistent despite the presence of correlated random coefficients (cf., Chamberlain, 1982, p. 11). Arellano and Bonhomme (2009) study the identification and estimation of higher-order moments of the distribution of the random coefficients. Unlike us, they maintain Chamberlain’s (1992) regularity conditions as well as impose additional assumptions.

Chamberlain (1982) showed that when \mathbf{X}_t is discretely valued the APE is generally not identified (p. 13). However, Chernozhukov, Fernández-Val, Hahn and Newey (2009), working with more general forms for $\mathbb{E}[Y_t|\mathbf{X}, A]$, show that when Y_t has bounded support the APE is partially identified and propose a method of estimating the identified set.⁶ In contrast, in our setup we show that the APE is point identified when at least one component of \mathbf{X}_t is continuously-valued.

Section 1 presents our identification results. We begin by (i) briefly reviewing the approach of Chamberlain (1992) and (ii) characterizing irregularity in the CRC model. We then present our method of irregular identification. Section 2 outlines our estimator as well as its large sample properties. Section 3 discusses various extensions of our basic approach.

In Section 4 we use our methods to estimate the average elasticity of calorie demand with respect to total household resources. Our sample is drawn from a population that participated in a pilot of the Nicaraguan conditional cash transfer program Red de Protección Social (RPS). Hunger is widespread in the communities from which our sample is drawn; we estimate that immediately prior to the start of the RPS program over half of households had less than the required number of calories needed for all their members to engage in ‘light activity’ on a daily basis.⁷

A stated goal of the RPS program is to reduce childhood malnutrition, and consequently increase human capital, by directly augmenting household income. The efficacy of this approach to reducing childhood malnutrition largely depends on the size of the average elasticity of calories demanded with respect to income across poor households.⁸ While most estimates of the elasticity of calorie demand are significantly positive, several recent estimates are small in value and/or imprecisely estimated, casting doubt on the value of income-oriented anti-hunger programs (Behrman and Deolalikar, 1987).⁹

⁶They consider the probit and logit models with unit-specific intercepts (in the index) in detail. They show how to construct bounds on the APE despite the incidental parameters problem and provide conditions on the distribution of \mathbf{X}_t such that these bounds shrink as T grows.

⁷We use Food and Agricultural Organization (FAO, 2001) gender- and age-specific energy requirements for ‘light activity’, as reported in Appendix 8 of Smith and Subandoro (2007), and our estimates of total calories available at the household-level to calculate the fraction of households suffering from ‘food insecurity’. Worldwide, the FAO estimates that 854 million people suffered from protein-energy malnutrition in 2001-03 (FAO, 2006). Halving this number by 2015, in proportion to the world’s total population, is the first United Nations Millennium Development Goal. Chronic malnutrition, particularly in early childhood, may adversely affect cognitive ability and economic productivity in the long-run (e.g., Dasgupta, 1993).

⁸Another motivation for studying this elasticity has to do with its role in theoretical models of nutrition-based poverty traps (see Dasgupta (1993) for a survey).

⁹Wolfe and Behrman (1983), using data from Somoza-era Nicaragua, estimate a calorie elasticity of just 0.01. Their estimate, if accurate, suggests that the income supplements provided by the RPS program should have little effect on caloric intake.

Disagreement about the size of the elasticity of calorie demand has prompted a vigorous methodological debate in development economics. Much of this debate has centered, appropriately so, on issues of measurement and measurement error (e.g., Bouis and Haddad, 1992; Bouis, 1994; Subramanian and Deaton, 1996). The implications of household-level correlated heterogeneity in the underlying elasticity for estimating its average, in contrast, have not been examined. If, for example, a households' food preferences, or preferences towards child welfare, co-vary with those governing labor supply, then its elasticity will be correlated with total household resources. An estimation approach which presumes the absence of such heterogeneity will generally be inconsistent for the parameter of interest. Our statistical model and corresponding estimator provides an opportunity, albeit in a specific setting, for assessing the relevance these types of heterogeneities.

We compare our CRC estimates of the elasticity of calorie demand with those estimated using standard panel data estimators (e.g., Behrman and Deolalikar, 1987; Bouis and Haddad, 1992), as well as those derived from cross-sectional regression methods as in Strauss and Thomas (1990, 1995), Subramanian and Deaton (1996), and others. Our preferred CRC elasticity estimates are 10 to 20 percent smaller than their corresponding textbook linear 'fixed effects' estimates (FE-OLS). Our results are consistent with the presence of modest 'correlated random coefficients bias'.

Section 5 summarizes and suggests areas for further research. Proofs are in the Appendix. The notation $\mathbf{0}_T$, $\mathbf{1}_T$, I_T and $\stackrel{D}{=}$ respectively denotes a $T \times 1$ vector of zeros, a $T \times 1$ vector of ones, the $T \times T$ identity matrix, and equality in distribution.

1 Identification

Our benchmark data generating process combines (1) with the following assumption.

Assumption 1.1 (STATIONARITY AND COMMON TRENDS)

- (i) $b_t(A, U_t) = b^*(A, U_t) + d_t(U_{2t})$ for $t = 1, \dots, T$ and $U_t = (U'_{1t}, U'_{2t})'$
- (ii) $U_t | \mathbf{X}, A \stackrel{D}{=} U_s | \mathbf{X}, A$ for $t = 1, \dots, T$, $t \neq s$
- (iii) $U_{2t} | \mathbf{X}, A \stackrel{D}{=} U_{2t}$
- (iv) $\mathbb{E}[b_t(A, U_t) | \mathbf{X} = \mathbf{x}]$ exists for all $t = 1, \dots, T$ and $\mathbf{x} \in \mathbb{X}^T$.

Part (i) of Assumption 1.1 implies that the random coefficient consists of a 'stationary' and 'nonstationary' component. The stationary part, $b^*(A, U_t)$, does not vary over time so that if $U_t = U_s$ we have $b^*(A, U_t) = b^*(A, U_s)$. The non-stationary part, which is a function of the subvector U_{2t} alone, may vary over time so that even if $U_{2t} = U_{2s}$ we may have $d_t(U_{2t}) \neq d_s(U_{2s})$.

Part (ii) imposes marginal stationarity of U_t given \mathbf{X} and A (cf., Manski, 1987). Marginal stationarity, while allowing for serial dependence in U_t , is restrictive. For example it rules out time-varying heteroscedasticity. Part (iii) requires that U_{2t} is independent of both \mathbf{X} and A . Maintaining (ii) and (iii) is weaker than assuming that U_t is i.i.d. over time and independent of \mathbf{X} and A as is often done in nonlinear panel data research (e.g., Chamberlain, 1980). Part (iv) is a technical condition. Note that Assumption 1.1 does not restrict the joint distribution of \mathbf{X} and A . Our model is a 'fixed effects' one.

Under Assumption 1.1 we have

$$\begin{aligned}
\mathbb{E}[b_t(A, U_t) | \mathbf{X}] &= \mathbb{E}[b^*(A, U_t) | \mathbf{X}] + \mathbb{E}[d_t(U_{2t}) | \mathbf{X}] \\
&= \mathbb{E}[b^*(A, U_1) | \mathbf{X}] + \mathbb{E}[d_t(U_{2t})] \\
&= \beta_0(\mathbf{X}) + \delta_{0t}, \quad t = 1, \dots, T,
\end{aligned} \tag{6}$$

where first equality uses part (i) of Assumption 1.1, the second parts (ii) and (iii), and the third establishes the notation $\beta_0(\mathbf{X}) = \mathbb{E}[b^*(A, U_1) | \mathbf{X}]$ and $\delta_{0t} = \mathbb{E}[d_t(U_{2t})]$. In what follows we normalize $\delta_{01} = \underline{0}$.

Equation (6) is a ‘common trends’ assumption. To see this consider two subpopulations with different regressor histories ($\mathbf{X} = \mathbf{x}$ and $\mathbf{X} = \mathbf{x}'$). Restriction (6) implies that¹⁰

$$\begin{aligned}
\mathbb{E}[b_t(A, U_t) | \mathbf{x}] - \mathbb{E}[b_s(A, U_s) | \mathbf{x}] &= \mathbb{E}[b_t(A, U_t) | \mathbf{x}'] - \mathbb{E}[b_s(A, U_s) | \mathbf{x}'] \\
&= \delta_{0t} - \delta_{0s}.
\end{aligned}$$

Now recall that a unit’s period t potential outcome function is $Y_t(\mathbf{x}_t) = \mathbf{x}'_t b_t(A, U_t)$. Let τ be any point in the support of both \mathbf{X}_t and \mathbf{X}_s , we have for all $\mathbf{x} \in \mathbb{X}^T$

$$\mathbb{E}[Y_t(\tau) - Y_s(\tau) | \mathbf{X} = \mathbf{x}] = \mathbb{E}[Y_t(\tau) - Y_s(\tau)] = \tau'(\delta_{0t} - \delta_{0s}). \tag{7}$$

Equation (7) implies that while the period t (linear) potential outcome functions may vary arbitrarily across subpopulations defined in terms of $\mathbf{X} = \mathbf{x}$, shifts in these functions over time are mean independent of \mathbf{X} . A variant of (7) is widely-employed in the program evaluation literature (e.g., Heckman, Ichimura, Smith and Todd, 1998; Angrist and Krueger, 1999). It is also satisfied by the linear panel data model featured in Chamberlain (1984).¹¹

Let the $(T - 1)p \times 1$ vector of aggregate shifts in the random coefficients $(\delta'_2, \dots, \delta'_T)'$ be denoted by $\boldsymbol{\delta}$ with the corresponding $T \times (T - 1)p$ matrix of time shifters given by

$$\mathbf{W} = \begin{pmatrix} \underline{0}'_p & \underline{0}'_p \\ \mathbf{X}'_2 & \underline{0}'_p \\ & \ddots \\ \underline{0}'_p & \mathbf{X}'_T \end{pmatrix}. \tag{8}$$

Under Assumption 1.1 we can write the conditional expectation of \mathbf{Y} given \mathbf{X} as:

$$\mathbb{E}[\mathbf{Y} | \mathbf{X}] = \mathbf{W}\boldsymbol{\delta}_0 + \mathbf{X}\beta_0(\mathbf{X}). \tag{9}$$

¹⁰We use the notation $\mathbb{E}[Y | \mathbf{X} = \mathbf{x}] = \mathbb{E}[Y | \mathbf{x}]$.

¹¹In an NBER working paper we show how to weaken (6) while still getting positive identification results. As we do not use these additional results when considering estimation they are omitted.

In some cases it will be convenient to impose a priori zero restrictions on $\boldsymbol{\delta}_0$ (which would imply restrictions on how $\mathbb{E}[Y_t(\mathbf{x}_t)]$ is allowed to vary over time). In order to accommodate such situations (without introducing additional notation) we can simply redefine \mathbf{W} and $\boldsymbol{\delta}_0$ accordingly. For example a model which allows only the intercept of $\mathbb{E}[Y_t(\mathbf{x}_t)]$ to shift over time is given by (9) above with $\mathbf{W} = (\mathbf{0}_{T-1}, I_{T-1})'$ and $\boldsymbol{\delta}_0$ equal to the $T-1$ vector of intercept shifts. To accommodate a range of options we hereon assume that \mathbf{W} is a known $T \times q$ function of \mathbf{X} .

Equation (9), which specifies a semiparametric mean regression function for \mathbf{Y} given \mathbf{X} , is the fundamental building block of the results that follow. Our identification results are based solely on different implications of (9). The role of equation (1) and Assumption 1.1 is to provide primitive restrictions on F_0 which imply (9). We emphasize that our results neither hinge on, nor necessarily fully exploit, all of these assumptions. Rather they flow from just one of their implications.

1.1 Regular identification

The partially linear form of (9) suggests identifying $\boldsymbol{\delta}_0$ using the conditional variation in \mathbf{W} given \mathbf{X} as in, for example, Engle, Granger, Rice and Weiss (1986).¹² In our benchmark model, however, \mathbf{W} is a $T \times q$ function of \mathbf{X} and hence no such conditional variation is available. Nevertheless Chamberlain (1992) has shown that $\boldsymbol{\delta}_0$ may be identified using the panel structure.

Let $\Phi(\mathbf{X})$ be some function of \mathbf{X} mapping into $T \times T$ positive definite matrices (in practice $\Phi(\mathbf{X}) = I_T$ will often suffice) and define the $T \times T$ idempotent ‘residual maker’ matrix:

$$M_\Phi(\mathbf{X}) = I_T - \mathbf{X} [\mathbf{X}'\Phi^{-1}(\mathbf{X})\mathbf{X}]^{-1} \mathbf{X}'\Phi^{-1}(\mathbf{X}). \quad (10)$$

Using the fact that $M_\Phi(\mathbf{X})\mathbf{X} = 0$ Chamberlain (1992) derived, for $T > p$ and other regularity conditions, the pair of moment restrictions

$$\mathbb{E} \left[\begin{array}{c} \mathbf{W}'\Phi^{-1}(\mathbf{X}) M_\Phi(\mathbf{X}) (\mathbf{Y} - \mathbf{W}\boldsymbol{\delta}_0) \\ [\mathbf{X}'\Phi^{-1}(\mathbf{X})\mathbf{X}]^{-1} \mathbf{X}'\Phi^{-1}(\mathbf{X}) (\mathbf{Y} - \mathbf{W}\boldsymbol{\delta}_0) - \boldsymbol{\beta}_0 \end{array} \right] = 0,$$

which identify $\boldsymbol{\delta}_0$ and $\boldsymbol{\beta}_0$ by

$$\boldsymbol{\delta}_0 = \mathbb{E} \left[\overline{\mathbf{W}}'_\Phi \Phi^{-1}(\mathbf{X}) \overline{\mathbf{W}}_\Phi \right]^{-1} \times \mathbb{E} \left[\overline{\mathbf{W}}'_\Phi \Phi^{-1}(\mathbf{X}) \overline{\mathbf{Y}}_\Phi \right] \quad (11)$$

$$\boldsymbol{\beta}_0 = \mathbb{E} \left[(\mathbf{X}'\Phi^{-1}(\mathbf{X})\mathbf{X})^{-1} \mathbf{X}'\Phi^{-1}(\mathbf{X}) (\mathbf{Y} - \mathbf{W}\boldsymbol{\delta}_0) \right], \quad (12)$$

where $\overline{\mathbf{W}}_\Phi = M_\Phi(\mathbf{X})\mathbf{W}$ and $\overline{\mathbf{Y}}_\Phi = M_\Phi(\mathbf{X})\mathbf{Y}$.

Note that $M_\Phi(\mathbf{X})$ may be viewed as a generalization of the within-group transform. To see this

¹²To be specific if $\widetilde{\mathbf{W}} = \mathbf{W} - \mathbb{E}[\mathbf{W}|\mathbf{X}]$ has a covariance matrix of full rank, then $\boldsymbol{\delta}_0 = \mathbb{E} \left[\widetilde{\mathbf{W}}'\widetilde{\mathbf{W}} \right]^{-1} \times \mathbb{E} \left[\widetilde{\mathbf{W}}'\mathbf{Y} \right]$.

note that premultiplying (9) by $M_\Phi(\mathbf{X})$ yields

$$\begin{aligned}\mathbb{E}[\overline{\mathbf{Y}}_\Phi | \mathbf{X}] &= \overline{\mathbf{W}}_\Phi \boldsymbol{\delta}_0 + M_\Phi(\mathbf{X}) \mathbf{X} \boldsymbol{\beta}_0(\mathbf{X}) \\ &= \overline{\mathbf{W}}_\Phi \boldsymbol{\delta}_0,\end{aligned}$$

so that $M_\Phi(\mathbf{X})$ ‘differences away’ the unobserved correlated effects, $\boldsymbol{\beta}_0(\mathbf{X})$. Equation (11) shows that $\boldsymbol{\delta}_0$ is identified by the remaining ‘within-group’ variation in \mathbf{W}_t .

With $\boldsymbol{\delta}_0$ asymptotically known, the APE is then identified by the (population) mean of the unit-specific generalized least squares (GLS) fits

$$\widehat{\boldsymbol{\beta}}_i = (\mathbf{X}'_i \Phi^{-1}(\mathbf{X}_i) \mathbf{X}_i)^{-1} \mathbf{X}'_i \Phi^{-1}(\mathbf{X}_i) (\mathbf{Y}_i - \mathbf{W}_i \boldsymbol{\delta}_0). \quad (13)$$

Chamberlain (1992) showed that setting $\Phi(\mathbf{X}) = \Sigma(\mathbf{X}) = \mathbb{V}(\mathbf{Y} | \mathbf{X})$ is optimal; resulting in estimators with asymptotic sampling variances equal to the variance bounds:

$$\mathcal{I}(\boldsymbol{\delta}_0)^{-1} = \mathbb{E} \left[\overline{\mathbf{W}}'_\Sigma \Sigma^{-1}(\mathbf{X}) \overline{\mathbf{W}}_\Sigma \right]^{-1} \quad (14)$$

$$\mathcal{I}(\boldsymbol{\beta}_0)^{-1} = \mathbb{V}(\boldsymbol{\beta}_0(\mathbf{X})) + \mathbb{E} \left[(\mathbf{X}' \Sigma^{-1}(\mathbf{X}) \mathbf{X})^{-1} \right] + K \mathcal{I}(\boldsymbol{\delta}_0)^{-1} K', \quad (15)$$

where $K = \mathbb{E} \left[(\mathbf{X}' \Sigma^{-1}(\mathbf{X}) \mathbf{X})^{-1} \mathbf{X}' \Sigma^{-1}(\mathbf{X}) \mathbf{W} \right]$.

1.2 Irregularity of the CRC panel data model

Chamberlain’s approach requires nonsingularity of $\mathcal{I}(\boldsymbol{\delta}_0)$ and $\mathcal{I}(\boldsymbol{\beta}_0)$. In this section we discuss when this condition might not hold and, consequently, no regular \sqrt{N} consistent estimator exists. We begin by noting that singularity $\mathcal{I}(\boldsymbol{\delta}_0)$ is generic if $T = p$, our primary case. The following proposition specializes Proposition 1 of Chamberlain (1992) to our problem.

Proposition 1.1 (ZERO INFORMATION) *Suppose that (i) $(F_0, \boldsymbol{\delta}_0, \boldsymbol{\beta}_0(\cdot))$ satisfies (9), (ii) $\Sigma(\mathbf{x})$ is positive definite for all $\mathbf{x} \in \mathbb{X}^T$, (iii) $\mathbb{E}[\mathbf{W}' \Sigma^{-1}(\mathbf{X}) \mathbf{W}] < \infty$ and (iv) $T = p$, then $\mathcal{I}(\boldsymbol{\delta}_0) = 0$.*

Proof. From Chamberlain (1992) the information bound for $\boldsymbol{\delta}_0$ is given by

$$\mathcal{I}(\boldsymbol{\delta}_0) = \mathbb{E} \left[\overline{\mathbf{W}}'_\Sigma \Sigma^{-1}(\mathbf{X}) \overline{\mathbf{W}}_\Sigma \right]$$

so that $\alpha' \mathcal{I}(\boldsymbol{\delta}_0) \alpha = 0$ is equivalent to $\overline{\mathbf{W}}_\Sigma \alpha = 0$ with probability one. If $T = p$, then \mathbf{X} is square so that

$$\overline{\mathbf{W}}_\Sigma = \mathbf{W} \left(I_T - \mathbf{X} [\mathbf{X}' \Sigma^{-1}(\mathbf{X}) \mathbf{X}]^{-1} \mathbf{X}' \Sigma^{-1}(\mathbf{X}) \mathbf{X} \right) = 0$$

such that $\overline{\mathbf{W}}_\Sigma \alpha = 0$. ■

An intuition for Proposition 1.1 is that when $T = p$ Chamberlain’s generalized within-group transform of \mathbf{W} eliminates *all* residual variation in \mathbf{W}_t over time. This is because the p predictors \mathbf{X}_t perfectly (linearly) predict each element of \mathbf{W}_t when $T = p$. Consequently the deviation of

\mathbf{W}_t from its ‘within-group mean’ is identically equal to zero; any approach based on within-unit variation will necessarily fail.

As a simple example consider the one period ($T = p = 1$) ‘panel data’ model where, suppressing the t subscript,

$$Y = \boldsymbol{\delta}_0 + Xb(A, U), \quad (16)$$

with X scalar. Under Assumption 1.1 this gives (9) with $\mathbf{W} = 1$ and $\mathbf{X} = X$. The generalized within-group operator for this model is $M_I(\mathbf{X}) = 1 - X \left(\frac{X^2}{\Phi(\mathbf{X})} \right)^{-1} \frac{X}{\Phi(\mathbf{X})} = 0$. Consequently $\bar{\mathbf{Y}}_I = \bar{\mathbf{W}}_I = 0$ and (11) does not identify $\boldsymbol{\delta}_0$. By Proposition 1.1 $\mathcal{I}(\boldsymbol{\delta}_0) = 0$. We show that $\boldsymbol{\delta}_0$ and $\boldsymbol{\beta}_0$ are irregularly identified in this model below.

