

Botosaru, Irene; Muris, Chris; Sokullu, Senay

Working Paper

Time-varying linear transformation models with fixed effects and endogeneity for short panels

cemmap working paper, No. CWP06/22

Provided in Cooperation with:

Institute for Fiscal Studies (IFS), London

Suggested Citation: Botosaru, Irene; Muris, Chris; Sokullu, Senay (2022) : Time-varying linear transformation models with fixed effects and endogeneity for short panels, cemmap working paper, No. CWP06/22, Centre for Microdata Methods and Practice (cemmap), London, <https://doi.org/10.47004/wp.cem.2022.0622>

This Version is available at:

<https://hdl.handle.net/10419/260387>

Standard-Nutzungsbedingungen:

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

Terms of use:

Documents in EconStor may be saved and copied for your personal and scholarly purposes.

You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.

If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.

Time-varying linear transformation models with fixed effects and endogeneity for short panels

Irene Botosaru
Chris Muris
Senay Sokullu

The Institute for Fiscal Studies
Department of Economics, UCL

cemmap working paper CWP06/22



Economic
and Social
Research Council

Time-Varying Linear Transformation Models with Fixed Effects and Endogeneity for Short Panels*

Irene Botosaru, Chris Muris, and Senay Sokullu[†]

January 22, 2022

Abstract

This paper considers a class of fixed- T nonlinear panel models with time-varying link function, fixed effects, and endogenous regressors. We establish sufficient conditions for the identification of the regression coefficients, the time-varying link function, the distribution of counterfactual outcomes, and certain (time-varying) average partial effects. We propose estimators for these objects and study their asymptotic properties. We show the relevance of our model by estimating the effect of teaching practices on student attainment as measured by test scores on standardized tests in mathematics and science. We use data from the Trends in International Mathematics and Science Study, and show that both traditional and modern teaching practices have positive effects of similar magnitudes on the performance of U.S. students on standardized tests in math and science.

*We would like to thank Stéphane Bonhomme, Xavier D'Haultfoeuille, Jean-Pierre Florens, Laura Liu, Thierry Magnac for comments and suggestions, and participants at various conferences and seminars. We would also like to thank Jan Bietenbeck for sharing the data.

[†]Contact: Botosaru and Muris are at McMaster University, Department of Economics, Canada, and Sokullu is at Bristol University, School of Economics, UK. Emails: botosari@mcmaster.ca, muerisc@mcmaster.ca, senay.sokullu@bristol.ac.uk. We gratefully acknowledge financial support from the Social Sciences and Humanities Research Council of Canada under grants IG 435-2021-0778. This research was undertaken, in part, thanks to funding from the Canada Research Chairs Program.

1 Introduction

We consider a nonlinear panel model with endogeneity, where the outcome for individual i at time t can be written as a time-varying transformation of a latent linear variable with endogenous regressors. That is, the observed outcome is specified as:

$$Y_{it} = h_t(Y_{it}^*) = h_t(\alpha_i + \bar{X}_{it}\bar{\beta} + U_{it}), i = 1, \dots, n, t = 1, \dots, T, \quad (1)$$

where $Y_{it} \in \mathcal{Y}_t \subseteq \mathbb{R}$ is a continuous random variable, $h_t : \mathbb{R} \rightarrow \mathcal{Y}_t$ is an unknown, strictly monotonic transformation that varies with t , $\alpha_i \in \mathbb{R}$ is an unobserved individual effect, $\bar{\beta} \in \mathbb{R}^{k+1}$ is a vector of regression coefficients,¹ $U_{it} \in \mathbb{R}$ is a stochastic error, and

$$\bar{X}_{it} = (X_{0it}, X_{1it}, \dots, X_{kit}) \in \mathcal{X}_t \subseteq \mathbb{R}^{k+1}$$

is a vector of explanatory variables that may be endogenous in the sense that:

$$E(U_{it} | X_{i1}, \dots, X_{iT}) \neq 0 \text{ for all } t = 1, \dots, T. \quad (2)$$

In particular, the covariates are not required to be strictly exogenous. The individual effect α_i is a fixed effect, in the sense that it can be arbitrarily correlated with \bar{X}_{it} . The distribution of U_{it} is left unspecified except for a conditional mean restriction in (5) below.

We are interested in the identification and estimation of $\bar{\beta}$ and h_t , $t = 1, \dots, T$, and of the distribution of the counterfactual outcomes, the average structural function, and certain partial effects. All these parameters are time-varying, whenever the transformation is time-varying.

An example of our framework is a nonlinear version of the standard linear dynamic panel model:

$$Y_{i0} = h_0(\phi_i, \tilde{X}_{i0}), \quad (3)$$

$$Y_{it} = h_t(\alpha_i + \tilde{X}_{it}\tilde{\beta} + \rho Y_{i,t-1} + U_{it}), t = 1, \dots, T, \quad (4)$$

¹In Appendix C, we allow for a nonparametric function of the covariates, e.g. $\rho(\bar{X}_{it})$.

where ϕ_i captures additional individual-specific unobserved heterogeneity in the initial condition. This specification is nested by (1) by setting $\bar{X}_{it} = (\tilde{X}_{it}, Y_{i,t-1})$, $\bar{\beta} = (\tilde{\beta}, \rho)$, and it nests the linear dynamic panel model when $h_t(v) = v$ for $t \geq 1$. To the best of our knowledge, this is the first paper showing identification of $(\tilde{\beta}, \rho, (h_t)_{t \geq 1})$ and the distribution of the counterfactual outcomes for this class of models. We study identification of this model in Appendix D.

Our solution to the endogeneity problem relies on the existence of instrumental variables, $Z_{it} \in \mathcal{Z} \subseteq \mathbb{R}^q$, $q \geq k + 1$, that satisfy the following mean independence condition: for any $z \in \mathcal{Z}$,

$$E(U_{it} - U_{it-1} | Z_{it} = z) = 0 \text{ for all } t = 2, \dots, T. \quad (5)$$

This conditional mean restriction is mild: it allows the instrumental variable at time t to affect the level of the errors at time t as long as it does so in a time-homogeneous way;² and it does not impose any restrictions on the serial dependence or heteroskedasticity of the stochastic errors.

First, we identify the regression coefficient and h_t using insights from the non-parametric instrumental variables (NPIV) literature, in particular Fève and Florens (2014) and Florens and Sokullu (2017). Second, we then identify the distribution of counterfactual outcomes, extending previous results in Botosaru and Muris (2017); Botosaru et al. (2021). This does not require knowledge or identification of the distribution of fixed effects, and only uses the regression coefficient and h_t identified in step 1. Partial effects are identified en passant. Third, we propose estimators based on Tikhonov regularization, and show that the regression coefficient estimator converges at the parametric rate, even if the link functions do not. Our estimators for the average partial effects also attain the \sqrt{n} rate. To the best of our knowledge, this is the first paper to derive such results for nonlinear transformation models with fixed effects and fixed- T .

²In contrast, in nonseparable models this type of conditional mean restriction is not sufficient for identification of the structural parameters or of the partial effects. In general, it is stronger assumptions, such as independence between U_t and X_t for each and all t that are maintained in those models.

We show the relevance of our approach by estimating the effect of teaching practices on student attainment. We use data from the Trends in International Mathematics and Science Study, where the test score Y_{it} of each student i is observed across two standardized tests t , one in mathematics ($t = 1$) and one in science ($t = 2$). Test scores are relative ranks, and thus do not have an ordinal scale. Therefore, it is important that approaches using test scores as outcome variables be invariant to monotone transformations, see, e.g., Cunha and Heckman (2008), Bonhomme and Sauder (2011). Our approach is invariant, and furthermore allows the monotone transformation to be different across math and science. On top of that, we accommodate fixed effects and endogenous regressors. Our application demands all of these features, as we explain in Section 7. We are not aware of any existing approach that delivers this combination of features. We find that both traditional and modern teaching practices have positive effects of similar magnitudes on test scores.³ This is different from other studies that use standardized test scores that find that modern teaching practices have almost nonexistent effect on test scores, see, e.g. Bietenbeck (2014).

Relative contribution. We are not aware of any existing work that combines the following four features: (1) fixed- T , (2) fixed effects, (3) nonlinearity via a time-varying link function, and (4) endogenous regressors. To the best of our knowledge, our identification and estimation results are therefore novel. Our paper contributes to at least four literatures: nonlinear panel models, dynamic panel models, transformation models, and partial effects in nonlinear panel models.

First, our results contribute to the literature on nonlinear panel models with fixed-effects and fixed- T . Within this literature, the maintained assumption for identification of the structural parameters, $\bar{\beta}, h_t$, has been of strict exogeneity of the regressors. For example, Abrevaya (1999) considered the outcome equation (1) and proposed an estimator for the regression coefficient. Botosaru et al. (2021) studied identification and estimation of the time-varying link function. In this paper, we relax the assumption of strictly and weakly exogenous regressors.

That it is challenging to deal with endogenous regressors in nonlinear panel mod-

³Freyberger (2018) uses the same data but studies a different issue.

els with fixed effects and fixed- T is evident from reviews of the literature in Arellano and Honoré (2001); Arellano and Bonhomme (2011). A notable exception is Altonji and Matzkin (2005), who consider an outcome equation that nests ours, and who also allow for endogenous regressors. Altonji and Matzkin (2005) make progress by imposing restrictions on the distribution of (α_i, X_i) . In contrast, we obtain identification without imposing such restrictions.

Within the nonlinear panel literature described above, Botosaru et al. (2021) is closest in spirit. The main differences with the specification there are that in this paper the link function in (1) is assumed to be strictly monotonic (so that the results of this paper apply to continuous outcomes only), and that the covariates are allowed to be endogenous.

Second, we also contribute to the literature on dynamic panel models. As a special case of our general result, we analyze a nonlinear version of the linear dynamic panel model, see 3-4 above, and Appendix D. There is a large literature on *linear* dynamic panels, see Bun and Sarafidis (2015) for a review. These models are very popular in applied practice. For example, the early key contributions by Arellano and Bond (1991); Blundell and Bond (1998) have 35,172 resp. 24,622 Google Scholar citations at the time of writing. That the combination of a dynamic structure with a nonlinear structure is difficult to handle is clear from the literature on dynamic discrete choice models with fixed effects, cf. Honoré and Kyriazidou (2000); Honoré and Weidner (2020); Muris et al. (2020); Honoré et al. (2021). For example, whether it is possible to accommodate time trends and endogenous regressors in such models is an open question. In contrast, our model allows for the transformations to vary over time in an arbitrary fashion. Furthermore, we accommodate additional endogenous regressors.

Third, we contribute to the literature on panel transformation models with endogenous regressors. We extend previous work by Florens et al. (2012), Fève and Florens (2014), and Florens and Sokullu (2017) to nonlinear panel models. The main difference with Fève and Florens (2014) is that our specification allows for a time-varying link function, while the analysis in Florens et al. (2012) and Florens and Sokullu (2017) applies to cross-sectional data (see also Fève and Florens (2010)). On the other hand, both Florens and Sokullu (2017) and Fève and Florens (2014) allow

for a nonparametric function of observed endogenous covariates, i.e. $\rho(X_{it})$ instead of $X_{it}\bar{\beta}$. In Appendix C, we explain how our analysis can be extended to allow for this possibility. Note that in this case, our framework nests that of Fève and Florens (2014).

Other related works have addressed the problem of endogeneity in transformation models via arguments based on special regressors, e.g. Chiappori et al. (2015), and control functions, e.g. Vanhems and Van Keilegom (2019). These papers consider a cross-sectional set-up, so the transformation function is not indexed by time and, importantly, there are no fixed effects. We consider a panel data setting and our identification argument uses instrumental variables.⁴

We adopt an inverse problem approach to derive sufficient conditions for the identification of $\bar{\beta}$ and h_t . As such, we use concepts from the NPIV literature such as invertibility of an operator, completeness, and measurable separability. Our proposed estimator is based on Tikhonov regularization and follows closely the procedure in Florens and Sokullu (2017).

Finally, there is a growing number of papers deriving conditions for the identification of marginal and partial effects for static and dynamic discrete choice models, e.g., Aguirregabiria and Carro (2021), Aguirregabiria et al. (2021), Davezies et al. (2021), Dobronyi et al. (2021), Liu et al. (2021), Pakel and Weidner (2021). None covers our specification with unknown and time-varying transformation. Botosaru and Muris (2017) consider partial effects for the class of models where the outcome equation is as in (1) with strictly exogenous X_t . Chernozhukov et al. (2013) consider a nonseparable outcome equation, but do not allow for endogenous regressors, arbitrary time-varyingness, and they require boundedness of the dependent variable and discreteness of the regressors. An important point made by recent papers starting with Botosaru and Muris (2017) is that, even in nonlinear models, average partial effects can be identified without identification of the distribution of the fixed effects.

Organization. The paper is organized as follows. In Section 2 we derive sufficient conditions for the identification of β and h_t , while in Section 3 we show that the

⁴Other related work, set in a cross-sectional setting, can be found in the literature review in e.g. Birke et al. (2017).

results of Section 2 are sufficient for the identification of a menu of partial effects, all of which are time-varying. In Section 4 we introduce our estimators for β and h_t and derive asymptotic results, while in Section 6 we present Monte-Carlo results that document the finite sample properties of our estimators. In Section 7 we present the empirical application to the effects of teaching practices on test scores. The Appendix contains all the proofs and extensions.

Notation. For a random variable V with support \mathcal{V} , we let $L_{\mathcal{V}}^2$ denote the space of functions $g : \mathcal{V} \rightarrow \mathbb{R}$ such that $E|g(V)|^2 < \infty$. We denote by g^{-1} the inverse of an arbitrary, invertible function $g : \mathbb{R} \rightarrow \mathbb{R}$. We use \otimes to denote the tensor product. We let $\mathcal{C} : \mathbb{R} \rightarrow \mathbb{R}$ denote a bounded, continuous, symmetric, univariate kernel function of order m , i.e. $\int \mathcal{C}(u)du = 1$, $\int u^j \mathcal{C}(u)du = 0$ for all $j = 1, \dots, m-1$, and $\int u^m \mathcal{C}(u)du < \infty$ and $\int \mathcal{C}^2(u)du < \infty$. We let $\mathcal{K} : \mathbb{R}^d \rightarrow \mathbb{R}$ denote a multivariate kernel function defined as the product kernel $\mathcal{K}(w) = \prod_{k=1}^d \mathcal{C}(w_k)$. For an operator K between two Hilbert spaces, we denote by $\mathcal{R}(K)$ the range of the operator and by $\mathcal{R}(K)^\perp$ its orthogonal complement.

2 Identification

Our identification results require at least two time periods, so we let $T = 2$ in what follows. We drop the i subscript in this section.

Assumption 1. *For each $t = 1, 2$, $h_t : \mathbb{R} \rightarrow \mathcal{Y}_t$ is strictly monotonic.*

Shape restrictions such as monotonicity are quite common in the literature on transformation models. Assumption 1 allows us to work with h_t^{-1} , the inverse of h_t , $t = 1, 2$.

Assumption 2. *(i) The first element of $\bar{\beta}$ is normalized to 1, i.e. $\bar{\beta} = (1, \beta)$, $\beta \in \mathbb{R}^k$. (ii) $E(h_1^{-1}(Y_1)) = 0$.*

Without parametric restrictions on h_t and on the distribution of U_t , the outcome in (1) follows a semiparametric single index specification. Therefore, both a scale

normalization, 2(i), and a location normalization, 2(ii), are needed for identification, see also Horowitz (2009).

Letting $X_{0t} \in \mathbb{R}$ denote the covariate associated with the normalized coefficient from Assumption 2(i) and

$$X_t \equiv (X_{1t}, X_{2t}, \dots, X_{kt}) \in \mathbb{R}^k,$$

the outcome equation (1) can then be written as

$$Y_t = h_t(\alpha + X_{0t} + X_t\beta + U_t).$$

Assumption 3. *There exist random variables $Z \in \mathcal{Z}$ such that for any $z \in \mathcal{Z}$,*

$$E(U_2 - U_1 | Z = z) = 0.$$

As Assumption 3 is made on the difference over time in the errors, it does not require that $E(U_t | Z) = 0$, thus allowing Z to enter the outcome equation, as long as it does so in a time-homogeneous way.

Let $\Delta X \equiv X_2 - X_1$, $\Delta X_0 \equiv X_{02} - X_{01}$, and $r(z) \equiv E(\Delta X_0 | Z = z)$. Assumptions 1, 2, and 3, obtain that, for any $z \in \mathcal{Z}$,

$$E(h_2^{-1}(Y_2) - h_1^{-1}(Y_1) - \Delta X\beta | Z = z) = r(z). \quad (6)$$

Equation (6) is an integral equation for the parameters of interest $(h_1^{-1}, h_2^{-1}, \beta)$. This shows that the three parameters can be characterized by the functional equation $K(h_1^{-1}, h_2^{-1}, \beta) = r$, where $K : \mathcal{H}_1 \rightarrow \mathcal{H}_2$ is a multilinear integral operator and $\mathcal{H}_1, \mathcal{H}_2$ are function spaces defined below.

Assumption 4. (i) $\mathcal{H}_1 = L_{\mathcal{Y}_1}^2 \otimes L_{\mathcal{Y}_2}^2 \otimes \mathbb{R}^k$ and $\mathcal{H}_2 = L_{\mathcal{Z}}^2$; (ii) $r \in L_{\mathcal{Z}}^2$; and (iii) *The joint distribution of $(Y_1, Y_2, \Delta X, Z)$ is dominated by the product of its marginal distributions, and its density is square integrable w.r.t. the product of marginals.*

A few remarks are in order. First, Assumption 4 allows us to define K as the

conditional expectation operator:

$$K : L_{\mathcal{Y}_1}^2 \otimes L_{\mathcal{Y}_2}^2 \otimes \mathbb{R}^k \rightarrow L_{\mathcal{Z}}^2, \quad (7)$$

that maps

$$(h_1^{-1}, h_2^{-1}, \beta) \mapsto E(h_2^{-1}(Y_2) - h_1^{-1}(Y_1) - \Delta X \beta \mid Z = \cdot).$$

Second, Assumptions 4(i) and (ii) are satisfied provided that the variances of $U_2 - U_1$ and ΔX are finite. Assumption 4(iii) guarantees that the conditional density $f_{Y_1, Y_2, \Delta X \mid Z}$ is well-defined and that the operator K is Hilbert-Schmidt.

If the model is correctly specified, then $r \in \mathcal{R}(K)$ and the functional equation in (6) has at least one solution for $(h_1^{-1}, h_2^{-1}, \beta)$. The following assumption guarantees uniqueness of the solution.