We do not provide a general result on when regular \sqrt{N} estimation of $\boldsymbol{\beta}_0$ is possible. However some insight into this question can be gleaned from a few examples. First, when $T = p$, it appears as though $\boldsymbol{\beta}_0$ will not be regularly identifiable unless $\boldsymbol{\delta}_0$ is known. This can be conjectured by the form of (15) which will generally be infinite if $\mathcal{I}(\boldsymbol{\delta}_0)^{-1}$ is. Even if $\boldsymbol{\delta}_0$ is known regular identification can be delicate. Consider the $T = p = 1$ model given above. In this model the right-hand-side of (12) above specializes to $\mathbb{E}[(Y - \boldsymbol{\delta}_0)/X]$, which will be undefined in X has positive density in the neighborhood of zero.

Less obviously, regular estimation may be impossible in heavily overidentified models (i.e., those where T substantially exceeds p).¹³ To illustrate again consider (16) with $\boldsymbol{\delta}_0$ known, but with $T \geq 2$. Assume further that $\Sigma(\mathbf{X}) = I_T$ and $X_t = S \cdot Z_t$ where

$$S \sim \mathcal{U}[a, b], \quad Z_t \stackrel{iid}{\sim} \mathcal{N}(0, 1).$$

Variation in X_t over time in this model is governed by S , which varies across units. For those units with S close to zero, X_t will vary little across periods. The unit specific design matrix in this model is given by $\mathbf{X}'\Sigma^{-1}(\mathbf{X})\mathbf{X} = \sum Z_t^2 \cdot S \sim \chi_T^2 \cdot \mathcal{U}[a, b]$. If $0 < a < b$ then

$$\mathbb{E} \left[\left(\mathbf{X}'\Sigma^{-1}(\mathbf{X})\mathbf{X} \right)^{-1} \right] = \begin{cases} \frac{\ln(b) - \ln(a)}{T-2} & T \geq 3 \\ \infty & T < 3 \end{cases},$$

so the right-hand-side of (12) will be well-defined if $T \geq 3$. If $a \leq 0$, then it is undefined *regardless of the number of time periods*. If $a \leq 0$ the support of S will contain zero, ensuring a positive density of units whose values of \mathbf{X}_t do not change over time. These ‘stayers’ will have singular design matrices in (13), causing the variance bound for $\boldsymbol{\beta}_0$ to be infinite.

To summarize regular identification of $\boldsymbol{\beta}_0$ requires sufficient within-unit variation in \mathbf{X}_t for all units. This is a very strong condition. Many microeconomic applications are characterized by a preponderance of stayers.¹⁴ While time series variation in \mathbf{X}_t is essential for identification,

¹³In contrast the variance bound for $\boldsymbol{\delta}_0$ will be finite when $T > p$ as long as there is *some* variation in \mathbf{X}_t over time.

¹⁴In Card’s (1996) analysis of the union wage premium, for example, less than 10 percent of workers switch between collective bargaining coverage and non-coverage across periods (Table V, p. 971).

persistence in its process is common in practice. This persistence will often imply that the right-hand-side of (12) is undefined.

1.3 Irregular identification

In this section we show that, under weak conditions, δ_0 and β_0 are irregularly identified when $T = p$. We show how to extend our methods to the irregular $T > p$ (overidentified) case in Section 3 below. Let

$$D = \det(\mathbf{X})$$

and

$$\mathbf{X}^* = \text{adj}(\mathbf{X})$$

respectively denote the determinant and adjoint of \mathbf{X} such that $\mathbf{X}^{-1} = \frac{1}{D}\mathbf{X}^*$ when the former exists. In what follows we will often refer to units where $D = 0$ as *stayers*. To motivate this terminology consider the case where $T = p = 2$ with \mathbf{W} and \mathbf{X} in (9) equal to

$$\mathbf{W} = \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} 1 & X_1 \\ 1 & X_2 \end{pmatrix},$$

with X_t scalar. This corresponds to a model with (i) a random intercept and slope coefficient and (ii) a common intercept shift between periods one and two. In this model

$$D = X_2 - X_1 = \Delta X;$$

hence $D = 0$ corresponds to $\Delta X = 0$, or a unit's value of X_t 'staying' fixed across the two periods. More generally $D = 0$ if two or more rows of \mathbf{X} coincide, which occurs if \mathbf{X}_t does not change across adjacent periods or reverts to an earlier value in a later period. Loosely-speaking, we may think of stayers as units whose value of \mathbf{X}_t changes little across periods.

Let $\mathbf{Y}^* = \mathbf{X}^*\mathbf{Y}$ and $\mathbf{W}^* = \mathbf{X}^*\mathbf{W}$ equal \mathbf{Y} and \mathbf{W} after premultiplication by the adjoint of \mathbf{X} . In the $T = p = 2$ example introduced above we have

$$\mathbf{X}^* = \begin{pmatrix} X_2 & -X_1 \\ -1 & 1 \end{pmatrix}, \quad \mathbf{Y}^* = \begin{pmatrix} X_2 Y_1 - X_1 Y_2 \\ \Delta Y \end{pmatrix}, \quad \mathbf{W}^* = \begin{pmatrix} -X_1 \\ 1 \end{pmatrix}.$$

In an abuse of notation let $\beta_0(d) = \mathbb{E}[\beta_0(\mathbf{X}) | D = d]$. Our identification result, in addition to (9), requires the following assumption.

Assumption 1.2 (SMOOTHNESS AND CONTINUITY)

- (i) For some $u_0 > 0$, $D = \det(\mathbf{X})$ has $\Pr(|D| < h) = \int_{-h}^h \phi(u) du$ with $\phi(u) > 0$ for all $h \leq u_0$;
- (ii) $\mathbb{E}[\|\mathbf{W}^*\|^2] < \infty$ and $\mathbb{E}[\mathbf{W}^{*\prime}\mathbf{W}^* | D = 0]$ is nonsingular; and
- (iii) the functions $\beta_0(u)$, $\phi(u)$, $\mathbb{E}[\mathbf{W}^* | D = u]$, and $\mathbb{E}[\mathbf{W}^{*\prime}\mathbf{W}^* | D = u]$ are continuous in u for

$$-u_0 \leq u \leq u_0.$$

Part (i) of Assumption 1.2 is essential as our approach involves conditioning on different values of D . While the requirement that D has positive density near zero is indispensable, the implication that $\Pr(D = 0) = 0$ can be relaxed. In Section 3 we show how to deal with the case where the distribution D has a point mass at zero. This may occur if the distribution of \mathbf{X}_t has mass points at a finite set of values, while being continuously distributed elsewhere. If there is overlap in the mass points of \mathbf{X}_t and \mathbf{X}_s ($t \neq s$), then the distribution of D will have a mass point at zero.

Part (ii) of Assumption 1.2 is required for identification of $\boldsymbol{\delta}_0$. It will typically hold in well-specified models and is straightforward to verify. Part (iii) is a smoothness assumption which, in conjunction with (i), allows us to trim without changing the estimand.

Identification of the aggregate time effects, $\boldsymbol{\delta}_0$: We begin by premultiplying (9) by \mathbf{X}^* to get

$$\mathbb{E}[\mathbf{Y}^* | \mathbf{X}] = \mathbf{W}^* \boldsymbol{\delta}_0 + D \boldsymbol{\beta}_0(\mathbf{X}),$$

where we use the fact that $D I_T = \mathbf{X}^* \mathbf{X}$. Conditioning on the subpopulation of ‘stayers’ yields

$$\mathbb{E}[\mathbf{Y}^* | \mathbf{X}, D = 0] = \mathbf{W}^* \boldsymbol{\delta}_0. \tag{17}$$

Under Assumption 1.2 equation (17) implies that $\boldsymbol{\delta}_0$ is identified by the conditional linear predictor (CLP)

$$\boldsymbol{\delta}_0 = \mathbb{E}[\mathbf{W}^{*'} \mathbf{W}^* | D = 0]^{-1} \times \mathbb{E}[\mathbf{W}^{*'} \mathbf{Y}^* | D = 0]. \tag{18}$$

Equation (18) shows that the subpopulation of stayers is used to tie down the aggregate time effects, $\boldsymbol{\delta}_0$. Since stayer’s correspond to units whose values of \mathbf{X}_t change little over time, the evolution of Y_t among these units is driven by the aggregate time effects. This approach to identifying $\boldsymbol{\delta}_0$ is reminiscent of Chamberlain’s (1986) ‘identification at infinity’ result for the intercept of the censored regression model (p. 205). Both approaches use a *small subpopulation* to tie down a feature of the *entire population*. An important difference is that our result does not require \mathbf{X}_t to have unbounded support. Consequently, our identification result is not sensitive to the ‘tail properties’ of the distribution of \mathbf{X} . Our key requirement, that D have positive density in a neighborhood about zero, is straightforward to verify. We do this in the empirical application by plotting a kernel density estimate of $\phi(d)$, the density of D .

In the $T = p = 2$ example we have, conditional on $D = 0$, the equality $\mathbf{Y}^* = \mathbf{W}^* \Delta Y$ so that (18) simplifies to, recalling that $D = \Delta X$,

$$\boldsymbol{\delta}_0 = \mathbb{E}[\Delta Y | \Delta X = 0]. \tag{19}$$

The common intercept shift is identified by the average change in Y_t in the subpopulation of stayers. Identification of $\boldsymbol{\delta}_0$ is irregular since $\Pr(D = 0) = 0$; $\boldsymbol{\delta}_0$ corresponds to the value of the nonparametric mean regression of ΔY given D at $D = 0$. Note the importance of the (verifiable)

requirement that $\phi_0 > 0$ for this result.

As a second example of (18) consider the one period ‘panel data’ model introduced above. From (16) we have

$$\mathbb{E}[Y|X = 0] = \delta_0,$$

or ‘identification at zero’.

Identification of the average partial effects, β_0 : Treating δ_0 as known we identify $\beta_0(\mathbf{x})$ for all \mathbf{x} such that d is non-zero by

$$\beta_0(\mathbf{x}) = \mathbb{E}[\mathbf{X}^{-1}(\mathbf{Y} - \mathbf{W}\delta_0)|\mathbf{X} = \mathbf{x}]. \quad (20)$$

It is instructive to consider the $T = p = 2$ case introduced above. In that model the second component of the right-hand side of (20), corresponding to the slope coefficient on X_t , evaluates to

$$\begin{aligned} \beta_{20}(\mathbf{x}) &= \frac{\mathbb{E}[\Delta Y|\mathbf{X} = \mathbf{x}] - \delta_0}{x_2 - x_1} \\ &= \frac{\mathbb{E}[\Delta Y|\mathbf{X} = \mathbf{x}] - \mathbb{E}[\Delta Y|\Delta X = 0]}{x_2 - x_1} \end{aligned} \quad (21)$$

where the second, difference-in-differences, equality follows by substituting in (19) above. Equation (21) indicates that the average slope coefficient, in a subpopulation homogenous in $\mathbf{X} = \mathbf{x}$, is equal to the *average* ‘rise’ – $\mathbb{E}[\Delta Y|\mathbf{X} = \mathbf{x}]$ – over the *common* ‘run’ – $x_2 - x_1$. The evolution of Y_t amongst stayers is used to eliminate the aggregate time effect from the average rise (i.e., to control for ‘common trends’) in this computation.

Using (21) we then might try, by appealing to the law of iterated expectations, to identify β_{20} by

$$\mathbb{E}\left[\frac{\Delta Y - \delta_0}{\Delta X}\right]. \quad (22)$$

An approach based on (22) was informally suggested by Mundlak (1961, p. 45). Chamberlain (1982) considered (22) with $\delta_0 = 0$, showing that it identifies β_{20} if $\mathbb{E}[|\Delta Y/\Delta X|] < \infty$. However, if ΔX has a positive, continuous density at zero – and if $\mathbb{E}[|\Delta Y| | \Delta X = d] - \delta_0$ does not vanish at $d = 0$ – then (22) will not be finite. For example, if ΔY and ΔX are independently and identically distributed according to the standard normal distribution, then $\Delta Y/\Delta X$ will be distributed according to the Cauchy distribution, whose expectation does not exist.

More generally the expectation

$$\mathbb{E}[\mathbf{X}^{-1}(\mathbf{Y} - \mathbf{W}\delta_0)]$$

will be undefined if the distribution of \mathbf{X} is such that D has a positive density in the neighborhood of $D = 0$ (i.e., there is a positive density of ‘stayers’). This will occur when, for example, at least two rows of \mathbf{X} ‘nearly’ coincide for ‘enough’ units (i.e., when part (i) of Assumption 1.2 holds).

To deal with the small denominator effects of stayers we trim. Under parts (i) and (iii) of Assumption 1.2 we have the equalities (see equation (42) in the Appendix)

$$\begin{aligned}
\boldsymbol{\beta}_0 &= \mathbb{E}[\boldsymbol{\beta}_0(\mathbf{X})] \\
&= \lim_{h \downarrow 0} \mathbb{E}[\boldsymbol{\beta}_0(\mathbf{X}) \cdot \mathbf{1}(|D| > h)] \\
&= \lim_{h \downarrow 0} \mathbb{E}[\mathbf{X}^{-1}(\mathbf{Y} - \mathbf{W}\boldsymbol{\delta}_0) \cdot \mathbf{1}(|D| > h)], \tag{23}
\end{aligned}$$

so that $\boldsymbol{\beta}_0$ is identified by the limit of the trimmed mean of $\mathbf{X}^{-1}(\mathbf{Y} - \mathbf{W}\boldsymbol{\delta}_0)$. Trimming eliminates those units with near-singular design matrices (i.e., stayers); by taking limits and exploiting continuity we avoid changing the estimand.

Note that if there is a point mass of stayers such that $\Pr(D = 0) = \pi_0 > 0$, then (23) does not equal $\boldsymbol{\beta}_0$, instead it equals

$$\boldsymbol{\beta}_0^M = \mathbb{E}[\boldsymbol{\beta}_0(\mathbf{X}) | D \neq 0],$$

or the movers average partial effect (MAPE). Let $\boldsymbol{\beta}_0^S = \mathbb{E}[\boldsymbol{\beta}_0(\mathbf{X}) | D = 0]$ equal the corresponding stayers average partial effect (SAPE). In Section 3 we show how to extend our results to identify the full average partial effect $\boldsymbol{\beta}_0 = \pi_0 \boldsymbol{\beta}_0^S + (1 - \pi_0) \boldsymbol{\beta}_0^M$ in this case.

The following proposition, which is proven in the Appendix as a by-product of the consistency part of Theorem 2.1 below, summarizes our main identification result.

Proposition 1.2 (IRREGULAR IDENTIFICATION) *Suppose that (i) $(F_0, \boldsymbol{\delta}_0, \boldsymbol{\beta}_0(\cdot))$ satisfies (9), (ii) $\Sigma(\mathbf{x})$ is positive definite for all $\mathbf{x} \in \mathbb{X}^T$, (iii) $T = p$, and (iv) Assumption 1.2 holds, then $\boldsymbol{\delta}_0$ and $\boldsymbol{\beta}_0$ are identified by, respectively, (18) and (23).*

2 Estimation

Our approach to estimation is to replace (18) and (23) with their sample analogs. We begin by discussing our estimator for the common parameters $\boldsymbol{\delta}_0$. Let h_N denote some bandwidth sequence such that $h_N \rightarrow 0$ as $N \rightarrow \infty$. We estimate $\boldsymbol{\delta}_0$ by the nonparametric conditional linear predictor fit:

$$\hat{\boldsymbol{\delta}} = \left[\frac{1}{Nh_N} \sum_{i=1}^N \mathbf{1}(|D_i| \leq h_N) \mathbf{W}_i^{*'} \mathbf{W}_i^* \right]^{-1} \times \left[\frac{1}{Nh_N} \sum_{i=1}^N \mathbf{1}(|D_i| \leq h_N) \mathbf{W}_i^{*'} \mathbf{Y}_i^* \right]. \tag{24}$$

This estimator has asymptotic properties similar to a standard (uniform) kernel regression fit for a one-dimensional problem. In particular, in the proof to Lemma A.2 in the Appendix we show that

$$\mathbb{V}(\hat{\boldsymbol{\delta}}) = O\left(\frac{1}{Nh_N}\right) \gg O\left(\frac{1}{N}\right),$$

so that its mean squared error (MSE) rate of convergence is slower than $1/N$ when $h_N \rightarrow 0$. We also show that the leading bias term in $\hat{\boldsymbol{\delta}}$ is quadratic in the bandwidth so that the fastest rate of

convergence of $\widehat{\boldsymbol{\delta}}$ to $\boldsymbol{\delta}_0$ will be achieved when the bandwidth sequence is of the form

$$h_N^* \propto N^{-1/5}.$$

To center the limiting distribution of $\sqrt{Nh_N}(\widehat{\boldsymbol{\delta}} - \boldsymbol{\delta}_0)$ at zero we use a bandwidth sequence that approaches zero faster than the MSE-optimal one. In particular we assume that $(Nh_N)^{1/2} h_N^2 \rightarrow 0$ as $N \rightarrow \infty$. We discuss our chosen bandwidth sequence in more detail below.

With $\widehat{\boldsymbol{\delta}}$ in hand we then estimate $\boldsymbol{\beta}_0$ using the trimmed mean¹⁵

$$\widehat{\boldsymbol{\beta}} = \frac{\frac{1}{N} \sum_{i=1}^N \mathbf{1}(|D_i| > h_N) \mathbf{X}_i^{-1} (\mathbf{Y}_i - \mathbf{W}_i \widehat{\boldsymbol{\delta}})}{\frac{1}{N} \sum_{i=1}^N \mathbf{1}(|D_i| > h_N)}. \quad (25)$$

To derive the asymptotic properties of $\widehat{\boldsymbol{\beta}}$ we begin by considering those of the infeasible estimator based on the true value of the time effects, $\boldsymbol{\delta}_0$:

$$\widehat{\boldsymbol{\beta}}_I = \frac{\frac{1}{N} \sum_{i=1}^N \mathbf{1}(|D_i| > h_N) \mathbf{X}_i^{-1} (\mathbf{Y}_i - \mathbf{W}_i \boldsymbol{\delta}_0)}{\frac{1}{N} \sum_{i=1}^N \mathbf{1}(|D_i| > h_N)}. \quad (26)$$

Like $\widehat{\boldsymbol{\delta}}$ the variance of $\widehat{\boldsymbol{\beta}}_I$ is of order $1/Nh_N$, however its asymptotic bias is linear, not quadratic, in h_N . The fastest feasible rate of convergence of $\widehat{\boldsymbol{\beta}}_I$ to $\boldsymbol{\beta}_0$ is consequently slower than the that of $\widehat{\boldsymbol{\delta}}$ to $\boldsymbol{\delta}_0$ ($N^{-2/3}$ versus $N^{-4/5}$). In order to center the limiting distribution of $\sqrt{Nh_N}(\widehat{\boldsymbol{\beta}}_I - \boldsymbol{\beta}_0)$ at zero we assume that $(Nh_N)^{1/2} h_N \rightarrow 0$ as $N \rightarrow \infty$. This is stronger than what is needed to appropriately center the distribution of the aggregate time effects.

The value of studying the large sample properties of $\sqrt{Nh_N}(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)$ is that our feasible estimator is a linear combination of $\widehat{\boldsymbol{\beta}}_I$ and $\widehat{\boldsymbol{\delta}}$:

$$\widehat{\boldsymbol{\beta}} = \widehat{\boldsymbol{\beta}}_I + \widehat{\Xi}_N (\widehat{\boldsymbol{\delta}} - \boldsymbol{\delta}_0), \quad (27)$$

with

$$\begin{aligned} \widehat{\Xi}_N &= \frac{\frac{1}{N} \sum_{i=1}^N \mathbf{1}(|D_i| > h_N) \mathbf{X}_i^{-1} \mathbf{W}_i}{\frac{1}{N} \sum_{i=1}^N \mathbf{1}(|D_i| > h_N)} \\ &= \frac{\frac{1}{N} \sum_{i=1}^N \mathbf{1}(|D_i| > h_N) D_i^{-1} \mathbf{W}_i^*}{\frac{1}{N} \sum_{i=1}^N \mathbf{1}(|D_i| > h_N)}. \end{aligned} \quad (28)$$

Note that $\widehat{\boldsymbol{\beta}}_I$ and $\widehat{\boldsymbol{\delta}}$ are respectively computed using the $|D_i| > h_N$ and $|D_i| \leq h_N$ subsamples, so they are conditionally independent given the $\{\mathbf{X}_i\}$. This independence exploits the fact that

¹⁵The denominator in (25) could be replaced by 1.

the same bandwidth sequence is used to estimate $\widehat{\boldsymbol{\delta}}$ and $\widehat{\boldsymbol{\beta}}$; it also results from our choice of the uniform kernel, which has bounded support. We proceed under these maintained assumption, acknowledging that it means that the rate of convergence of $\widehat{\boldsymbol{\delta}}$ to $\boldsymbol{\delta}_0$ is well below its optimal one. We view the gains from using the same bandwidth sequence for both $\widehat{\boldsymbol{\delta}}$ and $\widehat{\boldsymbol{\beta}}$ – in terms of simplicity and transparency of asymptotic analysis – as worth the cost in generality.

Lemma A.3 in the Appendix shows that

$$\widehat{\Xi}_N \xrightarrow{p} \Xi_0 \equiv \lim_{h \downarrow 0} \mathbb{E} [\mathbf{1}(|D_i| > h) D_i^{-1} \mathbf{W}_i^*].$$

We therefore recover the limiting distribution of the feasible estimator $\widehat{\boldsymbol{\beta}}$ from our results on $\widehat{\boldsymbol{\beta}}_I$ and $\widehat{\boldsymbol{\delta}}$ using a delta method type argument based on (27).

To formalize the above discussion and provide a precise result we require the following additional assumptions.

Assumption 2.1 (RANDOM SAMPLING) $\{(\mathbf{Y}_i, \mathbf{X}_i)\}_{i=1}^N$ are i.i.d. draws from a distribution F_0 which satisfies condition (9) above.

Assumption 2.2 (BOUNDED MOMENTS) $\mathbb{E} [\|\mathbf{X}_i^* \mathbf{Y}_i\|^4 + \|\mathbf{X}_i^* \mathbf{W}_i\|^4] < \infty$.

Assumption 2.3 (SMOOTHNESS) $\beta_0(u)$, $\phi(u)$, $\mathbb{E}[\mathbf{W}^* | D = u]$, $\mathbb{E}[\mathbf{W}^{*\prime} \mathbf{W}^* | D = u]$, $\mathbb{E}[\mathbf{X}^* \boldsymbol{\Sigma}(\mathbf{X}) \mathbf{X}^{*\prime} | D = u]$, and $m_r(u) = \mathbb{E}[\|\mathbf{X}_i^* \mathbf{Y}_i\|^r + \|\mathbf{X}_i^* \mathbf{W}_i\|^r | D = u]$ exist and are twice continuously differentiable for u in a neighborhood of zero and $0 \leq r \leq 4$.

Assumption 2.4 (LOCAL IDENTIFICATION) $\mathbb{E}[\mathbf{X}^* \boldsymbol{\Sigma}(\mathbf{X}) \mathbf{X}^{*\prime} | D = 0]$ is positive definite.

Assumption 2.5 (BANDWIDTH) As $N \rightarrow \infty$ we have $h_N \rightarrow 0$ such that $Nh_N \rightarrow \infty$ and $(Nh_N)^{1/2} h_N \rightarrow 0$.

Assumption 2.1 is a standard random sampling assumption. Our methods could be extended to consider other sampling schemes in the usual way. Assumptions 2.2 and 2.3 are regularity conditions that allow for the application of Liapunov’s central limit theorem for triangular arrays (e.g., Serfling, 1980). Assumption 2.5 is a bandwidth condition which ensures that $\sqrt{Nh_N} (\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)$ is asymptotically centered at zero with a finite variance as discussed below.

The smoothness imposed by Assumption 2.3 can be restrictive. For example if $T = p = 2$ with $\mathbf{X}_t = (1, X_t)'$ and X_1 and X_2 independent exponential random variables with parameter $1/\lambda$, then $D = \Delta X$ will be a Laplace(0, λ) random variable (the density of which is non-differentiable at zero). Non-differentiability of the density of D at $D = 0$ will prevent us from consistently estimating the common time effects, $\boldsymbol{\delta}$ (and, consequently, also $\boldsymbol{\beta}_0$). To gauge the restrictiveness of Assumption 2.3 note that twice continuous differentiability is required for nonparametric kernel estimation of, for example, $\phi(u)$ and $\mathbb{E}[\mathbf{W}^* | D = u]$, and is, consequently, a standard assumption in the literature on nonparametric density and conditional moment estimation (e.g., Pagan and Ullah, 1999; Chapters 2, 3).

Theorem 2.1 (LARGE SAMPLE DISTRIBUTION) *Suppose that (i) $(F_0, \boldsymbol{\delta}_0, \boldsymbol{\beta}_0(\cdot))$ satisfies (9), (ii) $\Sigma(\mathbf{x})$ is positive definite for all $\mathbf{x} \in \mathbb{X}^T$, (iii) $T = p$, and (iv) Assumptions 1.2 to 2.5 hold, then $\widehat{\boldsymbol{\delta}} \xrightarrow{p} \boldsymbol{\delta}_0$ and $\widehat{\boldsymbol{\beta}} \xrightarrow{p} \boldsymbol{\beta}_0$ with the normal limiting distribution*

$$\sqrt{Nh_N} \begin{pmatrix} \widehat{\boldsymbol{\delta}} - \boldsymbol{\delta}_0 \\ \widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0 \end{pmatrix} \xrightarrow{D} \mathcal{N}(0, \Omega_0), \quad \Omega_0 = \begin{pmatrix} \frac{\Lambda_0}{2\phi_0} & \frac{\Lambda_0 \Xi'_0}{2\phi_0} \\ \frac{\Xi_0 \Lambda_0}{2\phi_0} & 2\Upsilon_0 \phi_0 + \frac{\Xi_0 \Lambda_0 \Xi'_0}{2\phi_0} \end{pmatrix}$$

where

$$\Lambda_0 = \mathbb{E}[\mathbf{W}^{*\prime} \mathbf{W}^* | D = 0]^{-1} \mathbb{E}[\mathbf{W}^{*\prime} \mathbf{X}^* \Sigma(\mathbf{X}) \mathbf{X}^{*\prime} \mathbf{W}^* | D = 0] \mathbb{E}[\mathbf{W}^{*\prime} \mathbf{W}^* | D = 0]^{-1}$$

$$\Upsilon_0 = \mathbb{E}[\mathbf{X}^* \Sigma(\mathbf{X}) \mathbf{X}^{*\prime} | D = 0].$$

We comment that, in contrast to the irregularly identified semiparametric models discussed in Heckman (1990), Andrews and Schafgans (1998), and Khan and Tamer (2009), the rate of convergence for our estimator does not depend on delicate ‘relative tail conditions’. Our identification approach is distinct from the type of ‘identification and infinity’ arguments introduced by Chamberlain (1986) and leads to a somewhat simpler asymptotic analysis.

It is instructive to compare the asymptotic variances given in Theorem 2.1 with Chamberlain’s regular counterparts (given in (14) and (15) above). First consider the asymptotic variance of $\widehat{\boldsymbol{\delta}}$. In our setup \mathbf{W}^* plays a role analogous to the generalized within-group transformation of \mathbf{W} used by Chamberlain (i.e., $\overline{\mathbf{W}}_\Phi = M_\Phi(\mathbf{X}) \mathbf{W}$). Viewed in this light the form of Λ_0 is similar to that of $\mathcal{I}(\boldsymbol{\delta}_0)^{-1}$ in the regular case. The key difference is that (i) the expectations in Λ_0 are conditional on $D = 0$ (i.e., averages over the subpopulation of stayers) and (ii) the variance of $\widehat{\boldsymbol{\delta}}$ varies inversely with ϕ_0 . The greater the density of stayers, the easier it is to estimate $\widehat{\boldsymbol{\delta}}$. We comment that we could estimate $\widehat{\boldsymbol{\delta}}$ more precisely if we replaced (24) with a weighted least squares estimator. We do not pursue this idea here as it would require pilot estimation of $\Sigma(\mathbf{X})$, a high dimensional object, and hence is unlikely to be useful in practice.

The asymptotic variance of $\widehat{\boldsymbol{\beta}}$ also parallels the form of $\mathcal{I}(\boldsymbol{\beta}_0)^{-1}$. The first term, $2\Upsilon_0 \phi_0$, plays the role of $\mathbb{E}\left[(\mathbf{X}' \Sigma^{-1}(\mathbf{X}) \mathbf{X})^{-1}\right]$ in (15). This term corresponds to the average of the conditional sampling variances of the unit specific slope estimates. The ‘better’ the typical unit-specific design matrix, the greater the precision of the average $\widehat{\boldsymbol{\beta}}$. In the irregular case $2\Upsilon_0 \phi_0$ captures a similar effect. There the average is conditional on $D = 0$. In contrast to the aggregate time effects, the first term in the variance of $\widehat{\boldsymbol{\beta}}$ varies linearly with ϕ_0 ; suggesting that a small density of stayers is better for estimation of $\boldsymbol{\beta}_0$.

The second term in $\widehat{\boldsymbol{\beta}}$ ’s variance is analogous to the $K\mathcal{I}(\boldsymbol{\delta}_0)^{-1}K'$ term in (15). This term captures the effect of sampling variation in $\widehat{\boldsymbol{\delta}}$ on that of $\widehat{\boldsymbol{\beta}}$. Note that K is equal to the average of the $p \times q$ matrix of coefficients associated with the unit-specific GLS fit of the $q \times 1$ vector \mathbf{W}_t given the $p \times 1$ vector of regressors \mathbf{X}_t . It is instructive to consider an example where there is no asymptotic penalty associated with not knowing $\boldsymbol{\delta}_0$. Let $\mathbf{W} = (\mathbf{0}_{T-1}, I_{T-1})'$ and $\mathbf{X}_t = (1, X_t)'$ with X_t scalar such that $p = 2$ and $\boldsymbol{\delta}_0$ corresponds to a $q = T - 1$ vector of time-specific intercept

shifts. If the distribution of X_t is stationary over time, then realizations of X_t cannot be used to predict the time period dummies. In that case each column of K will consist of a vector of zeros with the exception of the first element (which will equal $1/T$). The lower-right-hand element of $K\mathcal{I}(\boldsymbol{\delta}_0)^{-1}K'$ will equal zero so that ignorance of $\boldsymbol{\delta}_0$ does not affect the precision which the second component of $\boldsymbol{\beta}_0$, corresponding to the average slope, may be estimated.¹⁶

Now consider the irregular case where $T = p = 2$. We have

$$\Xi_0 = \lim_{h \downarrow 0} \mathbb{E} \left[\mathbf{1}(|\Delta X| > h) \frac{1}{\Delta X} \begin{pmatrix} -X_1 \\ 1 \end{pmatrix} \right],$$

so that the lower-right-hand element of $\Xi_0\Lambda_0\Xi_0'$ will equal zero if $\lim_{h \downarrow 0} \mathbb{E}[\mathbf{1}(|\Delta X| > h)/\Delta X] = 0$. This condition will hold if, for example, X_1 and X_2 are exchangeable, so that ΔX is symmetrically distributed about zero (at least for $|\Delta X|$ in a neighborhood of zero). This will ensure the asymptotic equivalence of the feasible estimator $\hat{\boldsymbol{\beta}}$ and its infeasible counterpart $\hat{\boldsymbol{\beta}}_I$.

Chamberlain's variance bound for $\boldsymbol{\beta}_0$ contains a third term the analog of which is not present in the irregular case. This term, $\mathbb{V}(\boldsymbol{\beta}_0(\mathbf{X}))$, captures the effect of heterogeneity in the conditional average of the random coefficients on the asymptotic variance of $\hat{\boldsymbol{\beta}}$. In the irregular case a term equal to $\mathbb{V}(\boldsymbol{\beta}_0(\mathbf{X}))/N$ also enters the expression for the sampling variance of $\hat{\boldsymbol{\beta}}$ (see the calculations immediately prior to Equation (45) in the Appendix). However this term is asymptotically dominated by the two terms listed in Theorem 2.1 (which are of order $1/Nh_N$). The variance estimator described in Theorem 2.2 below implicitly accounts for this asymptotically dominated component.

The conditions of Theorem 2.1 place only weak restrictions on the bandwidth sequence. As is common in the semiparametric literature we deal with bias by undersmoothing. The appendix shows that the fastest rate of convergence of $\hat{\boldsymbol{\beta}}$ for $\boldsymbol{\beta}_0$ in mean square is achieved by bandwidth sequences of the form,

$$h_N^* = C_0 N^{-1/3},$$

where the mean squared error minimizing choice of constant is

$$C_0 = \frac{1}{2} \left(\frac{1}{\phi_0} \right)^{1/3} \left\{ \left[(\boldsymbol{\beta}_0 - \boldsymbol{\beta}_0^S) (\boldsymbol{\beta}_0 - \boldsymbol{\beta}_0^S)' \right]^{-1} \times \left[2\Upsilon_0 + \frac{\Xi_0\Lambda_0\Xi_0'}{2\phi_0^2} \right] \right\}^{1/3}, \quad (29)$$

and $\boldsymbol{\beta}_0^S = \mathbb{E}[\boldsymbol{\beta}_0(\mathbf{X}) | D = 0]$ equals the average of the random coefficients in the subpopulation of stayers. While the bandwidth sequence h_N^* achieves the fastest rate of convergence for our estimator, the corresponding asymptotic normal distribution for $\hat{\boldsymbol{\beta}}(h_N^*)$ will be centered at a bias term of $2(\boldsymbol{\beta}_0 - \boldsymbol{\beta}_0^S)\phi_0$. To eliminate this bias Assumption 2.5 requires that $h_N \rightarrow 0$ fast enough such that $(Nh_N)^{1/2}h_N \rightarrow 0$ as $N \rightarrow \infty$, but slow enough such that $(Nh_N)^{1/2} \rightarrow \infty$. A bandwidth sequence which converges to zero slightly faster than h_N^* is sufficient for this purpose. In particular

¹⁶Sampling error in the estimated time effects does affect the precision with which the common intercept, the first component of $\boldsymbol{\beta}_0$, may be estimated

if

$$h_N = o(N^{-1/3}),$$

then $\sqrt{Nh_N}(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)$ will be asymptotically centered at zero.

An alternative to undersmoothing would be to use a plug in bandwidth based on a consistent estimate of (29), say \widehat{C} . Such an approach is taken by Horowitz (1992) in the context of smoothed maximum score estimation. Denote the resulting estimate by $\widehat{\boldsymbol{\beta}}_{\text{PI}}$ (PI for ‘plug in’). Let $\widehat{\boldsymbol{\beta}}$ be the consistent undersmoothed estimate of Theorem 2.1 and $\widehat{\boldsymbol{\beta}}^S$ and $\widehat{\phi}_0$ estimates of $\boldsymbol{\beta}_0^S$ and ϕ_0 . The bias corrected estimate is then

$$\widehat{\boldsymbol{\beta}}_{\text{BC}} = \widehat{\boldsymbol{\beta}}_{\text{PI}} - 2(\widehat{\boldsymbol{\beta}} - \widehat{\boldsymbol{\beta}}^S)\widehat{\phi}_0\widehat{C}N^{-1/3}.$$

Unlike undersmoothing, this does not slow down the rate of convergence of $\widehat{\boldsymbol{\beta}}_{\text{BC}}$ to $\boldsymbol{\beta}_0$. A disadvantage is that it is more computationally demanding. In the empirical application below we experiment with a number of bandwidth values. A more systematic analysis of bandwidth selection, while beyond the scope of this paper, would be an interesting topic for further research.

Computation and consistent variance estimation: The computation of $\widehat{\boldsymbol{\delta}}$ and $\widehat{\boldsymbol{\beta}}$ is facilitated by observing that the solutions to (24) and (25) above coincide with those of the linear instrumental variables fit

$$\widehat{\boldsymbol{\theta}} = \left[\frac{1}{N} \sum_{i=1}^N \mathbf{Q}'_i \mathbf{R}_i \right]^{-1} \times \left[\frac{1}{N} \sum_{i=1}^N \mathbf{Q}'_i \mathbf{Y}_i^* \right],$$

for $\boldsymbol{\theta} = (\boldsymbol{\delta}', \boldsymbol{\beta}')'$ and

$$\mathbf{Q}_i = \left(h_N^{-1} \mathbf{1}(|D_i| \leq h_N) \mathbf{W}_i^*, \frac{\mathbf{1}(|D_i| > h_N)}{D_i} I_p \right), \quad \mathbf{R}_i = (\mathbf{W}_i^*, \mathbf{1}(|D_i| > h_N) D_i I_p),$$

where the dependence of \mathbf{Q}_i and \mathbf{R}_i on h_N is suppressed.

Let $\boldsymbol{\theta}(h)$ denote the probability limit of $\widehat{\boldsymbol{\theta}}$ when the bandwidth is held fixed at h ; then, by standard GMM arguments,

$$\widehat{V}(h) = \left[\frac{1}{N} \sum_{i=1}^N \mathbf{Q}'_i \mathbf{R}_i \right]^{-1} \times \left[\frac{h}{N} \sum_{i=1}^N \mathbf{Q}'_i \widehat{\mathbf{U}}_i^+ \widehat{\mathbf{U}}_i^{+'} \mathbf{Q}_i \right] \times \left[\frac{1}{N} \sum_{i=1}^N \mathbf{Q}'_i \mathbf{R}_i \right]^{-1} \quad (30)$$

is a consistent estimate of the asymptotic covariance of $\sqrt{Nh}(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}(h))$ with

$$\widehat{\mathbf{U}}_i^+ = \mathbf{Y}_i^* - \mathbf{R}_i \widehat{\boldsymbol{\theta}}.$$

Conveniently, this covariance estimator remains valid when, as is required by Theorem 2.1, the bandwidth shrinks with N .

Theorem 2.2 *Suppose the hypotheses of Theorem 2.1 hold, and that $\mathbb{E}[\|\mathbf{X}_i^* \mathbf{Y}_i\|^8 + \|\mathbf{X}_i^* \mathbf{W}_i\|^8] <$*

∞ and that Assumption 2.3 holds for $r \leq 8$. Then $\widehat{V}_N \equiv \widehat{V}(h_N) \xrightarrow{p} \Omega_0$.

Relative to a direct estimate of Ω_0 , (30) implicitly includes estimates of terms that, while asymptotically negligible, may be sizeable in small samples. Consequently confidence intervals constructed using it may have superior properties (cf., Newey, 1994b; Graham, Imbens and Ridder, 2009).

Operationally estimation and inference may proceed as follows. Let \mathbf{Y}_{it}^* , \mathbf{R}_{it} and \mathbf{Q}_{it} denote the t^{th} rows of their corresponding matrices. Using standard software compute the linear instrumental variables fit of \mathbf{Y}_{it}^* onto \mathbf{R}_{it} using \mathbf{Q}_{it} as the instrument (exclude the default constant term from this calculation). By Theorem 2.2 the ‘robust/clustered’ (at the unit-level) standard errors reported by the program will be asymptotically valid under the conditions of Theorem 2.1.

3 Extensions

In this section we briefly develop three direct extensions of our basic results. In Section 5 we discuss other possible generalizations and avenues for future research.

3.1 Linear functions of $\beta_0(\mathbf{X})$

In some applications the elements of \mathbf{X}_t may be functionally related. For example

$$\mathbf{X}_t = \left(1, R_t, R_t^2, \dots, R_t^{p-1}\right)' . \quad (31)$$

In such settings β_0 indexes the average structural function (ASF) of Blundell and Powell (2003). To emphasize the functional dependence write $\mathbf{X}_t = \mathbf{x}_t(R_t)$, then

$$g_t(\tau) = \mathbf{x}_t(\tau)'(\beta_0 + \delta_{0t}) ,$$

gives the expected period t outcome if (i) a unit is drawn at random from the (cross sectional) population and (ii) she is exogenously assigned input level $R_t = \tau$. Similarly, differences of the form

$$g_t(\tau') - g_t(\tau) ,$$

give the average period t outcome difference across two counterfactual policies: one where all units are exogenously assigned input level $R_t = \tau'$ and another where they are assigned $R_t = \tau$. Since it is a linear function of β_0 and δ_{0t} Theorem 2.1 can be used to conduct inference on $g_t(\tau)$.