Assumption 5. *K is such that for any $(\delta_1^{-1}, \delta_2^{-1}, b) \in L_{\mathcal{Y}_1}^2 \otimes L_{\mathcal{Y}_2}^2 \otimes \mathbb{R}^k$,*

$$K(\delta_1^{-1}, \delta_2^{-1}, b) = 0 \text{ a.s. } F_Z \implies \delta_1^{-1} = 0 \text{ a.s. } F_{Y_1}, \delta_2^{-1} = 0 \text{ a.s. } F_{Y_2}, b = 0.$$

Assumption 5 is an injectivity-like assumption on the multilinear operator K .⁵ It is possible to state sufficient assumptions for it in terms of linear operators by using the linearity of the expectation operator. To this end, define the following linear operators:

$$K_{y_t} : L_{\mathcal{Y}_t}^2 \rightarrow L_{\mathcal{Z}}^2 : h_t^{-1} \mapsto E(h_t^{-1}(Y_t) \mid Z = \cdot), \quad t = 1, 2, \quad (8)$$

$$K_x : \mathbb{R}^k \rightarrow L_{\mathcal{Z}}^2 : \beta \mapsto E(\Delta X \beta \mid Z = \cdot), \quad (9)$$

so that $K(h_1^{-1}, h_2^{-1}, \beta) = K_{y_2} h_2^{-1} - K_{y_1} h_1^{-1} - K_x \beta$. The following Assumption 6 is sufficient for Assumption 5.

Assumption 6. *(i) Each operator K_{y_1}, K_{y_2}, K_x is injective; (ii) $\mathcal{R}(K_{y_1}) \cap \mathcal{R}(K_{y_2}) \cap \mathcal{R}(K_x) = \{0\}$.*

⁵Florens et al. (2012) make a similar assumption on a bilinear operator in the proof of their Theorem 2.1.

Assumption 6(i) is equivalent to assuming that each Y_t is strongly identified by Z , i.e. for any $\gamma \in L^2_{Y_t}$,

$$E(\gamma(Y_t)|Z) = 0 \text{ a.s. } F_Z \implies \gamma(Y_t) = 0 \text{ a.s. } F_{Y_t}, t = 1, 2,$$

and that the matrix $E[E(\Delta X|Z)E(\Delta X'|Z)]$ has full rank. The strong identification assumption is the standard L^2 -completeness assumption usually invoked in the NPIV literature.⁶

Assumption 6(ii) is implied by $Y_1, Y_2, \Delta X$ each being strongly identified by Z , and by measurable separability of $(Y_1, Y_2, \Delta X)$.⁷ Measurable separability is a high-level assumption that rules out a linear relationship between Y_1, Y_2 , and ΔX . The assumption fails if there exists an *additive* functional relationship between Y_1, Y_2 , and ΔX , see, e.g., Newey et al. (1999). Lemma 1 in Appendix A.1 establishes low-level assumptions for measurable separability.

Theorem 1 (Identification). *Suppose that $(Y_1, Y_2, X_{01}, X_{02}, X_1, X_2, Z)$ follow the model described by 1, 2, and 5. Let Assumptions 1, 2, 3, 4, and either 5 or 6 hold. Then h_1, h_2, β are identified.*

Proof. The proof can be found in Appendix B.1. □

⁶If each Y_t and Z are continuous, have the same dimension, and support equal to a rectangle then the completeness condition holds generically in the sense of Andrews (2017), see also Chen et al. (2014), Newey and Powell (2003). It may be possible to consider weaker sufficient conditions by adapting Proposition 2.2 in d'Haultfoeulle (2010) to the case of square integrable functions. Other papers that provide sufficient conditions for completeness are d'Haultfoeulle (2011), Andrews (2017), and Hu and Shiu (2018).

⁷The random variables $(Y_1, Y_2, \Delta X)$ are measurably separable when, e.g., for any $\delta_t \in L^2_{Y_t}$ and any $b \in \mathbb{R}^k$ if

$$\delta_2(Y_2) - \delta_1(Y_1) - \Delta X b = 0 \text{ a.s. } F_{Y_1, Y_2, \Delta X},$$

then there exist constants $c_t \in \mathbb{R}$, $t = 1, 2$, such that

$$\delta_t(Y_t) = c_t \text{ a.s. } F_{Y_t}, \quad t = 1, 2.$$

3 Partial effects

Building on results on partial effects for nonlinear panel models in Botosaru and Muris (2017, Section 3.2), we show that identification of the structural parameters (β, h_1, h_2) implies identification of certain partial effects. The main differences between the current setting and that of Botosaru and Muris (2017) are that, here, (i) the transformation function is assumed to be invertible, and (ii) the regressors are allowed to be endogenous. We first show that the distribution of the counterfactual outcome at t is identified, and then we show identification of the average partial effects.

Denote the *counterfactual outcome* by

$$Y_{it}(x) = h_t(\alpha_i + x\beta + U_{it}), \quad (10)$$

which is the outcome of person i at time t under $X_{it} = x$, while holding (α_i, U_{it}) fixed. Note that $Y_{it} = Y_{it}(X_{it})$. The distribution of the counterfactual outcome at time t is defined as:

$$P(Y_{it}(x) \leq y) = P(h_t(\alpha_i + x\beta + U_{it}) \leq y), \quad (11)$$

for any value of (y, x) .

Because the distribution of (Y_{it}, X_{it}) is identified from the data, and (h_t, β) have been identified previously, this counterfactual quantity is also identified, and is given by:

$$\begin{aligned} P(Y_{it}(x) \leq y) &= P(h_t(\alpha_i + x\beta + U_{it}) \leq y) \\ &= P(\alpha_i + x\beta + U_{it} \leq h_t^{-1}(y)) \\ &= P(\alpha_i + X_{it}\beta + U_{it} \leq h_t^{-1}(y) - (x - X_{it})\beta) \\ &= P(h_t(\alpha_i + X_{it}\beta + U_{it}) \leq h_t(h_t^{-1}(y) - (x - X_{it})\beta)) \\ &= P(Y_{it} \leq h_t(h_t^{-1}(y) - (x - X_{it})\beta)). \end{aligned}$$

The first equality uses the outcome equation for our model; the second uses strict increasingness of h_t ; the third adds $(X_{it} - x)\beta$ on both sides of the inequality; the fourth applies the strictly monotone function to both sides; and the final equality substitutes the observed Y_{it} .

In our empirical illustration, we build on this result to obtain regressor effects. Rather than looking at a fixed value of the regressors x , we will look at a counterfactual value of the covariates that change the k th covariate by 1 unit. This counterfactual value of the regressors can be written as $X_{it} + e_k$, where e_k is the unit vector of the appropriate length, with a 1 in entry k , and zeros elsewhere. Following the sequence of equalities above, we obtain

$$P(Y_{it}(X_{it} + e_k) \leq y) = P(Y_{it} \leq h_t(h_t^{-1}(y) - \beta_k)),$$

where β_k is the value of the k th coefficient. Then, the difference in distributions

$$\begin{aligned} \tau_{k,t}(y) &\equiv P(Y_{it}(X_{it} + e_k) \leq y) - P(Y_{it} \leq y) \\ &= P(Y_{it} \leq h_t(h_t^{-1}(y) - \beta_k)) - P(Y_{it} \leq y) \end{aligned} \quad (12)$$

is the partial effect for regressor k .

In our empirical illustration, we are interested in the average effect rather than the distribution of the counterfactual:

$$\delta_{k,t} \equiv E[Y_{it}(X_{it} + e_k) - Y_{it}] \quad (13)$$

$$= E[h_t(h_t^{-1}(Y_{it}) + \beta_k)] - E[Y_{it}]. \quad (14)$$

The final expression depends only on observable or identified quantities, and is therefore identified. The resulting $\delta_{k,t}$ is the average partial effect (APE): it is the average change in the outcome at time t when the k th covariate changes by one unit at time t , ceteris paribus. The expression in (13) suggests using the sample analog of the right hand side as an estimator for the APE.

4 Estimation

Our identification argument naturally gives rise to a system of three normal equations based on the three linear operators K_{y_1}, K_{y_2}, K_x and their adjoints. However, the resulting system of normal equations is unnecessarily complicated.⁸ Instead, we follow Florens and Sokullu (2017) and work with the following bilinear operator:

$$K_y : \tilde{L}_{y_1}^2 \otimes L_{y_2}^2 \rightarrow L_Z^2 : (h_1^{-1}, h_2^{-1}) \mapsto K_{y_2} h_2^{-1} - K_{y_1} h_1^{-1},$$

where $\tilde{L}_{y_1}^2 \equiv \{h_1^{-1} \in L_{y_1}^2 : E(h_1^{-1}(Y_1) = 0)\}$,⁹ so that Assumption 2(ii) holds. We define the following dual operators or adjoints:

$$K_y^* : L_Z^2 \rightarrow \tilde{L}_{y_1}^2 \otimes L_{y_2}^2 : \psi \mapsto \begin{pmatrix} E[\psi(Z) | Y_2 = \cdot] \\ -\mathbb{P}E[\psi(Z) | Y_1 = \cdot] \end{pmatrix},$$

$$K_x^* : L_Z^2 \rightarrow \mathbb{R}^k : \psi \mapsto E[\psi(Z) \Delta X],$$

where \mathbb{P} is the operator that projects functions from $L_{y_1}^2$ to $\tilde{L}_{y_1}^2$. We can then write (6) as:

$$K_y(h_1^{-1}, h_2^{-1}) - K_x \beta = r, \quad (15)$$

and we can project the problem in (15) onto the parameter spaces, $\tilde{L}_{y_1}^2 \otimes L_{y_2}^2$ and \mathbb{R}^k , using the dual operators above. The functions $(h_1^{-1}, h_2^{-1}, \beta)$ are then characterized as solutions to the following system of normal equations:

$$K_y^* K_y(h_1^{-1}, h_2^{-1}) = K_y^* r + K_y^* K_x \beta, \quad (16)$$

$$K_x^* K_y(h_1^{-1}, h_2^{-1}) = K_x^* r + K_x^* K_x \beta. \quad (17)$$

Letting I be the identity operator in $L_{y_1}^2 \otimes L_{y_2}^2$, and $P_x \equiv K_x (K_x^* K_x)^{-1} K_x^*$ and $P_y \equiv K_y (K_y^* K_y)^{-1} K_y^*$ be the orthogonal projection operators onto the closure of

⁸We show in Appendix A.2 that the system of three normal equations based on the three linear operators that we use for identification is identical to the system with two normal equations based on a bilinear operator that we use for estimation.