In the presence of functional dependence across the elements of \mathbf{X}_t the derivative of $g_t(\tau)$ with respect to τ does not correspond to an average partial effect (APE).¹⁷ Instead such derivatives characterize the local curvature of the ASF. In such settings the average effect of a population-wide

¹⁷We thank a referee for several helpful comments on this point.

unit increase in R_t (i.e., the APE) is instead given by

$$\begin{aligned}\gamma_{0t} &\stackrel{def}{=} \mathbb{E} \left[\frac{\partial Y_t(\mathbf{x}_t(r_t))}{\partial r_t} \right] = \mathbb{E} \left[\left(\frac{\partial \mathbf{x}_t(r_t)}{\partial r_t} \right)' b_t(A, U_t) \right] \\ &= \mathbb{E} \left[\left(\frac{\partial \mathbf{x}_t(r_t)}{\partial r_t} \right)' (\boldsymbol{\beta}_0(\mathbf{x}_t(R_t)) + \boldsymbol{\delta}_{0t}) \right],\end{aligned}\tag{32}$$

where the second equality follows from iterated expectations and Assumption 1.1. Because $\partial \mathbf{x}_t(R_t)/\partial r_t$ may covary with $\boldsymbol{\beta}_0(\mathbf{x}_t(R_t))$ Theorem 2.1 cannot be directly applied to conduct inference on γ_{0t} . Fortunately it is straightforward to extend our methods to identify and consistently estimate parameters of the form

$$\begin{aligned}\gamma_{0t} &= \mathbb{E} [\Pi(\mathbf{X}) (\boldsymbol{\beta}_0(\mathbf{X}) + \boldsymbol{\delta}_{0t})] \\ &= \gamma_0 + \mathbb{E} [\Pi(\mathbf{X})] \boldsymbol{\delta}_{0t},\end{aligned}$$

where $\Pi(\mathbf{x})$ is a known function of \mathbf{x} and $\gamma_0 = \mathbb{E} [\Pi(\mathbf{X}) \boldsymbol{\beta}_0(\mathbf{X})]$. If \mathbf{X}_t is given by (31), for example, then to estimate the APE we would choose

$$\Pi(\mathbf{x}) = \frac{\partial \mathbf{x}_t(r_t)}{\partial r_t} = (0, 1, 2r_t, 3r_t^2, \dots, (p-1)r_t^{p-2}).$$

In order to estimate γ_{0t} we proceed as follows. First, identification and estimation of $\boldsymbol{\delta}_0$ is unaffected. Second, using (20) gives for any \mathbf{x} with $d \neq 0$

$$\gamma_0(\mathbf{x}) = \Pi(\mathbf{x}) \mathbb{E} [\mathbf{X}^{-1} (\mathbf{Y} - \mathbf{W}\boldsymbol{\delta}_0) | \mathbf{X} = \mathbf{x}],$$

so that

$$\gamma_0 = \lim_{h \downarrow 0} \mathbb{E} [\Pi(\mathbf{X}) \mathbf{X}^{-1} (\mathbf{Y} - \mathbf{W}\boldsymbol{\delta}_0) \cdot \mathbf{1}(|D| > h)].$$

This suggests the analog estimator

$$\hat{\gamma} = \frac{\frac{1}{N} \sum_{i=1}^N \mathbf{1}(|D_i| > h_N) \Pi(\mathbf{X}_i) \mathbf{X}_i^{-1} (\mathbf{Y}_i - \mathbf{W}_i \hat{\boldsymbol{\delta}})}{\frac{1}{N} \sum_{i=1}^N \mathbf{1}(|D_i| > h_N)}.$$

An argument essentially identical to that justifying Theorem 2.1 then gives

$$\sqrt{Nh_N} \begin{pmatrix} \hat{\boldsymbol{\delta}} - \boldsymbol{\delta}_0 \\ \hat{\gamma} - \gamma_0 \end{pmatrix} \xrightarrow{D} \mathcal{N} \left(0, \begin{pmatrix} \frac{\Lambda_0}{2\phi_0} & \frac{\Lambda_0 \Xi'_{\Pi,0}}{2\phi_0} \\ \frac{\Xi_{\Pi,0} \Lambda_0}{2\phi_0} & 2\Upsilon_{\Pi,0} \phi_0 + \frac{\Xi_{\Pi,0} \Lambda_0 \Xi'_{\Pi,0}}{2\phi_0} \end{pmatrix} \right),\tag{33}$$

with

$$\begin{aligned}\Upsilon_{\Pi,0} &= \mathbb{E} [\Pi(\mathbf{X}) \mathbf{X}^* \Sigma(\mathbf{X}) \mathbf{X}^{*'} \Pi(\mathbf{X})' | D = 0] \\ \Xi_{\Pi,0} &= \lim_{h \downarrow 0} \mathbb{E} [\mathbf{1}(|D| > h) \Pi(\mathbf{X}) \mathbf{X}^{-1} \mathbf{W}].\end{aligned}$$

We can then estimate γ_{0t} by

$$\widehat{\gamma}_t = \widehat{\gamma} + \widehat{\Pi}\widehat{\delta},$$

with $\widehat{\Pi} = \sum_{i=1}^N \Pi(\mathbf{X}_i)/N$. To conduct inference on γ_{0t} we use the delta method treating $\widehat{\Pi}$ as known. We may ignore the effects of sampling variability in $\widehat{\Pi}$ since its rate of convergence to $\mathbb{E}[\Pi(\mathbf{X})]$ is $1/N$.

3.2 Density of D has a point mass at $D = 0$

In some settings a positive fraction of the population may be stayers such that $\pi_0 \stackrel{def}{=} \Pr(D = 0) > 0$. This may occur even if all elements of X_t are continuously-valued. If the only continuous component of \mathbf{X}_t is the logarithm of annual earnings, for example, a positive fraction of individuals may have the same earnings level in each sampled period. This may be especially true if many workers are salaried.

A point mass at $D = 0$ simplifies estimation of δ_0 and complicates that of β_0 . When $\pi_0 > 0$ the estimator

$$\widehat{\delta} = \left[\frac{1}{N} \sum_{i=1}^N \mathbf{1}(D_i = 0) \mathbf{W}_i^{*'} \mathbf{W}_i^* \right]^{-1} \times \left[\frac{1}{N} \sum_{i=1}^N \mathbf{1}(D_i = 0) \mathbf{W}_i^{*'} \mathbf{Y}_i^* \right],$$

will be \sqrt{N} -consistent and asymptotically normal for $\widehat{\delta}$, as would be the (asymptotically equivalent) estimator described in Section 2 above.

The large sample properties of the infeasible estimator $\widehat{\beta}^I$ – see equation (26) – are unaffected by the point mass at $D = 0$ with two important exceptions. First its probability limit is no longer β_0 , the (full population) average partial effect, but $\beta_0^M = \mathbb{E}[\beta_0(\mathbf{X}) | D \neq 0]$, the movers' APE introduced in Section 1. Second, its asymptotic variance is scaled up by $1 - \pi_0$, the population proportion of movers. This gives $\sqrt{Nh_N} (\widehat{\beta}^I - \beta_0^M) \xrightarrow{D} \mathcal{N} \left(0, \frac{2\Upsilon_0\phi_0}{1-\pi_0} \right)$.

Reflecting the change of plims let $\widehat{\beta}^M$ equal the feasible estimator defined by (25). Using decomposition (27) we have

$$\begin{aligned} \sqrt{Nh_N} (\widehat{\beta}^M - \beta_0^M) &= \sqrt{Nh_N} (\widehat{\beta}^I - \beta_0^M) + \widehat{\Xi}_N \sqrt{Nh_N} (\widehat{\delta} - \delta_0) \\ &= \sqrt{Nh_N} (\widehat{\beta}^I - \beta_0^M) + \Xi_0 O_p(\sqrt{h_N}) \\ &= \sqrt{Nh_N} (\widehat{\beta}^I - \beta_0^M) + o_p(1), \end{aligned}$$

so that the sampling properties of $\widehat{\beta}^M$ are unaffected by those of $\widehat{\delta}$. In particular, adapting the argument used to show Theorem 2.1 yields

$$\sqrt{Nh_N} (\widehat{\beta}^M - \beta_0^M) \xrightarrow{D} \mathcal{N} \left(0, \frac{2\Upsilon_0\phi_0}{1-\pi_0} \right).$$

If a consistent estimator of the stayers effect

$$\beta_0^S \equiv \mathbb{E}[\beta_0(\mathbf{X}) | D \neq 0]$$

can be constructed, a corresponding consistent estimator of the APE $\beta_0 = \pi_0 \beta_0^S + (1 - \pi_0) \beta_0^M$ would be

$$\hat{\beta} \equiv \hat{\pi} \hat{\beta}^S + (1 - \hat{\pi}) \hat{\beta}^M,$$

where $\hat{\pi} \equiv \sum_{i=1}^N \mathbf{1}(|D_i| \leq h_N) / N$ is a \sqrt{N} -consistent estimator for π_0 .

Inspection of the equation immediately preceding (17) in Section 1 suggests one possible estimator for β_0^S . We have

$$\mathbb{E}[\mathbf{Y}^* | \mathbf{X}] = \mathbf{W}^* \delta_0 + D \beta_0(\mathbf{X}),$$

so that

$$\beta_0^S = \lim_{h \downarrow 0} \frac{\mathbb{E}[\mathbf{Y}^* | D = h] - \mathbb{E}[\mathbf{Y}^* | D = 0]}{h},$$

which suggests the estimator

$$\begin{pmatrix} \bar{\delta} \\ \hat{\beta}^S \end{pmatrix} = \arg \min_{\delta, \beta^S} \sum_{i=1}^N \mathbf{1}(|D_i| \leq h_N) (\mathbf{Y}_i^* - \mathbf{W}_i^* \delta - D_i \beta^S)' (\mathbf{Y}_i^* - \mathbf{W}_i^* \delta - D_i \beta^S),$$

with $\bar{\delta}$ an alternative \sqrt{N} -consistent estimator for δ_0 . Since the rate of convergence of a nonparametric estimator of the derivative of a regression function is lower than for its level, the rate of convergence of the combined estimator $\hat{\beta} \equiv \hat{\pi} \hat{\beta}^S + (1 - \hat{\pi}) \hat{\beta}^M$ will coincide with that of $\hat{\beta}^S$, and the asymptotic distribution of the latter would dominate the asymptotic distribution of $\hat{\beta}$ in this setting. We comment that part (iii) of assumption 1.2 may be less plausible in settings where $\Pr(D = 0) > 0$.¹⁸ In such settings ‘stayers’ may be very different from ‘near stayers’ such that a local linear regression approach to estimating $\hat{\beta}^S$ would be problematic.

3.3 Overidentification ($T > p$)

When $T > p$ the vector of common parameters δ_0 may be \sqrt{N} consistently estimated, as first suggested by Chamberlain (1992), by the sample counterpart of (11) above:

$$\hat{\delta} = \left[\frac{1}{N} \sum_{i=1}^N \overline{\mathbf{W}}'_{\Phi_i} \Phi_i \overline{\mathbf{W}}_{\Phi_i}^{-1} \right]^{-1} \times \left[\frac{1}{N} \sum_{i=1}^N \overline{\mathbf{W}}'_{\Phi_i} \Phi_i(\mathbf{X}) \overline{\mathbf{Y}}_{\Phi_i} \right],$$

with $\Phi_i \equiv \Phi(\mathbf{X}_i)$ positive definite with probability one.

The discussion in Section 1, however, suggests that Chamberlain’s (1992) proposed estimate of β_0 , the sample average of

$$\hat{\beta}_i \equiv (\mathbf{X}_i' \Phi_i^{-1} \mathbf{X}_i)^{-1} \mathbf{X}_i' \Phi_i^{-1} (\mathbf{Y}_i - \mathbf{W}_i \hat{\delta}),$$

¹⁸We thank a referee for this observation.

for $\widehat{\boldsymbol{\delta}}$ a \sqrt{N} -consistent estimator of $\boldsymbol{\delta}_0$, may behave poorly and will be formally inconsistent when $\mathcal{I}(\boldsymbol{\beta}_0) = 0$.

Adapting the trimming scheme introduced for the just identified $T = p$ case, a natural modification of Chamberlain's (1992) estimator is

$$\widehat{\boldsymbol{\beta}} = \frac{\sum_{i=1}^N \mathbf{1}(\det(\mathbf{X}'_i \Phi_i^{-1} \mathbf{X}_i) > h_N) \cdot (\mathbf{X}'_i \Phi_i^{-1} \mathbf{X}_i)^{-1} \mathbf{X}'_i \Phi_i^{-1} (\mathbf{Y}_i - \mathbf{W}_i \widehat{\boldsymbol{\delta}})}{\sum_{i=1}^N \mathbf{1}(\det(\mathbf{X}'_i \Phi_i^{-1} \mathbf{X}_i) > h_N)}.$$

If

$$\mathbb{E} \left[\frac{1}{\det(\mathbf{X}'_i \Phi_i^{-1} \mathbf{X}_i)} \right] < \infty, \quad (34)$$

then the introduction of trimming is formally unnecessary but may still be helpful in practice. It is straightforward to show asymptotic equivalence of the (infeasible) trimmed mean

$$\widehat{\boldsymbol{\beta}} = \frac{\sum_{i=1}^N \mathbf{1}(\det(\mathbf{X}'_i \Phi_i^{-1} \mathbf{X}_i) > h_N) \cdot (\mathbf{X}'_i \Phi_i^{-1} \mathbf{X}_i)^{-1} \mathbf{X}'_i \Phi_i^{-1} (\mathbf{Y}_i - \mathbf{W}_i \boldsymbol{\delta}_0)}{\sum_{i=1}^N \mathbf{1}(\det(\mathbf{X}'_i \Phi_i^{-1} \mathbf{X}_i) > h_N)},$$

with Chamberlain's (1992) proposal when $\mathbb{E}[\boldsymbol{\beta}(\mathbf{X}) | \det(\mathbf{X}'_i \Phi_i^{-1} \mathbf{X}_i) \leq h]$ is smooth (Lipschitz-continuous) in h , condition (34) holds, and $h_N = o(1/\sqrt{N})$. Since $\widehat{\boldsymbol{\beta}}$ will still be consistent for $\boldsymbol{\beta}$ even when (34) fails, a feasible version of the trimmed mean $\widehat{\boldsymbol{\beta}}$ may be better behaved in finite samples if the design matrix $(\mathbf{X}'_i \Phi_i^{-1} \mathbf{X}_i)$ is nearly singular for some observations.

4 Application

In this section we use our methods to estimate an elasticity of calorie demand. Our goal is to provide a concrete illustration of our methods, to compare them with alternatives which presume the absence of any nonseparable correlated heterogeneity, and to highlight the practical importance of trimming.

Model specification: We assume that the logarithm of total household calorie availability per capita in period t , $\ln(\text{Cal}_t)$, varies according to

$$\ln(\text{Cal}_t) = b_{1t}(A, U_t) + b_{2t}(A, U_t) \ln(\text{Exp}_t), \quad (35)$$

where Exp_t denotes real household expenditure per capita in year t and $b_{1t}(A, U_t)$ and $b_{2t}(A, U_t)$ are random coefficients; the latter equals the household-by-period-specific elasticity of calorie demand. Let $b_t(A, U_t) = (b_{1t}(A, U_t), b_{2t}(A, U_t))'$, $\mathbf{X}_t = (1, \ln(\text{Exp}_t))'$, and $Y_t = \ln(\text{Cal}_t)$ with \mathbf{X} and \mathbf{Y} as defined above. We allow for common intercept and slope shifts over time (i.e., Assumption 1.1). The 2000 to 2002 period coincided with a coffee crisis in Nicaragua, so there is some *a priori* reason to believe that macro-shifts in the demand elasticity may be important.

Defining \mathbf{W} as in (8) gives a general semiparametric regression model of

$$\mathbb{E}[\mathbf{Y}|\mathbf{W}, \mathbf{X}] = \mathbf{W}\boldsymbol{\delta} + \mathbf{X}\boldsymbol{\beta}(\mathbf{X}),$$

with $\boldsymbol{\delta}$ containing the common intercept and slope time shifts. The second element of $\boldsymbol{\beta}(\mathbf{x})$ equals the average (base year) elasticity of calorie demand in the subpopulation of households with $\mathbf{X} = \mathbf{x}$ (i.e., a subpopulation with a common income/aggregate expenditure history).

Relative to prior work, the distinguishing feature of our model is that it allows for the elasticity of calorie demand to vary across households in a way that may co-vary with total outlay. This allows household expenditures to co-vary with the unobserved determinants of calorie demand. For example both expenditures and calorie consumption are likely to depend on labor supply decisions (cf., Strauss and Thomas, 1990). Allowing the calorie demand curve to vary across households also provides a ‘nonparametric’ way to control for differences in household composition; a delicate modelling decision in this context (e.g., Subramanian and Deaton, 1996).¹⁹

Data description: We use data collected in conjunction with an external evaluation of the Nicaraguan conditional cash transfer program Red de Protección Social (RPS) (see IFPRI, 2005). The RPS evaluation sample is a panel of 1,581 households from 42 rural communities in the departments of Madriz and Matagalpa, located in the northern part of the Central Region of Nicaragua. Each sampled household was first interviewed in August/September 2000 with follow-ups attempted in October of both 2001 and 2002. Here we analyze a balanced panel of 1,358 households from all three waves.²⁰

The survey was fielded using an abbreviated version of the 1998 Nicaraguan Living Standards Measurement Survey (LSMS) instrument. As such it includes a detailed consumption module with information on household expenditure, both actual and in kind, on 59 specific foods and several dozen other common budget items (e.g., housing and utilities, health, education, and household goods). The responses to these questions were combined to form an annualized consumption aggregate, C_{it} . In forming this variable we followed the algorithm outlined by Deaton and Zaidi (2002).

In addition to recording food expenditures, actual quantities of foods acquired are available. Using conversion factors listed in the World Bank (2002) and Instituto Nacional de Estadísticas y Censos (2005) (henceforth INEC) we converted all food quantities into grams. We then used the caloric content and edible percent information in the Instituto de Nutrición de Centro América y

¹⁹A limitation of our model is its presumption of linearity at the household level. Strauss and Thomas (1990) argue that the elasticity of demand should *structurally* decline with household income. As we have three periods of data we can, in principal, include an additional function of Exp_t in the \mathbf{X}_t vector. We briefly explore this possibility below.

²⁰A total of 1,359 households were successfully interviewed in all three waves. One of these households reports zero food expenditures (and hence calorie availability) in one wave and is dropped from our sample. The preparation of our estimation sample from the raw public release data files involved some complex and laborious data-processing. We outline the procedures used in this section. A sequence of heavily commented STATA do files, which read in the IFPRI (2005) release of the data and output a text file of our estimation sample is available online at <https://files.nyu.edu/bsg1/public/>

Panamá (2000) (henceforth INCAP) food composition tables to construct a measure of daily total calories available for each household.²¹ The logarithm of this measure divided by household size, Y_{it} , serves as the dependent variable in our analysis.

The combination of both expenditure and quantity information at the household-level also allowed us to estimate unit prices for foods. These unit values were used to form a Paasche cost-of-living index for the i^{th} household in year t of

$$I_{it} = \left[S_{it} \left\{ \sum_{f=1}^F W_{f,it} \left(P_f^b / P_{f,it} \right) \right\} + (1 - S_{it}) J_{it} \right]^{-1}, \quad (36)$$

where S_{it} is the fraction of household spending devoted to food, $W_{f,it}$ the share of overall food spending devoted to the f^{th} specific food, $P_{f,it}$ the year t unit price paid by the household for food f , and P_f^b its ‘base’ price (equal to the relevant 2001 sample median price). We use 2001 as our base year since it facilitates comparison with information from a nationwide LSMS survey fielded that year. Following the suggestion of Deaton and Zaidi (2002) we replace household-level unit prices with village medians in order to reduce noise in the price data. In the absence of price information on nonfood goods we set J_{it} equal to one in 2001 and to the national consumer price index (CPI) in 2000 and 2002. Our independent variable of interest is real per capita consumption in thousands of Cordobas: $\text{Exp}_{it} = ([C_{it}/I_{it}])/M_{it}$; M_{it} is total household size.

Tables 1 and 2 summarize some key features of our estimation sample. Panel A of Table 1 gives the share of total food spending devoted to each of eleven broad food categories. Spending on staples (cereals, roots and pulses) accounts for about half of the average household’s food budget and over two thirds of its calories (Tables 1 and 2). Among the poorest quartile of households an average of around 55 percent of budgets are devoted to, and over three quarters of calories available derived from, staples. Spending on vegetables, fruit, and meat accounts for less than 15 percent of the average household’s food budget and less than 3 percent of calories available. That such a large fraction of calories are derived from staples, while not good dietary practice, is not uncommon in poor households elsewhere in the developing world (cf., Smith and Subandoro, 2007).

Panel B of the table lists real annual expenditure in Cordobas per adult equivalent and per capita. Adult equivalents are defined in terms of age- and gender-specific FAO (2001) recommended energy intakes for individuals engaging in ‘light activity’ relative to prime-aged males. As a point of reference the 2001 average annual expenditure per capita across all of Nicaragua was a nominal C\$7,781, while amongst rural households it was C\$5,038 (World Bank. 2003). The 42 communities in our sample, consistent with their participation in an anti-poverty demonstration experiment, are considerably poorer than the average Nicaraguan rural community.²²

²¹In forming our measure of calorie availability we followed the general recommendations of Smith and Subandoro (2007).

²²In October of 2001 the Cordoba-to-US\$ exchange rate was 13.65. Therefore per capita consumption levels in our sample averaged less than US\$ 300 per year.

Panel A:	Expenditure Shares (%)								
	All			Lower 25%			Upper 25%		
	2000	2001	2002	2000	2001	2002	2000	2001	2002
Cereals	49.1	36.0	32.7	53.3	40.9	35.7	45.7	31.6	29.4
Roots	1.3	3.1	2.7	1.3	2.6	2.0	1.5	3.6	3.6
Pulses	11.6	12.5	13.6	11.2	13.8	16.5	10.6	10.7	11.3
Vegetables	3.2	4.9	4.5	2.8	4.3	3.4	3.8	5.8	5.3
Fruit	0.6	0.9	1.1	0.5	0.7	0.9	0.8	1.2	1.2
Meat	3.1	6.9	7.7	2.2	4.0	5.1	5.3	9.9	10.4
Dairy	11.2	14.7	17.3	9.0	12.0	15.0	13.1	16.8	19.2
Oil	4.0	5.0	5.0	3.5	5.2	5.0	3.9	4.7	4.7
Other foods	15.8	16.0	15.4	16.2	16.7	16.5	15.4	15.7	14.9
Staples [◇]	62.1	51.6	49.0	65.8	57.3	54.1	57.8	45.9	44.3
Panel B:	Total Real Expenditure & Calories								
Expenditure per adult [♭]	5,506	4,679	4,510	2,503	2,397	2,200	9,481	7,578	7,460
(Expenditure per capita)	(4,277)	(3,764)	(3,887)	(2,016)	(2,130)	(2,102)	(7,302)	(5,845)	(6,114)
Food share	71.2	69.2	68.8	73.8	69.1	68.6	67.0	67.9	67.6
Calories per adult [♭]	2,701	3,015	2,948	1,706	2,127	2,013	3,737	3,849	3,758
(Calories per capita)	(2,086)	(2,435)	(2,529)	(1,351)	(1,854)	(1,873)	(2,842)	(2,962)	(3,041)
Percent energy deficient [♮]	51.0	39.3	39.7	85.0	69.7	76.2	19.8	14.5	13.0

Table 1: Real expenditure food budget shares of RPS households from 2000 to 2002

NOTES: Authors' calculations based on a balanced panel of 1,358 households from the RPS evaluation dataset (see IFPRI (2005)). Real household expenditure equals total annualized nominal outlay divided by a Paasche cost-of-living index. Base prices for the price index are 2001 sample medians. The nominal exchange rate in October of 2001 was 13.65 Cordobas per US dollar. Total calorie availability is calculated using the RPS food quantity data and the calorie content and edible portion information contained in INCAP (2000). Lower and upper 25 percent refers to the bottom and top quartiles of households based on the average of year 2000, 2001 and 2002 real consumption per adult equivalent and thus contains the same set of households in all three years.

[◇] Sum of cereal, roots and pulses.

[♭] "Adults" correspond to adult equivalents based on FAO (2001) recommended energy requirements for light activity.

[♮] Percentage of households with estimated calorie availability less than FAO (2001) recommendations for light activity given household demographics.

	Calorie Shares (%)								
	All			Lower 25%			Upper 25%		
	2000	2001	2002	2000	2001	2002	2000	2001	2002
Cereals	57.7	60.3	59.9	60.7	63.9	62.0	55.5	57.1	57.4
Roots	1.5	1.5	1.6	1.9	1.5	1.2	1.6	1.8	2.1
Pulses	13.1	11.3	12.8	12.1	11.3	13.3	13.1	11.0	12.1
Vegetables	0.7	0.7	0.6	0.6	0.6	0.4	0.8	0.9	0.8
Fruit	0.3	0.5	0.4	0.3	0.3	0.4	0.5	0.7	5.8
Meat	0.7	1.3	1.3	0.5	0.7	0.7	1.3	1.9	1.9
Dairy	4.1	4.3	4.5	3.4	3.0	3.4	4.7	5.2	5.5
Oil	6.9	7.6	7.5	5.8	6.9	6.7	7.4	8.1	8.0
Other foods	15.0	12.6	11.4	14.7	11.9	11.9	15.2	13.2	11.5
Staples [◇]	72.3	73.1	74.3	74.7	76.7	76.6	70.2	69.9	71.7

Table 2: Calorie shares of RPS households from 2000 to 2002

NOTES: Authors' calculations based on a balanced panel of 1,358 households from the RPS evaluation dataset (see IFPRI (2005)). Total calorie availability is calculated using the RPS food quantity data and the calorie content and edible portion information contained in INCAP (2000). Lower and upper 25 percent refers to the bottom and top quartiles of households based on the average of year 2000, 2001 and 2002 real consumption per adult equivalent and thus contains the same set of households in all three years.

[◇] Sum of cereal, roots and pulses.

Using the FAO (2001) energy intake recommendations for 'light activity' we categorized each household, on the basis of its demographic structure, as energy deficient or not. By this criterion approximately 40 percent of households in our sample are energy deficient each period. Amongst the poorest quartile this fraction rises to over 75 percent. These figures are reported in Panel B of Table 1.

Results: Table 3 reports our point estimates. Our first estimate corresponds to the pooled ordinary least squares (OLS) fit of $\ln(\text{Cal}_t)$ onto $\ln(\text{Exp}_t)$ using all three waves of the RPS data. Aggregate shifts in the intercept and slope coefficient are included. Also included in the model, to control for variation in food prices across markets, is a vector of 42 village-specific intercepts. Variants of this specification are widely employed in empirical work (e.g, Subramanian and Deaton, 1996; Table 2). The pooled OLS calorie elasticities are reported in Column 1. The elasticity approximately equals 0.7 in 2000 and 0.6 in both 2001 and 2002. All three elasticities are precisely determined. The estimates are high relative to others in the literature, but realistic given the extreme poverty of the households in our sample.

Column 2 augments the first model by allowing the intercept to vary across households. This 'fixed effects' estimator (FE-OLS) is also widely used in empirical work when panel data are available (e.g., Behrman and Deolalikar, 1987; Bouis and Haddad, 1992). Allowing for household-specific intercepts increases the elasticity by about 10 percent in all three years. The standard errors almost double in size.

	Panel A: Calorie Demand Elasticities					Panel B: Sensitivity to Trimming		
	(1) OLS	(2) FE	(3) R-CRC	(4) MDLK	(5) I-CRC	(1) I-CRC	(2) I-CRC	(3) I-CRC
2000 Elasticity	0.6837 (0.0305)	0.7550 (0.0441)	0.6617 (0.0424)	-0.0133 (0.2571)	0.7125 (0.0763)	0.6867 (0.0906)	0.6352 (0.0639)	0.6993 (0.0522)
2001 Elasticity	0.6105 (0.0383)	0.6635 (0.0608)	0.5861 (0.0565)	—	—	—	—	—
2002 Elasticity	0.5959 (0.0245)	0.6466 (0.0416)	0.5521 (0.0476)	—	0.5062 (0.0791)	0.5299 (0.0928)	0.4690 (0.0583)	0.5368 (0.0421)
Percent trimmed	—	—	—	0	6.5	5	10	20
Time shifters?	Yes	Yes	Yes	No	Yes	Yes	Yes	Yes

Table 3: Estimates of the calorie Engel curve: linear case

NOTES: Estimates based on the balanced panel of 1,358 households described in the main text. "OLS" denotes least squares applied to the pooled 2000, 2001, and 2002 samples, "FE-OLS" least squares with household-specific intercepts, "R-CRC" Chamberlain's (1992) estimator with identity weight matrix, "MDLK" the Mundlak (1961)/Chamberlain (1982) estimator described in the main text, and "I-CRC" our irregular correlated random coefficients estimator (using the 2000 and 2002 waves only). All models, with the exception of "MDLK", include common intercept and slope shifts across periods. The standard errors are computed in a way that allows for arbitrary within-village correlation in disturbances across households and time.

In Column 3 we use Chamberlain’s (1992) regular correlated random coefficients (R-CRC) estimator with an identity weight matrix. Since we have three years of data and only two random coefficients his methods, at least in principle, apply. The top panel of Figure 1 plots a histogram of $\det(\mathbf{X}'\mathbf{X})$, which shows a reasonable amount of density in the neighborhood of zero. This suggests that the right-hand-side of (12) may be undefined in the population. In practice the R-CRC estimator generates ‘sensible’ point estimates with estimated standard errors approximately equal to those of the corresponding FE-OLS estimates. The R-CRC point estimates are smaller than both the OLS and FE-OLS ones.

Columns 4 and 5 are based on only the 2000 and 2002 waves of data. By dropping the middle wave of data we artificially impose that $T = p = 2$; this ensures irregularity (Proposition 1.1). The lower panel of Figure 1 plots a histogram of $D = \ln(\text{Exp}_{2002}) - \ln(\text{Exp}_{2000})$; the figure indicates a substantial amount of density in the neighborhood of zero. Column 4 reports the ‘Mundlak’ estimate of the demand elasticity

$$\frac{1}{N} \sum_{i=1}^N \frac{\ln(\text{Cal}_{2002}) - \ln(\text{Cal}_{2000})}{\ln(\text{Exp}_{2002}) - \ln(\text{Exp}_{2000})}.$$

This average, as expected, is poorly behaved. It implies a nonsensical elasticity estimate with a very large standard error. Column 5 implements our estimator (I-CRC) with a bandwidth of $h_N = c_D N^{-1/3}$ where $c_D = \min(s_D, r_D/1.34)$ is a robust estimate of the sample standard deviation of D (s_D is the sample standard deviation and r_D the interquartile range). This implies that we trim, or categorize as ‘stayers’, about 7 percent of our sample. The I-CRC point estimate is sensible and well-determined. While the estimated year 2000 elasticity is close to its OLS counterpart, the 2002 elasticity is 20 percent lower in magnitude. Panel B of the table explores the sensitivity of our point estimates to trimming. Overall we find that the Column 5 point estimates are insensitive to modest variations in the bandwidth.

A nonlinear model: As we have three periods of data we can modify our model to allow the calorie elasticity to vary non-linearly with income. Nonlinearity in the calorie demand curve has been emphasized by Strauss and Thomas (1990, 1995) and Subramanian and Deaton (1996). We consider the model

$$\ln(\text{Cal}_t) = b_{1t}(A, U_t) + b_{2t}(A, U_t) \ln(\text{Exp}_t) + b_{3t}(A, U_t) \text{Exp}_t^{-1},$$

so that a household’s period-specific demand elasticity is given by $b_{2t}(A, U_t) - b_{3t}(A, U_t) \text{Exp}_t^{-1}$. We estimate the average of this elasticity in 2000, 2001 and 2002 using the approach outlined in Section 3. Figure 2 plots the a histogram of D with $\mathbf{X}_t = (1, \ln(\text{Exp}_t), \text{Exp}_t^{-1})'$. Because $\ln(\text{Exp}_t)$ and Exp_t^{-1} are highly correlated within-units, the density of D is substantial in the neighborhood of zero.

Table 4 reports average elasticity estimates based on the extended model. The average elasticities

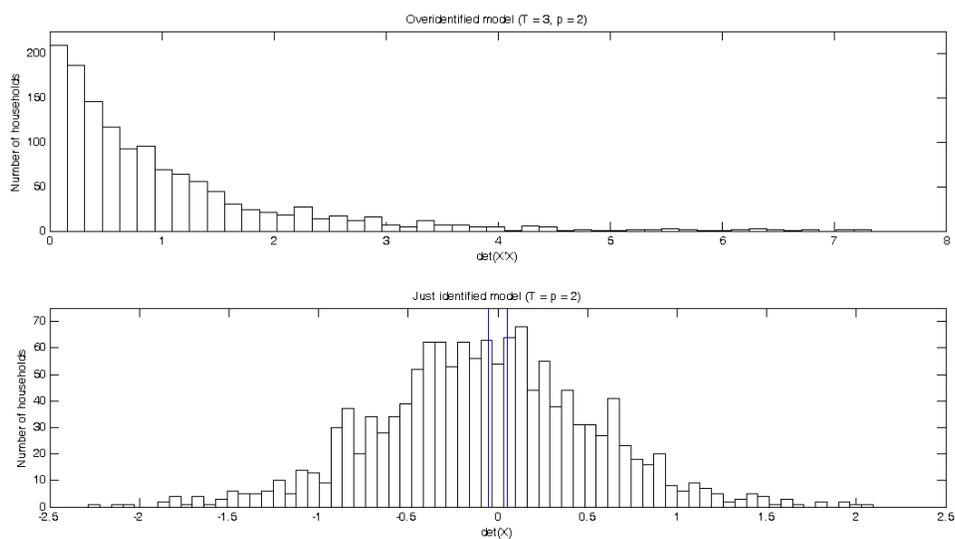


Figure 1: Histogram of the distribution of $\det(\mathbf{X}'\mathbf{X})$ (top panel, $T = 3, p = 2$) and D (bottom panel, $T = p = 2$)

NOTES: The two vertical blue lines in the lower panel correspond to the portion of the sample that is trimmed in our preferred estimates (Table 3, Column 5).

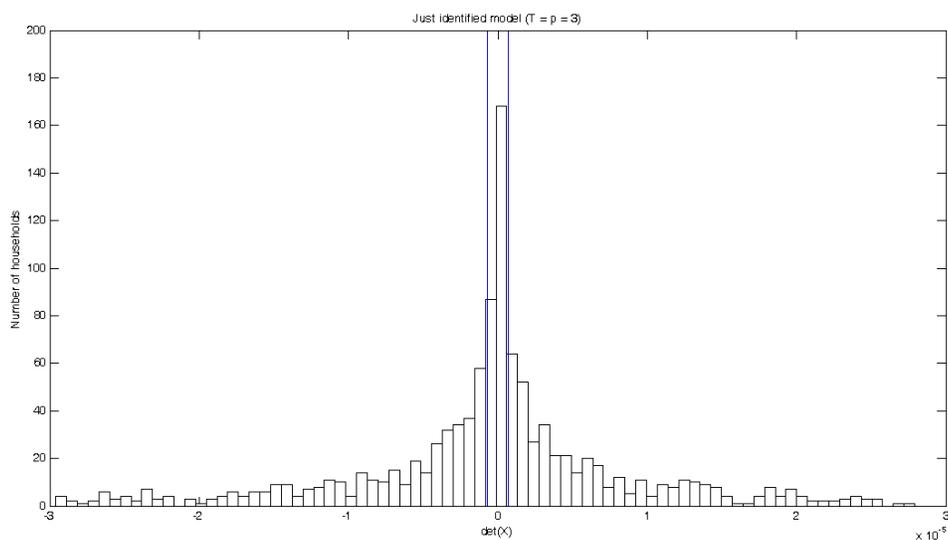


Figure 2: Histogram of the distribution of D ($T = p = 3$)

NOTES: The smallest and largest 10 percent of the D_i 's are excluded from the histogram. The two vertical blue lines correspond to the portion of the sample that is trimmed in our preferred estimates (Table 4, Column 3).

Panel A: Calorie Demand Elasticities					
	(1)	(2)	(3)	(4)	(5)
	OLS	FE	I-CRC	I-CRC	I-CRC
2000 Elasticity	0.6671 (0.0262)	0.7385 (0.0370)	0.2180 (0.2846)	0.1429 (0.2872)	0.2319 (0.2900)
2001 Elasticity	0.5992 (0.0380)	0.6514 (0.0584)	0.8889 (0.6342)	0.6150 (1.2170)	1.1262 (0.8696)
2002 Elasticity	0.5943 (0.0238)	0.6425 (0.0396)	0.5570 (0.5739)	0.6629 (0.3741)	0.4954 (0.5162)
Percent trimmed	—	—	18.8	15	25
Time shifters?	Yes	Yes	Yes	Yes	Yes

Table 4: Estimates of the calorie Engel curve: nonlinear case

NOTES: Estimates based on the balanced panel of 1,358 households described in the main text. "OLS" denotes least squares applied to the pooled 2000, 2001, and 2002 samples, "FE-OLS" least squares with household-specific intercepts, and "I-CRC" our irregular correlated random coefficients estimator (now using all three waves). All models include common intercept and slope shifts across periods. The standard errors are computed in a way that allows for arbitrary within-village correlation in disturbances across households and time. The average elasticity estimates in the OLS and FE-OLS columns are computed using the delta method. Those in the I-CRC columns as described in Section 3.

ties associated with the OLS and FE-OLS parameter estimates of the nonlinear model are virtually identical to their linear model Table 3 counterparts. Although the coefficients on Exp_t^{-1} and its interactions with the 2001 and 2002 time dummies are jointly significant in both models (not reported), the effect of their inclusion on the average elasticity estimates is negligible. Column 3 reports I-CRC estimates with $h_N = c_D N^{-1/3}$ (which, given the large density in the neighborhood of zero, results in the trimming of 20 percent of the sample). In contrast to the linear case, the I-CRC estimates are imprecisely determined; they are also more sensitive to variations in the bandwidth (Columns 3 and 4). We conclude that we are unable to reliably fit the nonlinear CRC model with the data available.

5 Conclusion

In this paper we have outlined a new estimator for the correlated random coefficients panel data model. Our estimator is designed for situations where the regularity conditions required for the method-of-moments procedure of Chamberlain (1992) do not hold. We illustrate the use of our methods in an exploration of the elasticity of demand for calories in a population of poor Nicaraguan households. This application is highly irregular, with many ‘near stayers’ in the sample. This implies that elasticity estimates based on the textbook FE-OLS estimator may be far from the relevant population average. We find that our methods work well in this setting, generating point estimates that are 10 to 20 percent smaller in magnitude than their FE-OLS counterparts (Table 3, Columns 5 versus 2).

While our procedure is simple to implement, it does require choosing a smoothing parameter. As in other areas of semiparametric econometrics, our theory places only weak restrictions on this choice. Developing an automatic, data-based, method of bandwidth selection would be useful.

Irregularity arises in other fixed effects panel data models (e.g., Manski, 1987; Honoré and Kyriazidou, 1997; Kyriazidou, 1997; Chamberlain 2010). It is an open question as to whether features of our approach could be extended to more complex nonlinear and/or dynamic panel data models. In ongoing work we are studying how to extend our methods to estimate quantile partial effects (e.g., unconditional quantiles of the distribution of the random coefficients) and to accommodate additional ‘triangular endogeneity’.

Appendix

A Proofs and derivations

A.1 Proof of Theorem 2.1

As noted in the main text our derivation of the limiting distribution of $\widehat{\beta}$ utilizes the decomposition

$$\widehat{\beta} = \widehat{\beta}_I + \widehat{\Xi}_N (\widehat{\delta} - \delta_0). \quad (37)$$

with $\widehat{\beta}_I$, $\widehat{\delta}$, and $\widehat{\Xi}_N$ respectively equal to (26), (24), and (28) of the main text. The proof proceeds in three steps. First we derive the limiting distribution of the infeasible estimator $\widehat{\beta}_I$. Second that of the common parameters $\widehat{\delta}$. Third we show that $\widehat{\Xi}_N$ has a well-defined probability limit. The limiting distribution of $\widehat{\beta}$ then follows from the delta method and the independence of $\widehat{\beta}_I$ and $\widehat{\delta}$.

Large sample properties of $\widehat{\beta}_I$: We begin with the infeasible estimator (26) which treats δ_0 as known. Recentering (26) yields

$$\widehat{\beta}_I - \beta_0 = \frac{\frac{1}{N} \sum_{i=1}^N \mathbf{1}(|D_i| > h_N) (\mathbf{X}_i^{-1} (\mathbf{Y}_i - \mathbf{W}_i \delta_0) - \beta_0)}{\frac{1}{N} \sum_{i=1}^N \mathbf{1}(|D_i| > h_N)}. \quad (38)$$

First consider the expected value of the term entering the summation in the denominator of (38):

$$\begin{aligned} \mathbb{E}[\mathbf{1}(|D_i| > h)] &= 1 - \Pr\{|D_i| \leq h\} \\ &= 1 - h \int_{-1}^1 \phi(uh) du \\ &= 1 - 2h\phi_0 + o(h) \\ &= 1 + o(1), \end{aligned} \quad (39)$$

where the third equality follows from Assumption 1.2 and the change of variables $u = t/h$ (with Jacobian $dt/du = h$).

Define $Z_{N,i}$ to be the term entering the summation in the numerator of (38):

$$Z_{N,i} \equiv \mathbf{1}(|D_i| > h) (\mathbf{X}_i^{-1} (\mathbf{Y}_i - \mathbf{W}_i \delta_0) - \beta_0). \quad (41)$$

Taking its expectation yields

$$\begin{aligned} \mathbb{E}[Z_{N,i}] &= \mathbb{E}[\mathbf{1}(|D_i| > h) \cdot (\beta_0(\mathbf{X}_i) - \beta_0)] \\ &= \mathbb{E}[\mathbf{1}(|D_i| \leq h) (\beta_0 - \beta_0(D_i))] \\ &= \int_{-h}^h (\beta_0 - \beta_0(t)) \phi(t) dt \\ &= 2 \left(\beta_0 - \beta_0^S \right) \phi_0 h + o(h), \end{aligned} \quad (42)$$

where $\beta_0^S \equiv \beta_0(0)$, again using Assumption 1.2.