⁹This operator is injective given Assumption 6.

the range of K_x and K_y , respectively, the above linear system is equivalent to:

$$K_y^* (I - P_x) r = K_y^* (I - P_x) K_y (h_1^{-1}, h_2^{-1}), \quad (18)$$

$$K_x^* (I - P_y) r = K_x^* (I - P_y) K_x \beta. \quad (19)$$

The parameters of interest can in principle be obtained from equations (18) and (19), after replacing K_x, K_y, K_x^*, K_y^* , and r by their sample analogues, call them $\hat{K}_x, \hat{K}_y, \hat{K}_x^*, \hat{K}_y^*$, and \hat{r} , respectively. For example,

$$\left(\hat{h}_1^{-1}, \hat{h}_2^{-1} \right)'_{naive} = \left(\hat{K}_y^* (I - \hat{P}_x) \hat{K}_y \right)^{-1} \hat{K}_y^* (I - \hat{P}_x) \hat{r}, \quad (20)$$

$$\hat{\beta}_{naive} = \left(\hat{K}_x^* (I - \hat{P}_y) \hat{K}_x \right)^{-1} \hat{K}_x^* (I - \hat{P}_y) \hat{r}. \quad (21)$$

It is well known in the literature on inverse problems, see, e.g. Carrasco and Florens (2011), Horowitz (2011), Centorrino et al. (2017), Florens and Sokullu (2017), Babii and Florens (2020), that estimating (h_1^{-1}, h_2^{-1}) and β by naively inverting the sample analogues of (18) and (19) as in (20) and (21) is a statistically ill-posed problem, in the sense that the naive estimators are not stable with respect to estimation error. We then use regularization to smooth out discontinuities due to inversion. In this paper, we will use Tikhonov regularization.

Letting γ_n be a regularization parameter, such that $\gamma_n \rightarrow 0$ at a rate defined in Assumption (11) below, the regularized estimators are given by:

$$\left(\hat{h}_1^{-1}, \hat{h}_2^{-1} \right)' = \left(\gamma_n I + \hat{K}_y^* (I - \hat{P}_x) \hat{K}_y \right)^{-1} \hat{K}_y^* (I - \hat{P}_x) \Delta X_0, \quad (22)$$

$$\hat{\beta} = \left(\hat{K}_x^* (I - \hat{P}_y^{\gamma_n}) \hat{K}_x \right)^{-1} \hat{K}_x^* (I - \hat{P}_y^{\gamma_n}) \hat{r}, \quad (23)$$

where $\hat{P}_y^{\gamma_n} \equiv \hat{K}_y (\gamma_n I + \hat{K}_y^* \hat{K}_y)^{-1} \hat{K}_y^*$ is the regularized projection operator P_y . Note that, although estimation of β is also affected by regularization, we show that $\hat{\beta}$ is \sqrt{n} -consistent and asymptotically normal.¹⁰

¹⁰Note that a single regularization parameter is introduced. Although it is possible to allow for two different regularization parameters, one for h_1^{-1} and one for h_2^{-1} , our asymptotic theory in the

Letting $\{Y_{1i}, Y_{2i}, \Delta X_{i0} = X_{02i} - X_{01i}, \Delta X_i = X_{2i} - X_{1i}, Z_i\}_{i=1}^n$ be a random sample from a population conformable to our assumptions in Section 2, we consider the following nonparametric estimators for the operators in (22) and (23):

$$\begin{aligned}\hat{K}_x \gamma(z) &= \frac{1}{nb_z^q} \frac{1}{\hat{f}_Z(z)} \sum_{i=1}^n \Delta X_i' \gamma \mathcal{K} \left(\frac{Z_i - z}{b_z} \right), \text{ for all } \gamma \in \mathbb{R}^k, \\ \hat{K}_y(g_1, g_2)(z) &= \frac{1}{\hat{f}_Z(z)} \left(\int g_2(y) \hat{f}_{Y_2, Z}(y, z) dy - \int g_1(y) \hat{f}_{Y_1, Z}(y, z) dy \right), \text{ for all } g_1, g_2 \in \tilde{L}_{Y_1}^2 \otimes L_{Y_2}^2, \\ \hat{K}_x^* g_3(z) &= \frac{1}{nb_z^q} \sum_{i=1}^n \Delta X_i \int g_3(z) \mathcal{K} \left(\frac{Z_i - z}{b_z} \right) dz, \text{ for all } g \in L_Z^2, \\ \hat{K}_y^* g_4(y_1, y_2) &= \left(\begin{array}{c} \frac{1}{\hat{f}_{Y_2}(y_2)} \int g_4(z) \hat{f}_{Y_2, Z}(y_2, z) dz \\ - \frac{1}{\hat{f}_{Y_1}(y_1)} \int g_4(z) \hat{f}_{Y_1, Z}(y_1, z) dz \end{array} \right), \text{ for all } g_4 \in L_Z^2, \\ \hat{r}(z) &= \frac{1}{nb_z^q} \frac{1}{\hat{f}_Z(z)} \sum_{i=1}^n \Delta X_{i0} \mathcal{K} \left(\frac{Z_i - z}{b_z} \right),\end{aligned}$$

where \mathcal{K} is a multivariate kernel function (see Notation in Section 1), b_z a bandwidth parameter that is assumed to be the same for each of the q components of Z and which approaches 0 as $n \rightarrow \infty$ at a rate specified in Assumption (11) below, and where:

$$\hat{f}_{Y_t, Z}(y, z) = \frac{1}{nb_{y_t} b_z^q} \sum_{i=1}^n \mathcal{C} \left(\frac{Y_{it} - y}{b_{y_t}} \right) \mathcal{K} \left(\frac{Z_i - z}{b_z} \right), t = 1, 2,$$

$$\hat{f}_{Y_t}(y) = \frac{1}{nb_{y_t}} \sum_{i=1}^n \mathcal{C} \left(\frac{Y_{it} - y}{b_{y_t}} \right), t = 1, 2,$$

next section requires them to converge to zero at the same rate. Hence, for the sake of exposition, we assume that the two regularization parameters are equal.

$$\hat{f}_Z(z) = \frac{1}{nb_z^q} \sum_{i=1}^n \mathcal{K}\left(\frac{Z_i - z}{b_z}\right),$$

where \mathcal{C} is a univariate kernel function (see Notation in Section 1) and b_{y_t} is a bandwidth parameter that approaches 0 as $n \rightarrow \infty$ at a rate specified in Assumption (11) below.

The estimators $(\hat{h}_1^{-1}, \hat{h}_2^{-1}, \hat{\beta})$ are then the solutions to (22) and (23), where the operators are replaced by their estimators defined above. Below, we describe how to implement the method and give an explicit expression for the estimators $(\hat{h}_1^{-1}, \hat{h}_2^{-1}, \hat{\beta})$.

Given $(\hat{h}_1^{-1}, \hat{h}_2^{-1}, \hat{\beta})$, the estimator for the APE $\delta_{k,t}$ defined in (14) is given by the sample analog of that expression, i.e.

$$\hat{\delta}_{k,t} = \frac{1}{n} \sum_{i=1}^n \left[\hat{h}_t \left(\hat{h}_t^{-1}(Y_{it}) + \hat{\beta}_k \right) - Y_{it} \right], \quad t = 1, 2. \quad (24)$$

4.1 Implementation of the estimation method

The estimators $(\hat{h}_1^{-1}, \hat{h}_2^{-1}, \hat{\beta})$ are constructed as follows.

Let A_{y_t} , $t = 1, 2$, and A_z be matrices with the (i, j) element given by:

$$A_{y_t}(i, j) = \frac{\mathcal{C}\left(\frac{Y_{ti} - Y_{tj}}{b_{y_t}}\right)}{\sum_{j=1}^n \mathcal{C}\left(\frac{Y_{ti} - Y_{tj}}{b_{y_t}}\right)}, \quad t = 1, 2,$$

$$A_z(i, j) = \frac{\mathcal{K}\left(\frac{Z_i - Z_j}{b_z}\right)}{\sum_{j=1}^n \mathcal{K}\left(\frac{Z_i - Z_j}{b_z}\right)},$$

for $i = 1, \dots, n$, where \mathcal{C} is the Gaussian kernel and b_{y_t} and b_z are bandwidths. We set them equal to $n^{-1/5}$ times the standard deviations of Y_t and Z (the “rule of thumb”), respectively.