Turning to the variance of $Z_{N,i}$ we use the ANOVA decomposition

$$\mathbb{V}(Z_{N,i}) = \mathbb{V}(\mathbb{E}[Z_{N,i}|D_i]) + \mathbb{E}[\mathbb{V}(Z_{N,i}|D_i)]. \quad (43)$$

The first term in (43) equals

$$\begin{aligned} \mathbb{V}(\mathbb{E}[Z_{N,i}|D_i]) &= \mathbb{E}[\mathbf{1}(|D_i| \leq h) (\beta_0(D_i) - \beta_0)^2] \\ &= \mathbb{V}(\beta_0(D_i)) + o(1). \end{aligned}$$

Now consider $\mathbb{V}(Z_{N,i}|D_i)$; using (42) above and recalling the equality $\mathbf{X}^{-1} = \frac{1}{D}\mathbf{X}^*$ when $|D| > 0$, we have

$$\begin{aligned} Z_{N,i} - \mathbb{E}[Z_{N,i}|D_i] &= \mathbf{1}(|D_i| > h) \{ \mathbf{X}_i^{-1} (\mathbf{Y}_i - \mathbf{W}_i \boldsymbol{\delta}) - \beta_0 - (\beta_0(D_i) - \beta_0) \} \\ &= \frac{\mathbf{1}(|D_i| > h)}{D_i} \{ \mathbf{Y}_i^* - \mathbf{W}_i^* \boldsymbol{\delta} - D_i \beta_0(D_i) \} \\ &= \frac{\mathbf{1}(|D_i| > h)}{D_i} \mathbf{X}_i^* (\mathbf{Y}_i - \mathbf{W}_i \boldsymbol{\delta} - \mathbf{X}_i \beta_0(D_i)). \end{aligned}$$

Again defining

$$\begin{aligned} \mathbf{U}_i &\equiv \mathbf{Y}_i - \mathbf{W}_i \boldsymbol{\delta} - \mathbf{X}_i \beta_0(\mathbf{X}_i) \\ &= \mathbf{Y}_i - \mathbf{W}_i \boldsymbol{\delta} - \mathbf{X}_i \beta_0(D_i) + \mathbf{X}_i (\beta_0(\mathbf{X}_i) - \beta_0(D_i)), \end{aligned}$$

it follows from iterated expectations that

$$\begin{aligned} \mathbb{V}(Z_{N,i}|D_i) &= \frac{\mathbf{1}(|D_i| > h_N)}{D_i^2} \mathbb{E}[\mathbf{X}_i^* (\mathbf{Y}_i - \mathbf{W}_i \boldsymbol{\delta} - \mathbf{X}_i \beta_0(D_i)) (\mathbf{Y}_i - \mathbf{W}_i \boldsymbol{\delta} - \mathbf{X}_i \beta_0(D_i))' \mathbf{X}_i^{*'} | D_i] \\ &= \frac{\mathbf{1}(|D_i| > h_N)}{D_i^2} \mathbb{E}[\mathbf{X}_i^* \boldsymbol{\Sigma}(\mathbf{X}_i) \mathbf{X}_i^{*'} | D_i] + \mathbf{1}(|D_i| > h_N) \mathbb{E}[(\beta_0(\mathbf{X}_i) - \beta_0(D_i)) (\beta_0(\mathbf{X}_i) - \beta_0(D_i))' | D_i] \\ &= \frac{\mathbf{1}(|D_i| > h_N)}{D_i^2} \mathbb{E}[\mathbf{X}_i^* \boldsymbol{\Sigma}(\mathbf{X}_i) \mathbf{X}_i^{*'} | D_i] + \mathbf{1}(|D_i| > h_N) \mathbb{V}(\beta_0(\mathbf{X}_i) | D_i), \end{aligned}$$

where $\boldsymbol{\Sigma}(\mathbf{X}_i) \equiv \mathbb{V}(\mathbf{U}_i | \mathbf{X}_i)$. Averaging the first term in $\mathbb{V}(Z_{N,i}|D_i)$ over the distribution of D_i gives

$$\begin{aligned} \mathbb{E}\left[\frac{\mathbf{1}(|D_i| > h_N)}{D_i^2} \mathbb{E}[\mathbf{X}_i^* \boldsymbol{\Sigma}(\mathbf{X}_i) \mathbf{X}_i^{*'} | D_i]\right] &= \int_{-\infty}^{-h} \frac{1}{t^2} \mathbb{E}[\mathbf{X}_i^* \boldsymbol{\Sigma}(\mathbf{X}_i) \mathbf{X}_i^{*'} | D_i = t] \phi(t) dt \\ &\quad + \int_h^{\infty} \frac{1}{t^2} \mathbb{E}[\mathbf{X}_i^* \boldsymbol{\Sigma}(\mathbf{X}_i) \mathbf{X}_i^{*'} | D_i = t] \phi(t) dt \\ &= \frac{1}{h} \int_{-\infty}^{-1} \frac{1}{u^2} \mathbb{E}[\mathbf{X}_i^* \boldsymbol{\Sigma}(\mathbf{X}_i) \mathbf{X}_i^{*'} | D_i = uh] \phi(uh) du \\ &\quad + \frac{1}{h} \int_1^{\infty} \frac{1}{u^2} \mathbb{E}[\mathbf{X}_i^* \boldsymbol{\Sigma}(\mathbf{X}_i) \mathbf{X}_i^{*'} | D_i = uh] \phi(uh) du \\ &= \frac{2\mathbb{E}[\mathbf{X}_i^* \boldsymbol{\Sigma}(\mathbf{X}_i) \mathbf{X}_i^{*'} | D_i = 0] \phi_0}{h} + O(1) \\ &= O(h^{-1}), \end{aligned} \quad (44)$$

where the third equality exploits Assumptions 1.2 and 2.3.

Averaging the second term over the distribution of D_i yields

$$\begin{aligned}\mathbb{E}[\mathbf{1}(|D_i| > h_N) \mathbb{V}(\boldsymbol{\beta}(\mathbf{X}_i) | D_i)] &\leq \mathbb{E}[\mathbb{V}(\boldsymbol{\beta}(\mathbf{X}_i) | D_i)] \\ &= \mathbb{V}(\boldsymbol{\beta}(\mathbf{X}_i)) \\ &= O(1).\end{aligned}$$

Thus, combing terms,

$$\mathbb{E}[\mathbb{V}(Z_{N,i} | D_i)] = \frac{2\mathbb{E}[\mathbf{X}_i^* \Sigma(\mathbf{X}_i) \mathbf{X}_i^{*'} | D_i = 0] \phi_0}{h} + O(1).$$

Combing this result with the expression for $\mathbb{V}(\mathbb{E}[Z_{N,i} | D_i])$ derived above yields a variance term of

$$\begin{aligned}\mathbb{V}(Z_{N,i}) &= \frac{2\mathbb{E}[\mathbf{X}_i^* \Sigma(\mathbf{X}_i) \mathbf{X}_i^{*'} | D_i = 0] \phi_0}{h} + O(1) \\ &= O(h^{-1}).\end{aligned}\tag{45}$$

Together (42), (45), and the independence generated by random sampling (Assumption 2.1) imply that

$$\mathbb{E}\left[\frac{1}{N} \sum_{i=1}^N Z_{N,i}\right] = o_p(1), \quad \mathbb{V}\left(\frac{1}{N} \sum_{i=1}^N Z_{N,i}\right) = O_p\left(\frac{1}{Nh_N}\right) = o_p(1),\tag{46}$$

under the bandwidth assumption (Assumption 2.5). This implies weak consistency of $\widehat{\boldsymbol{\beta}}_T$ for $\boldsymbol{\beta}_0$ (and indirectly Proposition 1.2).

To show asymptotic normality we check the conditions for Liapunov's Central Limit Theorem (CLT) for triangular arrays. Observe that

$$\begin{aligned}\mathbb{E}[\|Z_{N,i}\|^3] &= \mathbb{E}\left[\mathbf{1}(|D_i| > h) \left\|(\mathbf{X}_i^{-1}(\mathbf{Y}_i - \mathbf{W}_i \boldsymbol{\delta}_0) - \boldsymbol{\beta}_0)\right\|^3\right] \\ &= \mathbb{E}\left[\mathbf{1}(|D_i| > h) \left\|\frac{1}{D_i} \mathbf{X}_i^* (\mathbf{Y}_i - \mathbf{W}_i \boldsymbol{\delta}_0) - \boldsymbol{\beta}_0(D_i) - (\boldsymbol{\beta}_0 - \boldsymbol{\beta}_0(D_i))\right\|^3\right] \\ &= \mathbb{E}\left[\mathbf{1}(|D_i| > h) \left\|\frac{1}{D_i} \mathbf{X}_i^* \mathbf{U}_i + (\boldsymbol{\beta}_0(\mathbf{X}_i) - \boldsymbol{\beta}_0)\right\|^3\right].\end{aligned}$$

Application of the triangle inequality and the fact that $\mathbb{E}[\mathbf{X}^*(\mathbf{Y} - \mathbf{W}\boldsymbol{\delta}_0 - \mathbf{X}\boldsymbol{\beta}_0) | D] = 0$ yields

$$\begin{aligned}&\mathbb{E}\left[\left\|\frac{1}{D_i} \mathbf{X}_i^* (\mathbf{Y}_i - \mathbf{W}_i \boldsymbol{\delta}_0) - \boldsymbol{\beta}_0(D_i) - (\boldsymbol{\beta}_0 - \boldsymbol{\beta}_0(D_i))\right\|^3 \middle| D_i\right] \\ &\leq \mathbb{E}\left[\left(\left\|\frac{1}{D_i} \mathbf{X}_i^* (\mathbf{Y}_i - \mathbf{W}_i \boldsymbol{\delta}_0 - \mathbf{X}_i \boldsymbol{\beta}_0(D_i))\right\| + \|\boldsymbol{\beta}_0 - \boldsymbol{\beta}_0(D_i)\|\right)^3 \middle| D_i\right] \\ &= \left|\frac{1}{D_i^3}\right| \mathbb{E}\left[\left\|\mathbf{X}_i^* (\mathbf{Y}_i - \mathbf{W}_i \boldsymbol{\delta}_0 - \mathbf{X}_i \boldsymbol{\beta}_0(D_i))\right\|^3 \middle| D_i\right] + \|\boldsymbol{\beta}_0 - \boldsymbol{\beta}_0(D_i)\|^3 \\ &\quad + \frac{3\mathbb{E}\left[\left\|\mathbf{X}_i^* (\mathbf{Y}_i - \mathbf{W}_i \boldsymbol{\delta}_0 - \mathbf{X}_i \boldsymbol{\beta}_0(D_i))\right\|^2 \middle| D_i\right]}{D_i^2} \|\boldsymbol{\beta}_0 - \boldsymbol{\beta}_0(D_i)\|.\end{aligned}$$

Now note that

$$\begin{aligned}
\mathbb{E} \left[\|\mathbf{X}_i^* (\mathbf{Y}_i - \mathbf{W}_i \boldsymbol{\delta}_0 - \mathbf{X}_i \boldsymbol{\beta}_0(D_i))\|^2 \middle| D_i \right] &= \mathbb{E} \left[\text{Tr} (\mathbf{X}_i^* (\mathbf{Y}_i - \mathbf{W}_i \boldsymbol{\delta}_0 - \mathbf{X}_i \boldsymbol{\beta}_0(D_i)) (\mathbf{Y}_i - \mathbf{W}_i \boldsymbol{\delta}_0 - \mathbf{X}_i \boldsymbol{\beta}_0(D_i))' \mathbf{X}_i^{*'}) \middle| D_i \right] \\
&= \text{Tr} \left[\mathbb{E} \left[\mathbf{X}_i^* (\mathbf{Y}_i - \mathbf{W}_i \boldsymbol{\delta}_0 - \mathbf{X}_i \boldsymbol{\beta}_0(D_i)) (\mathbf{Y}_i - \mathbf{W}_i \boldsymbol{\delta}_0 - \mathbf{X}_i \boldsymbol{\beta}_0(D_i))' \mathbf{X}_i^{*'} \middle| D_i \right] \right] \\
&= \text{Tr} \left(\mathbb{E} \left[\mathbf{X}_i^* \Sigma (\mathbf{X}_i) \mathbf{X}_i^{*'} \middle| D_i \right] \right).
\end{aligned}$$

This gives

$$\begin{aligned}
&\mathbb{E} \left[\left\| \frac{1}{D_i} \mathbf{X}_i^* (\mathbf{Y}_i - \mathbf{W}_i \boldsymbol{\delta}_0) - \boldsymbol{\beta}_0(D_i) - (\boldsymbol{\beta}_0 - \boldsymbol{\beta}_0(D_i)) \right\|^3 \middle| D_i \right] \\
&\leq \left| \frac{1}{D_i^3} \right| m(D_i) + \|\boldsymbol{\beta}_0 - \boldsymbol{\beta}_0(D_i)\|^3 + \frac{3 \text{Tr} \left(\mathbb{E} \left[\mathbf{X}_i^* \Sigma (\mathbf{X}_i) \mathbf{X}_i^{*'} \middle| D_i \right] \right)}{D_i^2} \|\boldsymbol{\beta}_0 - \boldsymbol{\beta}_0(D_i)\|,
\end{aligned}$$

with $m(D_i)$ as defined in Assumption 2.3. Using the above inequality yields

$$\begin{aligned}
\mathbb{E} [\|Z_{N,i}\|^3] &\leq \int_{-\infty}^{\infty} \|\boldsymbol{\beta}_0 - \boldsymbol{\beta}_0(t)\|^3 \phi(t) dt - \int_{-h}^h \|\boldsymbol{\beta}_0 - \boldsymbol{\beta}_0(t)\|^3 \phi(t) dt \\
&\quad + \int_{-\infty}^{-h} \left(\left| \frac{1}{t^3} \right| m(t) + \frac{3 \text{Tr} \left(\mathbb{E} \left[\mathbf{X}_i^* \Sigma (\mathbf{X}_i) \mathbf{X}_i^{*'} \middle| D_i = t \right] \right)}{t^2} \|\boldsymbol{\beta}_0 - \boldsymbol{\beta}_0(t)\| \right) \phi(t) dt \\
&\quad + \int_h^{\infty} \left(\left| \frac{1}{t^3} \right| m(t) + \frac{3 \text{Tr} \left(\mathbb{E} \left[\mathbf{X}_i^* \Sigma (\mathbf{X}_i) \mathbf{X}_i^{*'} \middle| D_i = t \right] \right)}{t^2} \|\boldsymbol{\beta}_0 - \boldsymbol{\beta}_0(t)\| \right) \phi(t) dt \\
&= \mathbb{E} [\|\boldsymbol{\beta}_0 - \boldsymbol{\beta}_0(t)\|^3] - h \int_{-1}^1 \|\boldsymbol{\beta}_0 - \boldsymbol{\beta}_0(uh)\|^3 \phi(uh) du \\
&\quad + \int_{-\infty}^{-1} \left(\frac{1}{h^2} \left| \frac{1}{u^3} \right| m(uh) + \frac{3 \text{Tr} \left(\mathbb{E} \left[\mathbf{X}_i^* \Sigma (\mathbf{X}_i) \mathbf{X}_i^{*'} \middle| D_i = uh \right] \right)}{u^2 h} \|\boldsymbol{\beta}_0 - \boldsymbol{\beta}_0(uh)\| \right) \phi(uh) du \\
&\quad + \int_1^{\infty} \left(\frac{1}{h^2} \left| \frac{1}{u^3} \right| m(uh) + \frac{3 \text{Tr} \left(\mathbb{E} \left[\mathbf{X}_i^* \Sigma (\mathbf{X}_i) \mathbf{X}_i^{*'} \middle| D_i = uh \right] \right)}{u^2 h} \|\boldsymbol{\beta}_0 - \boldsymbol{\beta}_0(uh)\| \right) \phi(uh) du \\
&= \mathbb{E} [\|\boldsymbol{\beta}_0 - \boldsymbol{\beta}_0(D_i)\|^3] - 2 \left\| \boldsymbol{\beta}_0 - \boldsymbol{\beta}_0^S \right\|^3 \phi_0 h + m(0) \phi_0 h^{-2} 2 \int_1^{\infty} \frac{1}{u^3} du \\
&\quad + \frac{6 \text{Tr} \left(\mathbb{E} \left[\mathbf{X}_i^* \Sigma (\mathbf{X}_i) \mathbf{X}_i^{*'} \middle| D_i = 0 \right] \right)}{h} \left\| \boldsymbol{\beta}_0 - \boldsymbol{\beta}_0^S \right\| \phi_0 \int_1^{\infty} \frac{1}{u^2} du + o(h^{-2}) \\
&= m(0) \phi_0 h^{-2} + 6 \text{Tr} \left(\mathbb{E} \left[\mathbf{X}_i^* \Sigma (\mathbf{X}_i) \mathbf{X}_i^{*'} \middle| D_i = 0 \right] \right) \left\| \boldsymbol{\beta}_0 - \boldsymbol{\beta}_0^S \right\| \phi_0 h^{-1} \\
&\quad - 2 \left\| \boldsymbol{\beta}_0 - \boldsymbol{\beta}_0^S \right\| \phi_0 h + \mathbb{E} [\|\boldsymbol{\beta}_0 - \boldsymbol{\beta}_0(D_i)\|^3] \\
&= O(h^{-2}).
\end{aligned}$$

Using the above result we can verify the Liapunov condition. Let $a_N = \left(\frac{N}{h_N} \right)$, then

$$\frac{1}{a_N} \sum_{i=1}^N \mathbb{V}(Z_{N,i}) \rightarrow 2 \mathbb{E} \left[\mathbf{X}_i^* \Sigma (\mathbf{X}_i) \mathbf{X}_i^{*'} \middle| D_i = 0 \right] \phi_0,$$

and also

$$\begin{aligned}
\frac{\left(\sum_{i=1}^N \mathbb{E} [\|Z_{N,i} - \mathbb{E}[Z_{N,i}]\|^3]\right)^{1/3}}{(a_N)^{1/2}} &\leq \frac{\left(8 \sum_{i=1}^N \mathbb{E} [\|Z_{N,i}\|^3]\right)^{1/3}}{(a_N)^{1/2}} \\
&= O\left(\frac{(Nh^{-2})^{1/3}}{(Nh^{-1})^{1/2}}\right) \\
&= O\left((Nh)^{-1/6}\right) \\
&= o_p(1).
\end{aligned}$$

Application of the Liapunov CLT for triangular arrays, equation (39) above, and Slutsky's Theorem, then yields the follow Lemma.

Lemma A.1 *Suppose that (i) $(F_0, \boldsymbol{\delta}_0, \boldsymbol{\beta}_0(\cdot))$ satisfies (9), (ii) $\Sigma(\mathbf{x})$ is positive definite for all $\mathbf{x} \in \mathbb{X}^T$, (iii) $T = p$, and (iv) Assumptions 1.2 to 2.5 hold, then $\widehat{\boldsymbol{\beta}}_I \xrightarrow{p} \boldsymbol{\beta}_0$ with the normal limiting distribution*

$$\sqrt{Nh_N} \left(\widehat{\boldsymbol{\beta}}_I - \boldsymbol{\beta}_0\right) \xrightarrow{D} \mathcal{N}(0, 2\Upsilon_0\phi_0),$$

for $\Upsilon_0 = \mathbb{E}[\mathbf{X}_i^* \Sigma(\mathbf{X}_i) \mathbf{X}_i^{*'} | D_i = 0]$.

Large sample properties of $\widehat{\boldsymbol{\delta}}$: Recall that the non-random coefficients $\boldsymbol{\delta}_0$ are estimated by a uniform conditional linear predictor (CLP) estimator. Recentering (24) yields

$$\widehat{\boldsymbol{\delta}} - \boldsymbol{\delta}_0 = \left[\frac{1}{Nh} \sum_{i=1}^N \mathbf{1}(|D_i| \leq h) \mathbf{W}_i^{*'} \mathbf{W}_i^* \right]^{-1} \times \left[\frac{1}{Nh} \sum_{i=1}^N \mathbf{1}(|D_i| \leq h) \mathbf{W}_i^{*'} (D_i \boldsymbol{\beta}_0(\mathbf{X}_i) + \mathbf{U}_i^*) \right], \quad (47)$$

where

$$\begin{aligned}
\mathbf{U}^* &= \mathbf{Y}^* - \mathbf{W}^* \boldsymbol{\delta}_0 - D \boldsymbol{\beta}_0(\mathbf{X}) \\
&= \mathbf{X}^* (\mathbf{Y} - \mathbf{W} \boldsymbol{\delta}_0 - \mathbf{X} \boldsymbol{\beta}_0(\mathbf{X})) \\
&= \mathbf{X}^* \mathbf{U}.
\end{aligned} \quad (48)$$

First consider the expected value of the matrix being inverted in (47). Manipulations similar to those used to analyze $\widehat{\boldsymbol{\beta}}_I$ above yield

$$\begin{aligned}
\mathbb{E} [\mathbf{1}(|D_i| \leq h) \mathbf{W}_i^{*'} \mathbf{W}_i^*] &= \mathbb{E} [\mathbf{1}(|D_i| \leq h) \mathbb{E} [\mathbf{W}_i^{*'} \mathbf{W}_i^* | D_i]] \\
&= \int_{-h}^h \mathbb{E} [\mathbf{W}_i^{*'} \mathbf{W}_i^* | D_i = t] \phi(t) dt \\
&= h \int_{-1}^1 \mathbb{E} [\mathbf{W}_i^{*'} \mathbf{W}_i^* | D_i = uh] \phi(uh) du \\
&= 2\mathbb{E} [\mathbf{W}_i^{*'} \mathbf{W}_i^* | D_i = 0] \phi_0 h + o(h),
\end{aligned} \quad (49)$$

while for any fixed q -dimensional vector $\boldsymbol{\lambda}$ the variance of a quadratic form in that matrix satisfies

$$\begin{aligned}
\mathbb{V} \left[\frac{1}{Nh} \sum_{i=1}^N \mathbf{1}(|D_i| \leq h) (\boldsymbol{\lambda}' \mathbf{W}_i^{*'} \mathbf{W}_i^* \boldsymbol{\lambda}) \right] &\leq \frac{1}{Nh^2} \mathbb{E} \left[\mathbf{1}(|D_i| \leq h) \mathbb{E} \left[\|\mathbf{W}_i^*\|^4 \mid D_i \right] \right] \|\boldsymbol{\lambda}\|^4 \\
&= \frac{2\mathbb{E} \left[\|\mathbf{W}_i^*\|^4 \mid D_i = 0 \right] \phi_0 \|\boldsymbol{\lambda}\|^4}{Nh} + o\left(\frac{1}{Nh}\right) \\
&= o\left(\frac{1}{Nh}\right)
\end{aligned} \tag{50}$$

under Assumptions 2.3 and 2.5 so that

$$\frac{1}{Nh} \sum_{i=1}^N \mathbf{1}(|D_i| \leq h) \mathbf{W}_i^{*'} \mathbf{W}_i^* = 2\mathbb{E} \left[\mathbf{W}_i^{*'} \mathbf{W}_i^* \mid D_i = 0 \right] \phi_0 + o_p(1). \tag{51}$$

Now redefine $Z_{N,i}$ to equal the term entering the summation in the numerator of (47):

$$Z_{N,i} \equiv \mathbf{1}(|D_i| \leq h) \mathbf{W}_i^{*'} (D_i \boldsymbol{\beta}_0(\mathbf{X}_i) + \mathbf{U}_i^*).$$

Using the fact that $\mathbb{E}[\mathbf{W}_i^{*'} \mathbf{U}_i^* \mid D_i] = \mathbb{E}[\mathbf{W}_i^{*'} \mathbf{X}_i^* \mathbb{E}[\mathbf{U} \mid \mathbf{X}] \mid D_i] = 0$ yields an expected value of $Z_{N,i}$ equal to

$$\begin{aligned}
\mathbb{E}[Z_{N,i}] &= \mathbb{E} \left[\mathbf{1}(|D_i| \leq h) \mathbf{W}_i^{*'} (D_i \boldsymbol{\beta}_0(\mathbf{X}_i) + \mathbf{U}_i^*) \right] \\
&= \mathbb{E} \left[\mathbf{1}(|D_i| \leq h) D_i \mathbb{E} \left[\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) \mid D_i \right] \right] \\
&= \int_{-h}^h t \mathbb{E} \left[\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) \mid D_i = t \right] \phi(t) dt \\
&= h \int_{-1}^1 uh \mathbb{E} \left[\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) \mid D_i = uh \right] \phi(uh) du \\
&= \mathbb{E} \left[\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) \mid D_i = 0 \right] \phi_0 h^2 \int_{-1}^1 u du \\
&\quad + \left\{ \frac{\partial \mathbb{E} \left[\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) \mid D_i = 0 \right]}{\partial d} \phi_0 + \mathbb{E} \left[\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) \mid D_i = 0 \right] \phi_0' \right\} h^3 \int_{-1}^1 u^2 du + o(h^3) \\
&= \frac{2}{3} \left\{ \frac{\partial \mathbb{E} \left[\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) \mid D_i = 0 \right]}{\partial d} \phi_0 + \mathbb{E} \left[\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) \mid D_i = 0 \right] \phi_0' \right\} h^3 + o(h^3),
\end{aligned}$$

where we use the following Taylor approximation and Assumption 1.2 and 2.3 in deriving the second to last equality above:

$$\begin{aligned}
\mathbb{E} \left[\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) \mid D_i = uh \right] uh \phi(uh) &= 0 + \mathbb{E} \left[\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) \mid D_i = 0 \right] \phi_0 uh \\
&\quad + \left\{ \frac{\partial \mathbb{E} \left[\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) \mid D_i = 0 \right]}{\partial d} \phi_0 + \mathbb{E} \left[\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) \mid D_i = 0 \right] \phi_0' \right\} (uh)^2 + o(h^2).
\end{aligned}$$

The numerator (47) therefore equals

$$\begin{aligned}
&\frac{1}{Nh} \sum_{i=1}^N \mathbf{1}(|D_i| \leq h) \mathbf{W}_i^{*'} (D_i \boldsymbol{\beta}_0(\mathbf{X}_i) + \mathbf{U}_i^*) \\
&= \frac{2}{3} \left\{ \frac{\partial \mathbb{E} \left[\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) \mid D_i = 0 \right]}{\partial d} \phi_0 + \mathbb{E} \left[\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) \mid D_i = 0 \right] \phi_0' \right\} h^2 + o_p(h^2).
\end{aligned} \tag{52}$$

Using the ratio of (52) and (51) yields a bias expression for $\widehat{\boldsymbol{\delta}} - \boldsymbol{\delta}_0$ of

$$\widehat{\boldsymbol{\delta}} - \boldsymbol{\delta}_0 = \frac{1}{3} \mathbb{E} [\mathbf{W}_i^{*'} \mathbf{W}_i^* | D_i = 0]^{-1} \times \left\{ \frac{\partial \mathbb{E} [\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) | D_i = 0]}{\partial d} + \mathbb{E} [\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) | D_i = 0] \frac{\phi_0'}{\phi_0} \right\} h^2 + o_p(h^2). \quad (53)$$

This implies that we can center the asymptotic distribution of $\sqrt{N}h_N (\widehat{\boldsymbol{\delta}} - \boldsymbol{\delta}_0)$ at zero by choosing h_N such that $(Nh_N)^{1/2} h_N^2 \rightarrow 0$ (Assumption 2.5).

Now consider the variance of $Z_{N,i}$. As before we proceed by evaluating the two terms in the variance decomposition (43) separately. The first of the two terms evaluates to

$$\begin{aligned} \mathbb{V}(\mathbb{E}[Z_{N,i} | D_i]) &= \mathbb{V}(\mathbf{1}(|D_i| \leq h) D_i \mathbb{E}[\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) | D_i]) \\ &= \mathbb{E} \left[\mathbf{1}(|D_i| \leq h) D_i^2 \mathbb{E}[\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) | D_i] \mathbb{E}[\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) | D_i]' \right] \\ &\quad - \mathbb{E}[\mathbf{1}(|D_i| \leq h) D_i \mathbb{E}[\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) | D_i]] \mathbb{E}[\mathbf{1}(|D_i| \leq h) D_i \mathbb{E}[\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) | D_i]]'. \end{aligned}$$

Evaluating the two expectations entering the above expressions yields

$$\begin{aligned} \mathbb{E}[\mathbf{1}(|D_i| \leq h) D_i \mathbb{E}[\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) | D_i]] &= \int_{-h}^h t \mathbb{E}[\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) | D_i = t] \phi(t) dt \\ &= h^2 \int_{-1}^1 u \mathbb{E}[\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) | D_i = uh] \phi(uh) du \\ &= o(h^2), \end{aligned}$$

and

$$\begin{aligned} &\mathbb{E} \left[\mathbf{1}(|D_i| \leq h) D_i^2 \mathbb{E}[\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) | D_i] \mathbb{E}[\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) | D_i]' \right] \\ &= \int_{-h}^h t^2 \mathbb{E}[\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) | D_i = t] \mathbb{E}[\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) | D_i = t]' \phi(t) dt \\ &= h \int_{-1}^1 (uh)^2 \mathbb{E}[\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) | D_i = uh] \mathbb{E}[\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) | D_i = uh]' \phi(uh) du \\ &= \mathbb{E}[\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) | D_i = 0] \mathbb{E}[\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) | D_i = 0]' \phi_0 h^3 \int_{-1}^1 u^2 du + o(h^3) \\ &= \frac{2}{3} \mathbb{E}[\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) | D_i = 0] \mathbb{E}[\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) | D_i = 0]' \phi_0 h^3 + o(h^3). \end{aligned}$$

We conclude that $\mathbb{V}(\mathbb{E}[Z_{N,i} | D_i]) = o(h^3)$.

Now consider the second term in (43). The conditional variance, using the conditional moment restriction (9), is

$$\begin{aligned} \mathbb{V}(Z_{N,i} | D_i) &= \mathbf{1}(|D_i| \leq h) \mathbb{V}(D_i \mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) + \mathbf{W}_i^{*'} \mathbf{U}_i^* | D_i) \\ &= \mathbf{1}(|D_i| \leq h) D_i^2 \mathbb{V}(\mathbf{W}_i^{*'} \boldsymbol{\beta}_0(\mathbf{X}_i) | D_i) + \mathbf{1}(|D_i| \leq h) \mathbb{V}(\mathbf{W}_i^{*'} \mathbf{U}_i^* | D_i). \end{aligned}$$

Using an ANOVA decomposition to evaluate $\mathbb{V}(\mathbf{W}_i^{*'} \mathbf{U}_i^* | D_i)$ gives

$$\begin{aligned} \mathbb{V}(\mathbf{W}_i^{*'} \mathbf{U}_i^* | D_i) &= \mathbb{E}[\mathbb{V}(\mathbf{W}_i^{*'} \mathbf{U}_i^* | \mathbf{X}_i) | D_i] + \mathbb{V}(\mathbb{E}[\mathbf{W}_i^{*'} \mathbf{U}_i^* | \mathbf{X}_i] | D_i) \\ &= \mathbb{E}[\mathbf{W}_i^{*'} \boldsymbol{\Sigma}^* \mathbf{X} \mathbf{X}' \mathbf{W}_i^* | D_i] + 0, \end{aligned}$$

and hence

$$\begin{aligned}
\mathbb{E} [\mathbf{1} (|D_i| \leq h) \mathbb{E} [\mathbf{W}_i^{*'} \mathbf{X}^* \Sigma (\mathbf{X}) \mathbf{X}^{*'} \mathbf{W}_i^* | D_i]] &= \int_{-h}^h \mathbb{E} [\mathbf{W}_i^{*'} \mathbf{X}^* \Sigma (\mathbf{X}) \mathbf{X}^{*'} \mathbf{W}_i^* | D_i = t] \phi (t) dt \\
&= h \int_{-1}^1 \mathbb{E} [\mathbf{W}_i^{*'} \mathbf{X}^* \Sigma (\mathbf{X}) \mathbf{X}^{*'} \mathbf{W}_i^* | D_i = uh] \phi (uh) du \\
&= 2\mathbb{E} [\mathbf{W}_i^{*'} \mathbf{X}^* \Sigma (\mathbf{X}) \mathbf{X}^{*'} \mathbf{W}_i^* | D_i = 0] \phi_0 h + o(h).
\end{aligned}$$

Similarly

$$\begin{aligned}
\mathbb{E} [\mathbf{1} (|D_i| \leq h) D_i^2 \mathbb{V} (\mathbf{W}_i^{*'} \beta_0 (\mathbf{X}_i) | D_i)] &= \int_{-h}^h t^2 \mathbb{V} (\mathbf{W}_i^{*'} \beta_0 (\mathbf{X}_i) | D_i = t) \phi (t) dt \\
&= h \int_{-1}^1 (uh)^2 \mathbb{V} (\mathbf{W}_i^{*'} \beta_0 (\mathbf{X}_i) | D_i = uh) \phi (uh) du \\
&= \mathbb{V} (\mathbf{W}_i^{*'} \beta_0 (\mathbf{X}_i) | D_i = 0) \phi_0 h^3 \int_{-1}^1 u^2 du + o(h^3) \\
&= \frac{2}{3} \mathbb{V} (\mathbf{W}_i^{*'} \beta_0 (\mathbf{X}_i) | D_i = 0) \phi_0 h^3 + o(h^3).
\end{aligned}$$

Collecting terms we conclude that

$$\mathbb{V} (Z_{N,i}) = 2\mathbb{E} [\mathbf{W}_i^{*'} \mathbf{X}^* \Sigma (\mathbf{X}) \mathbf{X}^{*'} \mathbf{W}_i^* | D_i = 0] \phi_0 h + o(h). \quad (54)$$

Applying Liapunov's Central Limit Theorem for triangular arrays, we have

$$\frac{1}{\sqrt{N h_N}} \sum_{i=1}^N Z_{N,i} \xrightarrow{D} \mathcal{N} (0, 2\mathbb{E} [\mathbf{W}_i^{*'} \mathbf{X}^* \Sigma (\mathbf{X}) \mathbf{X}^{*'} \mathbf{W}_i^* | D_i = 0] \phi_0).$$

Slutsky's Theorem and (51) above then give the following Lemma.

Lemma A.2 *Suppose that (i) $(F_0, \delta_0, \beta_0 (\cdot))$ satisfies (9), (ii) $\Sigma (\mathbf{x})$ is positive definite for all $\mathbf{x} \in \mathbb{X}^T$, (iii) $T = p$, and (iv) Assumptions 1.2 to 2.5 hold, then $\widehat{\delta} \xrightarrow{P} \delta_0$ with the normal limiting distribution*

$$\sqrt{N h_N} (\widehat{\delta} - \delta_0) \xrightarrow{D} \mathcal{N} \left(0, \frac{\Lambda_0}{2\phi_0} \right),$$

where

$$\Lambda_0 = \mathbb{E} [\mathbf{W}_i^{*'} \mathbf{W}_i^* | D_i = 0]^{-1} \mathbb{E} [\mathbf{W}_i^{*'} \mathbf{X}^* \Sigma (\mathbf{X}) \mathbf{X}^{*'} \mathbf{W}_i^* | D_i = 0] \mathbb{E} [\mathbf{W}_i^{*'} \mathbf{W}_i^* | D_i = 0]^{-1}.$$

Large sample properties of $\widehat{\beta}$: The following lemma characterizes the probability limit of $\widehat{\Xi}_N$.

Lemma A.3 *If $(F_0, \delta_0, \beta_0 (\cdot))$ satisfies (9) and Assumptions 1.2 to 2.5 hold we have $\widehat{\Xi}_N \xrightarrow{P} \Xi_0$, where*

$$\begin{aligned}
\Xi_0 &= \lim_{h_N \downarrow 0} \mathbb{E} [\mathbf{1} (|D_i| > h_N) \mathbf{X}_i^{-1} \mathbf{W}_i]. \\
&\equiv \lim_{N \rightarrow \infty} \Xi_N.
\end{aligned}$$

To verify this result, we must first establish that the defined limit exists. For h_N sufficiently small, we can decompose the expectation Ξ_N as

$$\begin{aligned}\Xi_N &= \mathbb{E} [\mathbf{1}(|D_i| \geq u_0) \mathbf{X}_i^{-1} \mathbf{W}_i] + \mathbb{E} [\mathbf{1}(h_N < |D_i| < u_0) \mathbf{X}_i^{-1} \mathbf{W}_i] \\ &= O(1) + \mathbb{E} [\mathbf{1}(h_N < |D_i| < u_0) D_i^{-1} \mathbf{W}_i^*] \\ &= O(1) + \int_{h_N}^{u_0} \left(\frac{\xi(u) - \xi(-u)}{u} \right) du,\end{aligned}$$

where

$$\xi(u) \equiv \phi(u) \mathbb{E}[W_i^* | D_i = u]$$

is twice continuously differentiable for $|u| < u_0$ by Assumption 1.2. Using the Taylor's series expansion

$$\xi(u) - \xi(-u) = \xi(0) - \xi(0) + 2 \frac{d\xi(0)}{du} \cdot u + \left[\frac{d^2\xi(u^*)}{du^2} - \frac{d^2\xi(-u^*)}{du^2} \right] \cdot \left(\frac{u^2}{2} \right),$$

for u^* some intermediate value between h_N and u_0 , it follows that

$$\left| \mathbf{1}(h_N < u < u_0) \cdot \left(\frac{\xi(u) - \xi(-u)}{u} \right) \right| \leq \mathbf{1}(u \leq u_0) \left[2 \left\| \frac{d\xi(0)}{du} \right\| + \max_{|u| \leq u_0} \left\| \frac{d^2\xi(u)}{du^2} \right\| \cdot u_0 \right],$$

and since the right-hand side is integrable, $\Xi_N \rightarrow \Xi_0$ by dominated convergence. Then, taking λ to be an arbitrary (fixed) q -vector, verification that $\widehat{\Xi}_N \xrightarrow{p} \Xi_0$ follows from the convergence of the covariance matrix of the numerator of $\widehat{\Xi}_N \lambda$ to zero:

$$\begin{aligned}\left\| \mathbb{V} \left[\frac{1}{N} \sum_{i=1}^N \mathbf{1}(|D_i| > h_N) \mathbf{X}_i^{-1} \mathbf{W}_i \lambda \right] \right\| &= \left\| \mathbb{V} \left[\frac{1}{N} \sum_{i=1}^N \mathbf{1}(|D_i| > h_N) D_i^{-1} \mathbf{W}_i^* \lambda \right] \right\| \\ &\leq \frac{2 \|\lambda\|^2}{N} \mathbb{E} \left[\mathbf{1}(|D_i| > h_N) |D_i|^{-2} \|\mathbf{W}_i^*\|^2 \right] \\ &= \frac{2 \|\lambda\|^2}{N} \left[\mathbb{E} \left[\mathbf{1}(|D_i| \geq u_0) |D_i|^{-2} \|\mathbf{W}_i^*\|^2 \right] + \mathbb{E} \left[\mathbf{1}(h_N < |D_i| < u_0) |D_i|^{-2} \|\mathbf{W}_i^*\|^2 \right] \right] \\ &= O\left(\frac{1}{N}\right) + \frac{2 \|\lambda\|^2}{N} \int_{h_N}^{u_0} \left(\frac{\zeta(u) + \zeta(-u)}{u^2} \right) du,\end{aligned}$$

where $\|\mathbf{W}_i^*\|^2 \equiv \text{tr} [\mathbf{W}_i^* \mathbf{W}_i^*]$ and

$$\zeta(u) \equiv \phi(u) \mathbb{E}[\|\mathbf{W}_i^*\|^2 | D_i = u]$$

is bounded for $|u| < u_0$, so

$$\begin{aligned}\left\| \mathbb{V} \left[\frac{1}{N} \sum_{i=1}^N \mathbf{1}(|D_i| > h_N) \mathbf{X}_i^{-1} \mathbf{W}_i \lambda \right] \right\| &\leq O\left(\frac{1}{N}\right) + \frac{4 \|\lambda\|^2}{N} \left(\max_{|u| \leq u_0} \|\zeta(u)\| \right) \int_{h_N}^{u_0} \left(\frac{1}{u^2} \right) du \\ &= O\left(\frac{1}{N}\right) + \frac{4 \|\lambda\|^2}{N} \left(\max_{|u| \leq u_0} \|\zeta(u)\| \right) \left[\frac{1}{h_N} - \frac{1}{u_0} \right] \\ &= O\left(\frac{1}{N h_N}\right) \\ &= o(1)\end{aligned}$$

under Assumption 2.5.

Lemmas A.1, A.2, and A.3 as well as the decomposition (37) then give Theorem 2.1.

A.2 MSE-optimal bandwidth sequence:

The MSE-optimal bandwidth sequence given in equation (29) of the main text may be derived as follows. Using (42), (53) and Lemma A.3 yields a leading asymptotic bias term of $2(\boldsymbol{\beta}_0 - \boldsymbol{\beta}_0^S) \phi_0 h$. The asymptotic variance is given in the statement of Theorem 2.1. The asymptotic MSE of $\widehat{\boldsymbol{\beta}}$ is thus

$$4(\boldsymbol{\beta}_0 - \boldsymbol{\beta}_0^S) (\boldsymbol{\beta}_0 - \boldsymbol{\beta}_0^S)' \phi_0^2 h_N^2 + \frac{2\Upsilon_0 \phi_0 + \frac{\Xi_0 \Lambda_0 \Xi_0'}{2\phi_0}}{N h_N}.$$

Minimizing this object with respect to h_N gives the result in the main text.

A.3 Proof of Theorem 2.2:

Rewriting equation (30) we have

$$\widehat{V} = \left[\frac{1}{N} \sum_{i=1}^N \mathbf{Q}'_i \mathbf{R}_i \right]^{-1} \times \left[\frac{h_N}{N} \sum_{i=1}^N \mathbf{Q}'_i \widehat{\mathbf{U}}_i^+ \widehat{\mathbf{U}}_i^{+'} \mathbf{Q}_i \right] \times \left[\frac{1}{N} \sum_{i=1}^N \mathbf{Q}'_i \mathbf{R}_i \right]^{-1} \quad (55)$$

for $\mathbf{U}_i^+ = \mathbf{Y}_i^* - \mathbf{R}_i \boldsymbol{\theta}_0$, with $\widehat{\boldsymbol{\theta}} = \left[\frac{1}{N} \sum_{i=1}^N \mathbf{Q}'_i \mathbf{R}_i \right]^{-1} \times \left[\frac{1}{N} \sum_{i=1}^N \mathbf{Q}'_i \mathbf{Y}_i^* \right]$, for $\boldsymbol{\theta} = (\boldsymbol{\delta}', \boldsymbol{\beta}')'$ and

$$\mathbf{Q}_i = \left(h_N^{-1} \mathbf{1}(|D_i| \leq h_N) \mathbf{W}_i^*, \frac{\mathbf{1}(|D_i| > h_N)}{D_i} I_p \right), \quad \mathbf{R}_i = (\mathbf{W}_i^*, \mathbf{1}(|D_i| > h_N) D_i I_p).$$

The dependence of \mathbf{Q}_i and \mathbf{R}_i on h_N is suppressed to simplify the notation.