Letting \mathbb{P}_n be the $n \times n$ matrix with $\frac{n-1}{n}$ on the diagonal and $-\frac{1}{n}$ elsewhere used to impose Assumption 2 by projecting onto the space of functions of Y_1 where the mean is 0, (22) can be written as:

$$\begin{pmatrix} \gamma_n h_2^{-1} + A_{y_2} (I - \hat{P}_x) A_z h_2^{-1} - A_{y_2} (I - \hat{P}_x) A_z h_1^{-1} \\ -\gamma_n h_1^{-1} + \mathbb{P}_n A_{y_1} (I - \hat{P}_x) A_z h_2^{-1} - \mathbb{P}_n A_{y_1} (I - \hat{P}_x) A_z h_1^{-1} \end{pmatrix} = \hat{R}, \quad (25)$$

where

$$\hat{R} \equiv \begin{pmatrix} A_{y_2} (I - \hat{P}_x) A_z \Delta X_0 \\ \mathbb{P}_n A_{y_1} (I - \hat{P}_x) A_z \Delta X_0 \end{pmatrix},$$

and

$$\hat{P}_x = A_z \Delta X \left(\frac{\Delta X'}{n} A_z \Delta X \right)^{-1} \frac{\Delta X'}{n}.$$

Then the estimators $(\hat{h}_2^{-1}, \hat{h}_1^{-1})$ are given by:

$$\begin{pmatrix} \hat{h}_2^{-1} \\ \hat{h}_1^{-1} \end{pmatrix} = \begin{pmatrix} \gamma_n I + A_{y_2} (I - \hat{P}_x) A_z & -A_{y_2} (I - \hat{P}_x) A_z \\ \mathbb{P}_n A_{y_1} (I - \hat{P}_x) A_z & -(\gamma_n I + \mathbb{P}_n A_{y_1} (I - \hat{P}_x) A_z) \end{pmatrix}^{-1} \hat{R}. \quad (26)$$

and, given $(\hat{h}_2^{-1}, \hat{h}_1^{-1})$, $\hat{\beta}$ is given by:

$$\hat{\beta} = (\hat{K}_x^* \hat{K}_x)^{-1} \hat{K}_x^* [\hat{K}_y (\hat{h}_1^{-1}, \hat{h}_2^{-1}) - \hat{r}].$$

We suggest choosing the regularization parameter γ_n that minimizes the squared norm of residuals, following Florens and Sokullu (2017).

5 Asymptotic Properties

In this section, we derive assumptions for the \sqrt{n} -asymptotic normality of $\hat{\beta}$ and for the rate of convergence of $(\hat{h}_1^{-1}, \hat{h}_2^{-1})$.

In this section, a subscript of 0 will denote the true value of the parameter being

estimated.

Assumption 7. *The operator K_y is compact.*

This assumption allows us to use singular value decomposition (SVD) of the operator K_y .

Definition 1. Let $\mathcal{T} : \mathcal{E} \mapsto \mathcal{F}$ be a compact operator and let $\{\lambda_j, \phi_j, \psi_j\}$ be the singular system \mathcal{T} such that:

$$\mathcal{T}\phi_j = \lambda_j\psi_j \quad \text{and} \quad \mathcal{T}^*\psi_j = \lambda_j\phi_j,$$

where λ_j denotes the sequence of the nonzero singular values of the compact linear operator \mathcal{T} , and ϕ_j and ψ_j , for all $j \in \mathbb{N}$, are orthonormal sequences of functions in \mathcal{E} and \mathcal{F} , respectively. The *singular value decomposition* for each function $\varphi \in \mathcal{E}$ can be written as:

$$\mathcal{T}\varphi = \sum_{j=1}^{\infty} \lambda_j \langle \varphi, \phi_j \rangle \psi_j.$$

Given the definition above let $\{\lambda_j, \phi_j, \psi_j\}$ for $j \geq 1$ be the singular system of the operator K_y and let $\{\mu_l, e_l, \tilde{\psi}_l\}$ for $l = 1, 2, \dots, k$ be the singular system of the operator K_x , such that for each $\beta \in \mathbb{R}^k$ we can write:

$$K_x\beta = \sum_{l=1}^k \mu_l \langle \beta, e_l \rangle \tilde{\psi}_l.$$

Assumption 8. *Source Condition: There exists $\nu > 0$ and $\eta > 0$ such that:*

$$\sum_{j=1}^{\infty} \frac{\langle (h_1^{-1}, h_2^{-1}), \phi_j \rangle^2}{\lambda_j^{2\nu}} = \sum_{j=1}^{\infty} \frac{(\langle h_1^{-1}, \phi_{1,j} \rangle + \langle h_2^{-1}, \phi_{2,j} \rangle)^2}{\lambda_j^{2\nu}} < \infty,$$

and

$$\max_{l=1, \dots, k} \sum_{j=1}^{\infty} \frac{\langle \tilde{\psi}_l, \psi_j \rangle^2}{\lambda_j^{2\eta}} < \infty.$$

Assumption 8 is a common assumption in the NPIV literature and it defines a regularity space for the parameters of interest. The first equation in Assumption 8 defines a regularity space for (h_1^{-1}, h_2^{-1}) , in other words, this assumption adds a smoothness condition on the unknown functions. The second equation in Assumption 8 is about collinearity between Y_1, Y_2 and ΔX . As it is pointed out in Florens et al. (2012), η can be interpreted as a degree of collinearity between Y_1, Y_2 and ΔX measured through a projection on the instruments Z . For instance, when $\eta = \infty$, $\mathcal{R}(K_y)$ and $\mathcal{R}(K_x)$ are orthogonal to each other and the estimation of β is not affected by the existence of the nonparametric component as $K_y^* K_x$ and $K_x^* K_y$ vanish from the normal equations (16) and (17).

Assumption 9. *The parameters ν, η in Assumption 8 satisfy $\nu \leq 2$ and $\eta \leq 2$.*

Assumption 9 is for the sake of exposition and it is without loss of generality. In this paper, we solve the ill-posed inverse problem we encounter during estimation using Tikhonov regularization. Since Tikhonov regularization has a qualification of two, we cannot improve upon the rate of convergence when the functions we consider have regularity greater than 2, i.e., $\nu, \eta > 2$. Hence, under this assumption during the derivation of the rates, we can simply write ν or η instead of $\min\{\nu, 2\}$ or $\min\{\eta, 2\}$.

Assumption 10. *Let s be the minimum between the order of the kernel used in estimation and the order of the differentiability of densities $f(Y_1, Y_2, Z)$, $f(\Delta X, Z)$ and $f(\Delta X_0, Z)$ and assume that $s \geq 2$ and*

$$\begin{aligned} \|\hat{K}_y - K_y\|^2 &= O_p\left(\frac{1}{nb_n^{q+1}} + b_n^{2s}\right), \\ \|\hat{K}_y^* - K_y^*\|^2 &= O_p\left(\frac{1}{nb_n^{q+1}} + b_n^{2s}\right), \\ \|\hat{K}_y^* \hat{r} - \hat{K}_y^* \hat{K}_y(h_{1,0}^{-1}, h_{2,0}^{-1})\|^2 &= O_p\left(\frac{1}{n} + b_n^{2s}\right), \\ \|\hat{r} - r_0\|^2 &= O_p\left(\frac{1}{nb_n^q} + b_n^{2s}\right), \end{aligned}$$

where q is the dimension of the instrument vector Z and $b_{y_1} = b_{y_2} = b_z = b_n$ is the bandwidth.

Assumption 10 is a high-level assumption on the convergence rate of the estimated operators. Preliminary conditions leading to these rates have been studied in Darolles et al. (2011). Note that we set the bandwidths to be equal for exposition reasons. Below we state the rates we need for the smoothing parameters to converge to zero to obtain our final result.

Assumption 11. $\lim_{n \rightarrow \infty} \gamma_n \rightarrow 0$, $\lim_{n \rightarrow \infty} b_n^{2s} \rightarrow 0$, $\lim_{n \rightarrow \infty} n b_n^{q+1} \rightarrow \infty$, $\lim_{n \rightarrow \infty} n \gamma_n \rightarrow 0$, $\lim_{n \rightarrow \infty} n \gamma_n b_n^{2s} \rightarrow 0$, $\lim_{N \rightarrow \infty} \frac{\gamma_n}{b_n^{q+1}} \rightarrow 0$.

Assumption 12. $\mathcal{R}(K_y)^\perp = \mathcal{N}(K_y^*) \neq \{0\}$.

Assumption 12 implies that there exists an element ψ_j defined by the SVD of K_y such that $\psi_j \in \mathcal{R}(K_y)^\perp$. For example, this condition is satisfied in the joint nondegenerate normal case, i.e, if $(Y_1, Y_2, \Delta X, Z)$ is jointly distributed as a nondegenerate normal distribution. In such a case, the null space of K_y^* is $\{0\}$ if the range of the covariance with $(Y_1, Y_2, \Delta X)$ and Z is equal to the dimension of Z .

Assumption 13. For $\theta > 0$, we have: $E \left[|U_2 - U_1|^{2+\theta} |Z \right] = c$, for any $c \in \mathbb{R}$, and $E \left[|(I - P_y) K_x|^{2+\theta} \right] < \infty$.

Assumption 13 gives the conditions needed to satisfy the Liapounoff condition to apply the Liapounoff central limit theorem to obtain asymptotic normality of our estimators.

Using equation (23) we can show that:

$$\sqrt{n} \left(\hat{\beta} - \beta_0 \right) = \hat{M}_\gamma^{-1} \left\{ \sqrt{n} \left[K_x^* (I - P_y) \hat{E} (U_2 - U_1 | Z) \right] + O_p(1) \right\},$$

where

$$\begin{aligned} \hat{M}_\gamma &\equiv \hat{K}_x^* \hat{K}_y \left(\gamma_n I + \hat{K}_y^* \hat{K}_y \right)^{-1} \hat{K}_y^* \hat{K}_x - \hat{K}_x^* \hat{K}_x, \\ \hat{E} (U_2 - U_1 | Z) &\equiv r - \hat{K}_y (h_1^{-1}, h_2^{-1}) + \hat{K}_x \beta. \end{aligned}$$

This decomposition is useful for the following result.