Starting with the ‘Jacobian’ term in \widehat{V} we get, by the definitions of \mathbf{Q}_i , \mathbf{R}_i , and $\widehat{\Xi}_N$,

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N \mathbf{Q}'_i \mathbf{R}_i &= \frac{1}{N} \sum_{i=1}^N \begin{pmatrix} h_N^{-1} \mathbf{1}(|D_i| \leq h_N) \mathbf{W}_i^{*'} \\ \frac{\mathbf{1}(|D_i| > h_N)}{D_i} I_p \end{pmatrix} \begin{pmatrix} \mathbf{W}_i^* & \mathbf{1}(|D_i| > h_N) D_i I_p \end{pmatrix} \\ &= \begin{pmatrix} \frac{1}{N h_N} \sum_{i=1}^N \mathbf{1}(|D_i| \leq h_N) \mathbf{W}_i^{*'} \mathbf{W}_i^* & \mathbf{0}_Q \mathbf{0}'_P \\ \frac{1}{N} \sum_{i=1}^N \mathbf{1}(|D_i| > h_N) \mathbf{X}_i^{-1} \mathbf{W}_i & \frac{1}{N} \sum_{i=1}^N \mathbf{1}(|D_i| > h_N) I_p \end{pmatrix} \\ &\xrightarrow{p} \begin{pmatrix} 2\mathbb{E}[\mathbf{W}_i^{*'} \mathbf{W}_i^* | D_i = 0] \phi_0 & \mathbf{0} \\ \Xi_0 & I_p \end{pmatrix} \end{aligned}$$

by (51) and A.3. Decomposing the middle term $h_N \sum_{i=1}^N \mathbf{Q}'_i \widehat{\mathbf{U}}_i^+ \widehat{\mathbf{U}}_i^{+'} \mathbf{Q}_i / N$, yields

$$\begin{aligned} &\left\| \frac{h_N}{N} \sum_{i=1}^N \mathbf{Q}'_i \widehat{\mathbf{U}}_i^+ \widehat{\mathbf{U}}_i^{+'} \mathbf{Q}_i - \frac{h_N}{N} \sum_{i=1}^N \mathbf{Q}'_i \mathbf{U}_i^+ \mathbf{U}_i^{+'} \mathbf{Q}_i \right\| \\ &\leq \frac{2h_N}{N} \sum_{i=1}^N \|\mathbf{U}_i^+\| \|\mathbf{Q}_i\|^2 \|\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0\| + \frac{h_N}{N} \sum_{i=1}^N \|\mathbf{Q}_i\|^2 \|\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0\|^2 \\ &= \frac{2}{N h_N} \sum_{i=1}^N \left(\mathbf{1}(|D_i| \leq h_N) \|\mathbf{W}_i^*\|^2 \|\mathbf{U}_i^+\| + \mathbf{1}(|D_i| > h_N) \frac{h_N^2 \|\mathbf{U}_i^+\|}{|D_i|^2} \right) \|\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0\| \\ &\quad + \frac{1}{N h_N} \sum_{i=1}^N \left(\mathbf{1}(|D_i| \leq h_N) \|\mathbf{W}_i^*\|^2 + \mathbf{1}(|D_i| > h_N) \frac{h_N^2}{|D_i|^2} \right) \|\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0\|^2. \end{aligned} \quad (56)$$

By a similar argument as for (50) above,

$$\begin{aligned}\mathbb{E} \left[h_N^{-1} \mathbf{1}(|D_i| \leq h_N) \|\mathbf{W}_i^*\|^2 (\|\mathbf{U}_i^+\| + 1) \right] &= \mathbb{E} \left[h_N^{-1} \mathbf{1}(|D_i| \leq h_N) \mathbb{E} \left[\|\mathbf{W}_i^*\|^2 (\|\mathbf{U}_i^+\| + 1) \mid D_i \right] \right] \\ &\rightarrow 2\mathbb{E} \left[\|\mathbf{W}_i^*\|^2 (\|\mathbf{U}_i^+\| + 1) \mid D_i = 0 \right] \phi_0,\end{aligned}$$

and the same reasoning as (44) yields

$$\begin{aligned}\mathbb{E} \left[\mathbf{1}(|D_i| > h_N) \|\mathbf{W}_i^*\|^2 (\|\mathbf{U}_i^+\| + 1) \frac{h_N}{|D_i|^2} \right] &= \mathbb{E} \left[h_N^{-1} \mathbf{1}(|D_i| > h_N) \frac{h_N}{|D_i|^2} \mathbb{E} \left[\|\mathbf{W}_i^*\|^2 (\|\mathbf{U}_i^+\| + 1) \mid D_i \right] \right] \\ &\rightarrow 2\mathbb{E} \left[\|\mathbf{W}_i^*\|^2 (\|\mathbf{U}_i^+\| + 1) \mid D_i = 0 \right] \phi_0.\end{aligned}$$

Thus, by Markov's inequality, (56) yields

$$\begin{aligned}\frac{h_N}{N} \sum_{i=1}^N \mathbf{Q}_i' \widehat{\mathbf{U}}_i^+ \widehat{\mathbf{U}}_i^{+'} \mathbf{Q}_i &= \frac{h_N}{N} \sum_{i=1}^N \mathbf{Q}_i' \mathbf{U}_i^+ \mathbf{U}_i^{+'} \mathbf{Q}_i + \mathcal{O}_p \left(\|\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0\| \right) \\ &= \frac{h_N}{N} \sum_{i=1}^N \mathbf{Q}_i' \mathbf{U}_i^+ \mathbf{U}_i^{+'} \mathbf{Q}_i + o_p(\mathbf{1}).\end{aligned}$$

Finally,

$$\begin{aligned}\frac{h_N}{N} \sum_{i=1}^N \mathbf{Q}_i' \mathbf{U}_i^+ \mathbf{U}_i^{+'} \mathbf{Q}_i &= \left(\begin{array}{c} \frac{1}{Nh_N} \sum_{i=1}^N \mathbf{1}(|D_i| \leq h_N) \mathbf{W}_i^{*'} (\mathbf{Y}_i^* - \mathbf{W}_i^* \boldsymbol{\delta}_0) (\mathbf{Y}_i^* - \mathbf{W}_i^* \boldsymbol{\delta}_0)' \mathbf{W}_i^* \\ \mathbf{0}_P \mathbf{0}'_Q \\ \mathbf{0}_Q \mathbf{0}'_P \\ \frac{h_N}{N} \sum_{i=1}^N \mathbf{1}(|D_i| > h_N) \{ \mathbf{X}_i^{-1} (\mathbf{Y}_i - \mathbf{W}_i \boldsymbol{\delta}_0) - \boldsymbol{\beta}_0 \} \{ \mathbf{X}_i^{-1} (\mathbf{Y}_i - \mathbf{W}_i \boldsymbol{\delta}_0) - \boldsymbol{\beta}_0 \}' \end{array} \right).\end{aligned}$$

By the calculations yielding (54), the expected value of the first diagonal submatrix is

$$\mathbb{E} \left[h_N^{-1} \mathbf{1}(|D_i| \leq h_N) \mathbf{W}_i^{*'} (\mathbf{Y}_i^* - \mathbf{W}_i^* \boldsymbol{\delta}_0) (\mathbf{Y}_i^* - \mathbf{W}_i^* \boldsymbol{\delta}_0)' \mathbf{W}_i^* \right] = 2\mathbb{E} \left[\mathbf{W}_i^{*'} \mathbf{X}^* \Sigma(\mathbf{X}) \mathbf{X}^* \mathbf{W}_i^* \mid D_i = 0 \right] \phi_0 + o(1),$$

while the variance of any term in the matrix is $O((Nh_N)^{-1}) = o(1)$ by Assumptions 2.3 and 2.5 (with $r \leq 8$).

Similarly, the expectation of the second diagonal submatrix is

$$\begin{aligned}\mathbb{E} \left[h_N \mathbf{1}(|D_i| > h_N) \{ \mathbf{X}_i^{-1} (\mathbf{Y}_i - \mathbf{W}_i \boldsymbol{\delta}_0) - \boldsymbol{\beta}_0 \} \{ \mathbf{X}_i^{-1} (\mathbf{Y}_i - \mathbf{W}_i \boldsymbol{\delta}_0) - \boldsymbol{\beta}_0 \}' \right] &= 2\mathbb{E} \left[\mathbf{X}^* \Sigma(\mathbf{X}) \mathbf{X}^* \mid D = 0 \right] \phi_0 + o(1) \\ &= 2\Upsilon_0 \phi_0 + o(1),\end{aligned}$$

while similar calculations to those leading to (46) yield

$$\begin{aligned}\left\| \mathbb{V} \left[h_N \mathbf{1}(|D_i| > h_N) \{ \mathbf{X}_i^{-1} (\mathbf{Y}_i - \mathbf{W}_i \boldsymbol{\delta}_0) - \boldsymbol{\beta}_0 \} \{ \mathbf{X}_i^{-1} (\mathbf{Y}_i - \mathbf{W}_i \boldsymbol{\delta}_0) - \boldsymbol{\beta}_0 \}' \right] \right\| &\leq \frac{4m_8(0)\phi_0}{Nh_n} + o\left(\frac{1}{Nh_n} \right) \\ &= o(1),\end{aligned}$$

where $m_8(u)$ is defined in Assumption 2.3, Thus

$$\frac{h_N}{N} \sum_{i=1}^N \mathbf{Q}'_i \mathbf{U}_i^+ \mathbf{U}_i^{+'} \mathbf{Q}_i \xrightarrow{p} \begin{pmatrix} 2\mathbb{E}[\mathbf{W}_i^{*'} \mathbf{X}^* \Sigma(\mathbf{X}) \mathbf{X}^{*'} \mathbf{W}_i^* | D_i = 0] \phi_0 & \mathbf{Q}_Q \mathbf{Q}'_P \\ \mathbf{Q}_P \mathbf{Q}'_Q & 2\Upsilon_0 \phi_0 + o(1) \end{pmatrix},$$

ensuring $\widehat{V} \xrightarrow{p} V_0$.

References

- [1] Abrevaya, Jason. (2000). “Rank estimation of a generalized fixed-effects regression model,” *Journal of Econometrics* 95 (1): 1 - 23.
- [2] Altonji, Joseph G. and Rosa L. Matzkin (2005). “Cross section and panel data estimators for nonseparable models with endogenous regressors,” *Econometrica* 73 (4): 1053 - 1102.
- [3] Angrist, Joshua D. and Alan B. Krueger. (1999). “Empirical strategies in labor economics,” *Handbook of Labor Economics* 3 (1): 1277 - 1366 (O. C. Ashenfelter & D. Card, Eds.). Amsterdam: North-Holland.
- [4] Andrews, Donald W. K. and Marcia M. A. Schafgans. (1998). “Semiparametric estimation of the intercept of a sample selection model,” *Review of Economic Studies* 65 (3): 497 - 517.
- [5] Arellano, Manuel and Stephanie Bonhomme. (2009). “Identifying distributional characteristics in random coefficients panel data models,” *Mimeo*.
- [6] Arellano, Manuel and Raquel Carrasco. (2003). “Binary choice panel data models with predetermined variables,” *Journal of Econometrics* 115 (1): 125 - 157.
- [7] Arellano, Manuel and Bo Honoré. (2001). “Panel data models: some recent developments,” *Handbook of Econometrics* 5: 3229 - 3298 (J. Heckman & E. Leamer, Eds.). Amsterdam: North-Holland.
- [8] Behrman, Jere R. and Anil B. Deolalikar. (1987). “Will developing country nutrition improve with income? A case study for rural south India,” *Journal of Political Economy* 95 (3): 492 - 507.
- [9] Bester, C. Alan and Christian Hansen. (2009). “Identification of Marginal Effects in a Non-parametric Correlated Random Effects Model,” *Journal of Business and Economic Statistics* 27 (2): 235 - 250.
- [10] Blundell, Richard W. and James L. Powell. (2003). “Endogeneity in nonparametric and semi-parametric regression models,” *Advances in Economics and Econometrics: Theory and Applications II*: 312 - 357. (M. Dewatripont, L.P. Hansen, S. J. Turnovsky, Eds.). Cambridge: Cambridge University Press.

- [11] Bonhomme, Stephane. (2010). "Functional differencing," *Mimeo*.
- [12] Bouis, Howarth E. (1994). "The effect of income on demand for food in poor countries: are our food consumption databases giving us reliable estimates?" *Journal of Development Economics* 44(1): 199-226.
- [13] Bouis, Howarth E. and Lawrence J. Haddad. (1992). "Are estimates of calorie-income elasticities too high? A recalibration of the plausible range," *Journal of Development Economics* 39 (2): 333 - 364.
- [14] Browning, Martin and Jesus Carro. (2007). "Heterogeneity and microeconometrics modelling," *Advances in Economics and Econometrics: Theory and Applications III*: 47 - 74. (R. Blundell, W. Newey & T. Persson, Eds.). Cambridge: Cambridge University Press.
- [15] Card, David. (1996). "The effect of unions on the structure of wages: a longitudinal analysis," *Econometrica* 64 (4): 957 - 979.
- [16] Chamberlain, Gary. (1980). "Analysis of covariance with qualitative data," *Review of Economic Studies* 47 (1): 225 - 238.
- [17] Chamberlain, Gary. (1982). "Multivariate regression models for panel data," *Journal of Econometrics* 18 (1): 5 - 46.
- [18] Chamberlain, Gary. (1984). "Panel data," *Handbook of Econometrics 2*: 1247 - 1318 (Z. Griliches & M.D. Intriligator, Eds.). Amsterdam: North-Holland.
- [19] Chamberlain, Gary. (1986). "Asymptotic efficiency in semi-parametric models with censoring," *Journal of Econometrics* 32 (2): 189 - 218.
- [20] Chamberlain, Gary. (1992). "Efficiency bounds for semiparametric regression," *Econometrica* 60 (3): 567 - 596.
- [21] Chamberlain, Gary. (2010). "Binary response models for panel data: identification and information," *Econometrica* 78 (1): 159 - 168.
- [22] Chernozhukov, Victor, Iván Fernández-Val, Jinyong Hahn and Whitney Newey. (2008). "Identification and estimation of marginal effects in nonlinear panel data models," *CEMMAP Working Paper CWP25/08*.
- [23] Dasgupta, Partha. (1993). *An Inquiry into Well-Being and Destitution*. Oxford: Oxford University Press.
- [24] Deaton, Angus and Salman Zaidi. (2002). "Guidelines for constructing consumption aggregates for welfare analysis," *LSMS Working Paper Number 135*.

- [25] Engle, Robert F., C. W. J. Granger, John Rice, and Andrew Weiss. (1986). "Semiparametric estimates of the relation between weather and electricity sales," *Journal of the American Statistical Association* 81 (394): 310 - 320.
- [26] Food and Agricultural Organization (FAO). (2001). "Human energy requirements: report of a joint FAO/WHO/UNU expert consultation," *FAO Food and Nutrition Technical Report Series* 1.
- [27] Food and Agricultural Organization (FAO). (2006). *The State of Food Insecurity in the World 2006*. Rome: Food and Agricultural Organization.
- [28] Graham, Bryan S., Guido W. Imbens, Geert Ridder. (2009). "Complementarity and aggregate implications of assortative matching: a nonparametric analysis," *Mimeo*.
- [29] Heckman, James, Hidehiko Ichimura, Jeffrey Smith and Petra Todd. (1998). "Characterizing selection bias using experimental data," *Econometrica* 66 (5): 1017 - 1098.
- [30] Heckman, James J. (1990). "Varieties of selection bias," *American Economic Review* 80 (2): 313 - 18.
- [31] Honoré, Bo E. (1992). "Trimmed LAD and least squares estimation of truncated and censored regression models with fixed effects," *Econometrica* 60 (3): 533 - 565.
- [32] Honoré, Bo E. and Ekaterini Kyriazidou. (1997). "Panel data discrete choice models with lagged dependent variables," *Econometrica* 68 (4): 839 - 874.
- [33] Horowitz, Joel L. (1992). "A smoothed maximum score estimator for the binary response model," *Econometrica* 60 (3): 505 - 531.
- [34] Imbens, Guido W. (2007). "Nonadditive models with endogenous regressors," *Advances in Economics and Econometrics: Theory and Applications III*: 17 - 46. (R. Blundell, W. Newey & T. Persson, Eds.). Cambridge: Cambridge University Press.
- [35] Imbens, Guido W. and Whitney K. Newey. (2009). "Identification and estimation of triangular simultaneous equations models without additivity," *Econometrica* 77 (5): 1481 - 1512.
- [36] Instituto Nacional de Estadísticas y Censos (INEC). (2005). *Metodología de construcción del agregado de consumo, de las líneas de pobreza y del agregado de ingreso en Nicaragua*. Retrieved January 22, 2008, from the World Wide Web: <http://www.inec.gob.ni/Pobreza/pubmetodol.htm>.
- [37] Instituto de Nutrición de Centro América y Panamá (INCAP) and Organización Panamericana la Salud (OPS). (2000). *Tabla de Composición de Alimentos de Centroamérica*. Retrieved January 2008, from the World Wide Web: <http://www.tabladealimentos.net/tca/TablaAlimentos/inicio.html>

- [38] International Food Policy Research Institute (IFPRI). (2005). *Nicaraguan RPS evaluation data (2000-02): overview and description of data files (April 2005 Release)*. Washington D.C.: International Food Policy Research Institute.
- [39] Khan, Shakeeb and Elie Tamer. (2009). "Irregular identification, support conditions, and inverse weight estimation," *Mimeo*.
- [40] Kyriazidou, Ekaterini. (1997). "Estimation of a panel sample selection model," *Econometrica* 65 (6): 1335 - 1364.
- [41] Manski, Charles F. (1987). "Semiparametric analysis of random effects linear models from binary panel data," *Econometrica* 55 (2): 357 - 362.
- [42] Mundlak, Yair. (1961). "Empirical production function free of management bias," *Journal of Farm Economics* 43 (1): 44 - 56.
- [43] Mundlak, Yair. (1978a). "On the pooling of time series and cross section data," *Econometrica* 46 (1): 69 - 85.
- [44] Mundlak, Yair. (1978b). "Models with variable coefficients: integration and extension," *Annales de l'Insee* 30-31: 483 - 510.
- [45] Newey, Whitney K. (1994a). "The asymptotic variance of semiparametric estimators," *Econometrica* 62 (6): 1349 - 1382.
- [46] Newey, Whitney K. (1994b). "Kernel estimation of partial means and a general variance estimator," *Econometric Theory* 10 (2): 233 - 253.
- [47] Pagan, Adrian and Aman Ullah. (1999). *Nonparametric Econometrics*. Cambridge: Cambridge University Press.
- [48] Serfling, Robert J. (1980). *Approximation Theorems of Mathematical Statistics*. New York: John Wiley & Sons, Inc.
- [49] Smith, Lisa C. and Ali Subandoro. (2007). *Measuring food security using household expenditure surveys*. Washington D.C.: International Food Policy Research Institute.
- [50] Strauss, John and Duncan Thomas. (1990). "The shape of the calorie-expenditure curve," *Yale University Economic Growth Center Discussion Paper No. 595*.
- [51] Strauss, John and Duncan Thomas. (1995). "Human resources: empirical modeling of household and family decisions," *Handbook of Development Economics* 3 (1): 1883 - 2023. (J. Behrman & T.N. Srinivasan). Amsterdam: North-Holland.
- [52] Subramanian, Shankar and Angus Deaton. (1996). "The demand for food and calories," *Journal of Political Economy* 104 (1): 133 - 162.

- [53] Wolfe, Barbara L. and Jere R. Behrman. (1983). "Is income overrated in determining adequate nutrition?" *Economic Development and Cultural Change* 31 (3): 525 - 549.
- [54] Wooldridge, Jeffrey M. (2005a). "Unobserved heterogeneity and estimation of average partial effects," *Identification and Inference for Econometric Models: Essays in Honor of Thomas Rothenberg*: 27 - 55 (D.W.K. Andrews & J.H. Stock, Eds.). Cambridge: Cambridge University Press.
- [55] Wooldridge, Jeffrey M. (2005b). "Fixed-effects and related estimators for correlated-random coefficient and treatment-effect panel data models," *Review of Economics and Statistics* 87 (2): 385 - 390.
- [56] World Bank. (2002). *Basic information document: Nicaragua living standards measurement survey 1998*. Washington D.C.: The World Bank.
- [57] World Bank. (2003). *Nicaragua Poverty Assessment: Raising Welfare and Reducing Vulnerability*. Washington D.C.: The World Bank.