Theorem 2. *Assume that $\text{Var}(U_2 - U_1 | Z) = \sigma^2$. Moreover let Assumptions 8, 9, 10, 11, 12 and 13 hold. Then:*

$$\left\| \left(\hat{h}_1^{-1}, \hat{h}_2^{-1} \right)' - \left(h_{1,0}^{-1}, h_{2,0}^{-1} \right)' \right\|_{L^2}^2 = O_p \left(\frac{1}{\gamma_n^2} \left(\frac{1}{n} + b_n^{2s} \right) + \frac{1}{\gamma_n^2} \left(\frac{1}{nb_n^{q+1}} + b_n^{2s} \right) \gamma_n^\nu + \gamma_n^\nu \right),$$

and

$$\sqrt{n} \left(\hat{\beta} - \beta_0 \right) \rightarrow \mathcal{N}(0, V),$$

where

$$V \equiv \sigma^2 M^{-1} \left[\sum_j E(\Delta X \psi_j) E(\Delta X \psi_j)' \right] M^{-1}, \psi \in \mathcal{R}(K_y)^\perp,$$

$$M \equiv K_x^* K_y (K_y^* K_y)^{-1} K_y^* K_x - K_x^* K_x.$$

Proof. The proof can be found in Appendix B. □

Theorem 2 shows that a \sqrt{n} -convergence rate and asymptotic normality for $\hat{\beta}$ can be obtained, as well as showing the convergence rate of $\left(\hat{h}_1^{-1}, \hat{h}_2^{-1} \right)$. Note that the estimator for h_t^{-1} is not necessarily monotone in its argument. We can make the estimator monotone by rearrangement. The weak convergence result obtained remains valid for the estimator obtained by rearrangement since the rearrangement operator is Hadamard differentiable, see Chernozhukov et al. (2010).

Corollary 1. *Let Assumptions 8 to 11 hold, and assume that $s \geq 2(q+1)$ and $\gamma_n \sim n^{-\frac{3}{8}}$. Then*

$$\left\| \left(\hat{h}_1^{-1}, \hat{h}_2^{-1} \right)' - \left(h_{1,0}^{-1}, h_{2,0}^{-1} \right)' \right\|_{L^2}^2 = O_p(n^{-1/4}).$$

Proof. The proof can be found in the Appendix. □

Consider now the limiting distribution of the estimator of the APE defined in (24) above. The APE is characterized by the moment condition

$$E \left[h_t \left(h_t^{-1} (Y_{it}) - \beta_k \right) - Y_{it} - \delta_{k,t} \right] = 0.$$

Then, given a random sample $\{Y_{it}\}_{i=1}^n$ and estimators $\hat{\beta}, \hat{h}_t$, the APE $\delta_{k,t}$ can be estimated by the zero of the estimating equation below:

$$\frac{1}{n} \sum_{i=1}^n \left(\hat{h}_t \left(\hat{h}_t^{-1} (Y_{it}) - \hat{\beta}_k \right) - Y_{it} - \delta_{k,t} \right) = 0.$$

This shows that $\hat{\delta}_{k,t}$ is a plug-in two-step Z-estimator. That this estimator can be shown to be \sqrt{n} -asymptotically normal should be no surprise given the regularity conditions on h_t , the way that $\delta_{k,t}$ enters the estimating equation, and the rate results on $\hat{\beta}$ and \hat{h}_t^{-1} .

A general result on two-step Z-estimators can be found in Chen et al. (2003). In that paper, Theorems 1 and 2 state sufficient high-level conditions under which $\hat{\delta}_{k,t}$ can be shown to be consistent and \sqrt{n} -asymptotically normal.¹¹ Here we make high-level assumptions as in Theorem 2 in Chen et al. (2003), in order to state our result on the \sqrt{n} -asymptotic normality of $\hat{\delta}_{k,t}$. Our simulation studies in Section 6 provide suggestive evidence that $\hat{\delta}_{k,t}$ is \sqrt{n} -asymptotically normal, e.g. Figures 3 and 4.

Using the notation in Chen et al. (2003) and assuming that h_t is differentiable

¹¹See also van der Vaart and Wellner (1996).

on its support, define the following objects:

$$\begin{aligned}
M(\delta_{k,t}, h_t, \beta_k) &\equiv E[m(\delta_{k,t}, h_t, \beta_k)] \\
&\equiv E[h_t(h_t^{-1}(Y_t) - \beta_k) - Y_t - \delta_{k,t}], \\
M_n(\delta_{k,t}, h_t, \beta_k) &\equiv \frac{1}{n} \sum_{i=1}^n (h_t(h_t^{-1}(Y_{it}) - \beta_k) - Y_{it} - \delta_{k,t}), \\
\Gamma_1(\delta_{k,t}, h_t, \beta_k) &\equiv \frac{\partial}{\partial \delta_{k,t}} M(\delta_{k,t}, h_t, \beta_k) = -1, \\
\Gamma_2(\delta_{k,t}, h_t, \beta_k) [\bar{h}_t - h_t] &\equiv \frac{d}{d\gamma} M(\delta_{k,t}, h_t + \gamma(\bar{h}_t - h_t), \beta_k) \Big|_{\gamma=0} \\
&= E \left[\left(1 - \frac{h'_t(h_t^{-1}(Y_t) - \beta_k)}{h_t^2(Y_t)} \right) [\bar{h}_t(Y_t) - h_t(Y_t)] \right], \\
\Gamma_3(\delta_{k,t}, h_t, \beta_k) &\equiv \frac{\partial}{\partial \beta_k} M(\delta_{k,t}, h_t, \beta_k) \\
&= -E[h'_t(h_t^{-1}(Y_t) - \beta_k)],
\end{aligned}$$

where h'_t is the first derivative of h_t with respect to its argument.

Theorem 3. *Let the assumptions of Corollary 1 hold, and assume that (i) h_t is continuously differentiable on its support, and Lipschitz continuous with a uniformly bounded derivative for $t = 1, 2$; (ii) the density of Y_t is bounded away from zero and is bounded from above for $t = 1, 2$; (iii) for $t = 1, 2$,*

$$\begin{aligned}
&\|M(\delta_{k,t}, h_t, \beta_k) - M(\delta_{k,t}, h_{t0}, \beta_{k0}) - \Gamma_2(\delta_{k,t}, h_{t0}, \beta_{k0})[h_t - h_{t0}] - \Gamma_3(\delta_{k,t}, h_{t0}, \beta_{k0})\| \\
&\leq c(\|h_t - h_{t0}\|_{L^2}^2 + \|\beta_k - \beta_{k0}\|^2);
\end{aligned}$$

(iv) for $t = 1, 2$, and some finite matrix V_1 ,

$$\sqrt{n} \left(M_n(\delta_{k,t,0}, h_{t0}, \beta_{k0}) + \Gamma_2(\delta_{k,t,0}, h_{t0}, \beta_{k0}) [\hat{h}_t - h_{t0}] + \Gamma_3(\delta_{k,t,0}, h_{t0}, \beta_{k0}) \right) \rightarrow \mathcal{N}(0, V_1).$$

Then $t = 1, 2$,

$$\sqrt{n} (\hat{\delta}_{k,t} - \delta_{k,t}) \rightarrow \mathcal{N}(0, V_1).$$

Proof. The proof can be found in the Appendix. □

6 Simulation study

In this section we illustrate the small sample performance of our proposed estimator through Monte Carlo simulations. We consider the case of $T = 2$, and let

$$\begin{aligned}(Z_1, Z_2) &\sim \mathcal{N} \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \right), \\ \xi &\sim \mathcal{U}[0, 1], \\ (\omega_1, \omega_2) &\sim \mathcal{N} \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma_\omega^2 & 0 \\ 0 & \sigma_\omega^2 \end{pmatrix} \right), \quad \sigma_\omega^2 = 0.5, \\ (U_1, U_2) &\sim \mathcal{N} \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma_u^2 & 0 \\ 0 & \sigma_u^2 \end{pmatrix} \right), \quad \sigma_u^2 = 0.6,\end{aligned}$$

so that

$$\begin{aligned}X_{01} &= 0.7Z_1 + 0.5U_1 + \xi, \\ X_{02} &= 0.8Z_2 + 0.4U_2 + \xi + 20, \\ X_1 &= 0.8Z_1 + 0.7Z_2 + \omega_1 + U_1, \\ X_2 &= 0.7Z_1 + 0.8Z_2 + \omega_2 + U_2.\end{aligned}$$

Additionally, let

$$\alpha \sim \mathcal{N}(0, 1) + \frac{1}{2}(X_1 + X_2),$$

and $h_1(s) = s$, $h_2(s) = \log(s)$, $\beta = 1$, so that

$$\begin{aligned}Y_1 &= \alpha + X_{01} + \beta X_1 + U_1, \\ Y_2 &= \log(\alpha + X_{02} + \beta X_2 + U_2).\end{aligned}$$

We simulate the model 500 times for sample sizes $n \in \{100, 200, 500, 1000\}$. We

estimate the functions h_1, h_2 and the finite dimensional parameter β following the method described in Section 4. We impose monotonicity of the infinite dimensional parameters by rearrangement. We choose the regularization parameter that minimizes the squared norm of residuals, following Florens and Sokullu (2017). Figures 1 and 2 show the estimated functions \hat{h}_1^{-1} and \hat{h}_2^{-1} , respectively. The light gray shaded area shows the estimated curves obtained at each draw plotted pointwise, dark gray dots show the pointwise average across simulations of the estimated functions, i.e. $\frac{1}{500} \sum_{s=1}^{500} \hat{h}_{s,t}(y_t^*)$, $t = 1, 2$, whereas the black dots show the true (pointwise) function. Table 1 shows the mean and standard error of $\hat{\beta}$ for different sample sizes. As expected, both bias and standard deviation decrease with increasing sample size.

After obtaining \hat{h}_1^{-1} , \hat{h}_2^{-1} and $\hat{\beta}$, we compute $\hat{\delta}_{k,2}$ as in (24). Table 2 shows the mean, standard error, and RMSE of estimated average partial effects for different sample sizes as well as the true average partial effect at $t = 2$ which is calculated using true values of h_1, h_2 , and β . We show two different figures, one for $n = 500$ (Figure 3) and one for $n = 1000$ (Figure 4), which provide suggestive evidence that our estimator of the APE attains \sqrt{n} -asymptotic normality.

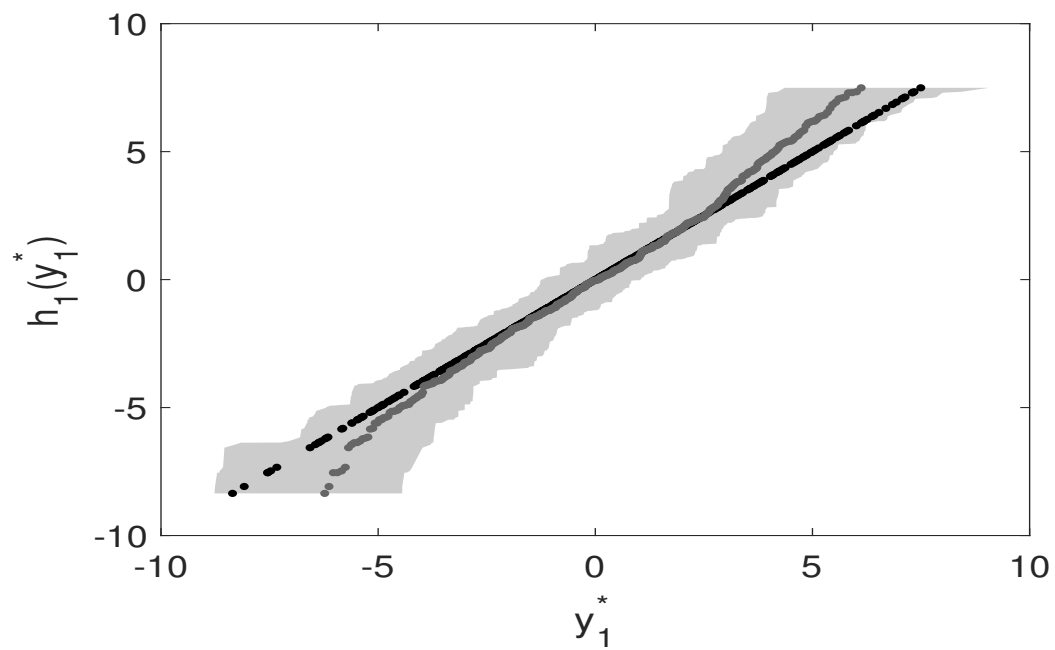


Figure 1: **Simulation result with 500 draws for h_1 , monotonicity imposed by rearrangement, $n = 500$.**

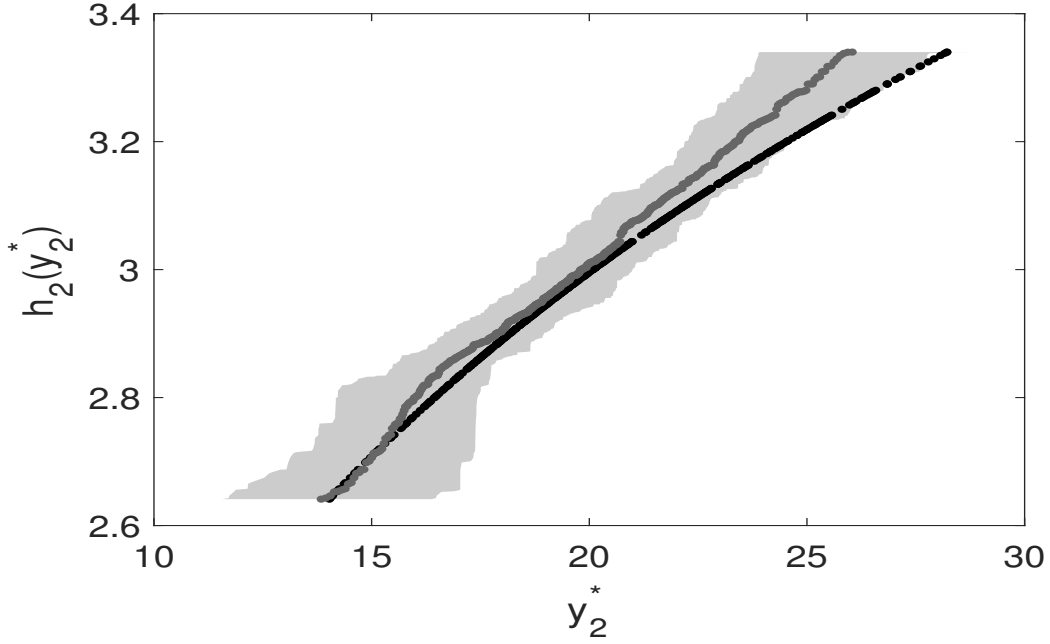


Figure 2: **Simulation result with 500 draws for h_2 , monotonicity imposed by rearrangement, $n = 500$.**

Table 1: Estimation results for β

| | Mean | Std. Err |
|------------|--------|----------|
| $n = 100$ | 0.8614 | 0.2767 |
| $n = 200$ | 0.9696 | 0.2145 |
| $n = 500$ | 1.0363 | 0.1736 |
| $n = 1000$ | 1.0583 | 0.1326 |

Table 2: Estimation results for APE

| | Mean | Std. Err | RMSE | True APE |
|------------|--------|----------|--------|----------|
| $n = 100$ | 0.0589 | 0.0165 | 0.0186 | 0.0505 |
| $n = 200$ | 0.0612 | 0.0113 | 0.0154 | 0.0506 |
| $n = 500$ | 0.0607 | 0.0084 | 0.0131 | 0.0506 |
| $n = 1000$ | 0.0597 | 0.0062 | 0.0109 | 0.0506 |

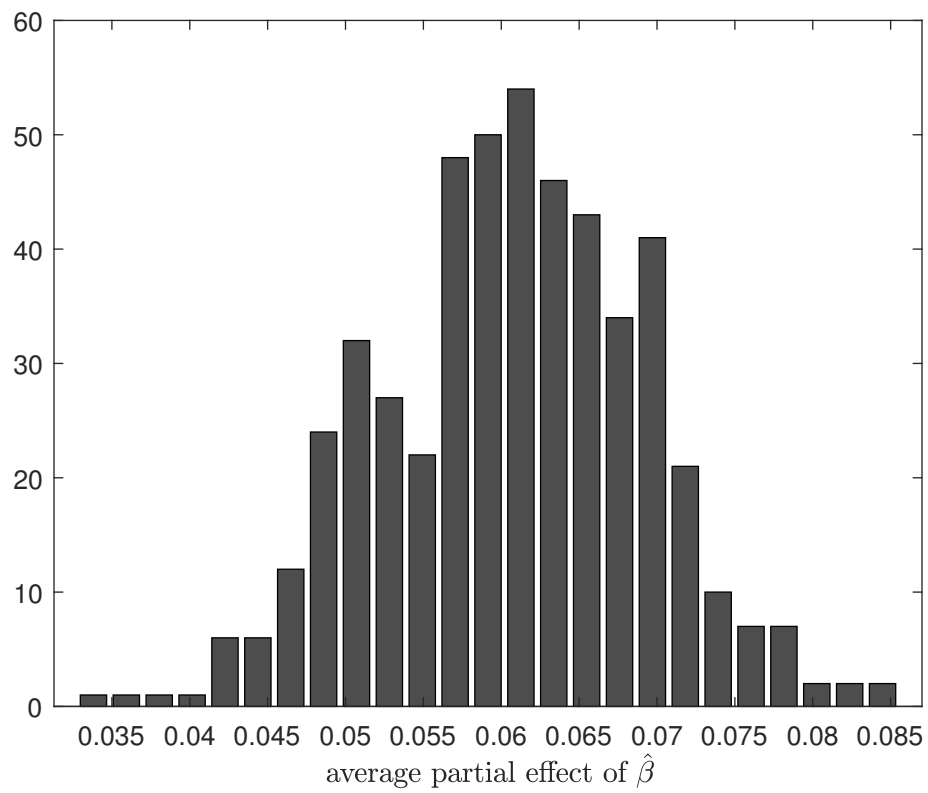


Figure 3: **Histogram of $A\hat{P}E$ for $n=500$.**

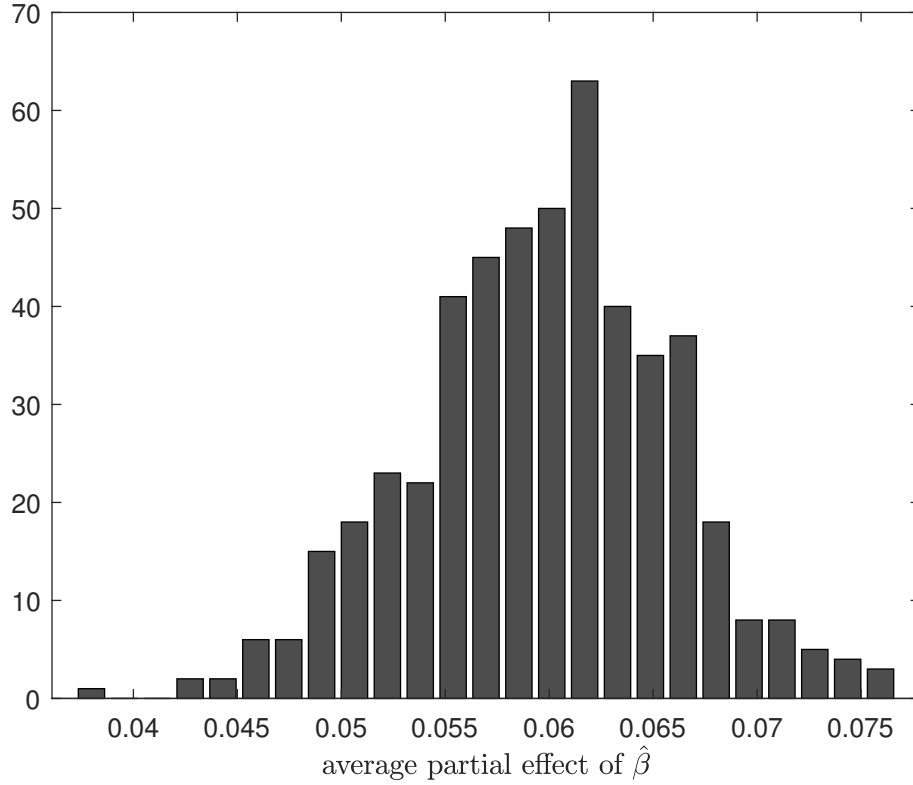


Figure 4: **Histogram of \hat{APE} for $n=1000$.**

7 Empirical Illustration

In this section, we analyze the effect of teaching practices on student achievement as measured by test scores on standardized tests in mathematics and science. Because test scores are relative ranks, any monotonic transformation of a test score is a valid score. Hence, our method is well suited to this application because it is invariant to monotonic transformations of the outcome variable. This allows us to avoid both arbitrary normalizations of test scores (see, e.g., Bonhomme and Sauder (2011)) and anchoring the scale of test scores to a measure with a well-defined cardinal scale (see,

e.g., Cunha and Heckman (2008)). In addition, our method allows for endogeneity in the factors that generate a student’s test scores beyond student-specific fixed effects, such as, for example, a measure of teaching practices based on the student’s own answer.

We use data from the Trends in International Mathematics and Science Study (TIMSS), which is an international assessment of mathematics and science knowledge of fourth and eight-grade students. Students in selected classes are administered standardized tests in mathematics and science, and background information is obtained from students and their teachers in both subjects via questionnaires. We use the 2007 wave of TIMSS for the US, that was used and described in detail in Bietenbeck (2014) and in Freyberger (2018). Our data set contains test scores of 6057 students in the eight grade on the two subjects, so that each student is observed twice: once in mathematics and once in science. Information on teaching practices comes from a questionnaire asking students how often they engaged in a range of classroom activities in each subject. Activities are classified as either traditional or modern, with the former relying on rote learning and individual work, and the latter relying on teamwork and involvement of students in discussions and presentations. We use the classification in Bietenbeck (2014), so that measurements of teaching practices are class-averages of the frequency (or percentage of lessons) of traditional or modern classroom activities. As Bietenbeck (2014) explains, these class-level indices do not add up to 100%, because teachers that use a variety of both traditional and modern teaching practices in all their lessons can score high on both indices.

Our empirical specification models student i ’s test score in subject $t \in \{\text{math, science}\}$ as the output of a production function that takes as inputs student-specific covariates:

$$Y_{it} = h_t \left(\alpha_i + \bar{R}_{it} + \bar{M}_{it}\beta + U_{it} \right), \quad (27)$$

where Y_{it} is the overall (raw) score of student i in subject t , h_t is an unknown monotonic function specific to subject t , α_i is a student fixed effect, e.g., a student’s initial endowment such as her cognitive ability, and \bar{R}_{it} and \bar{M}_{it} are class-level indices of, respectively, traditional and modern teaching practices in subject t as reported

by both student i and her classmates, and U_{it} are shocks to educational attainment of student i in subject t , that could reflect luck on an exam or an improvement or worsening in academic achievement in a particular subject relative to the long-run performance in that subject, see, e.g. Bonhomme and Sauder (2011).

Since class-averages \bar{R}_{it} and \bar{M}_{it} contain student i 's own response, there may be simultaneity issues. Following Bietenbeck (2014), we use class averages without student i 's response, $\bar{R}_{(-i)t}$ and $\bar{M}_{(-i)t}$, as instrumental variables,¹² i.e.

$$E(U_{i,math} - U_{i,science} | \bar{R}_{(-i)math}, \bar{M}_{(-i)math}, \bar{R}_{(-i)science}, \bar{M}_{(-i)science}) = 0, \quad (28)$$

which is Assumption 3 stated in the context of our application. Note that our assumptions in Section 2 allow student i 's shocks to educational attainment to be correlated across mathematics and science, so that if a student shows an improvement in academic achievement in mathematics in one year relative to her long-run academic performance, then she may show an academic improvement in science as well.

We are interested in estimating the following APEs: (i) the effect on mathematics and science test scores of increasing the traditional teaching index by 1 (from 0% to 100%, for example) while holding the modern teaching index unchanged, and (ii) the effect on test scores of increasing the modern teaching index by 1 while holding the traditional index unchanged. These partial effects correspond to counterfactuals associated to an increase in, respectively, traditional and modern teaching practices at the expense of practices that are neither traditional nor modern, such as reviewing an exam or homework. Using expression (13) in Section 3, we compute the following APEs:

$$\delta_{R,t} = E(Y_{it}(\bar{R}_{it} + 1)) - E(Y_{it}), \quad (29)$$

$$\delta_{M,t} = E(Y_{it}(\bar{M}_{it} + \beta)) - E(Y_{it}). \quad (30)$$

Before estimating h_t and β in our preferred outcome equation (27), which are

¹²In the linear equivalent to (27) (i.e. $h_t(x) = x$), Bietenbeck (2014) uses $\bar{R}_{(-i)t}$ and $\bar{M}_{(-i)t}$ as exogenous regressors.

needed to compute the APEs above, we run a few linear panel regression specifications in order to (i) compare our results to those in Bietenbeck (2014), (ii) establish that using \bar{R}_{it} and \bar{M}_{it} as covariates and $\bar{R}_{(-i)t}$ and $\bar{M}_{(-i)t}$ as instrumental variables replicates the effects of using $\bar{R}_{(-i)t}$ and $\bar{M}_{(-i)t}$ as covariates (as in Bietenbeck, 2014), and (iii) motivate arbitrary subject-specific monotonic transformations of the test scores, by showing that there are subject-specific effects and that using standardized test scores versus raw scores obtains different results.

First, we run the following linear panel data regressions where the outcomes \tilde{Y}_{it} are the standardized scores and λ_t are subject-specific effects. Specification (31) below is that of Table 3 column 3 in Bietenbeck (2014):

$$\tilde{Y}_{it} = \alpha_i + \lambda_t + \bar{R}_{(-i)t} + \bar{M}_{(-i)t}\beta + U_{it}. \quad (31)$$

Specification (32) uses \bar{R}_{it} and \bar{M}_{it} as covariates to establish that there is an endogeneity problem as explained in Bietenbeck (2014):

$$\tilde{Y}_{it} = \alpha_i + \lambda_t + \bar{R}_{it} + \bar{M}_{it}\beta + U_{it}, \quad (32)$$

while specification (33) below corrects the endogeneity problem by using as $\bar{R}_{(-i)t}$ and $\bar{M}_{(-i)t}$ as instrumental variables in:

$$\tilde{Y}_{it} = \alpha_i + \lambda_t + \bar{R}_{it} + \bar{M}_{it}\beta + U_{it}, \text{ and (28) holds.} \quad (33)$$

We then run the same regressions with the raw scores, Y_{it} .

The results of these six regressions can be found in Table 3. The results for the specifications using the standardized scores show that our specification in (31) reproduces those in Table 3 column 3 in Bietenbeck (2014);¹³ our specification in (32) suffers from an endogeneity problem, which is then corrected by the specification in (33). The specification in (33) recovers the results associated with the specification in

¹³We do not exactly reproduce the results since we do not control for teacher-specific covariates and we do not use student sampling weights as in Bietenbeck (2014).

(31), giving peace of mind about the validity of the instrumental variables.¹⁴ When we repeat the exercise with the raw scores Y_{it} , we find that our specifications corresponding to (31) and (33) obtain a positive and significant effect of modern teaching practices, which is different from the results using the standardized test scores \tilde{Y}_{it} . We also ran these specifications using a wide range of Box-Cox transformations of the raw scores, and obtained similar results. Our results then suggest that the type of transformation applied to test scores matters.

| | Standardized scores | | | Raw scores | | |
|----------------------------|---------------------|------------------|------------------|-------------------|-------------------|-------------------|
| | Model 1 | Model 2 | Model 3 | Model 1 | Model 2 | Model 3 |
| Traditional teaching index | 0.312 (0.053) | 0.374 (0.054) | 0.325 (0.055) | 16.251 (4.130) | 20.810 (4.211) | 16.999 (4.279) |
| Modern teaching index | 0.047 (0.045) | 0.076 (0.046) | 0.050 (0.047) | 12.009 (3.546) | 14.500 (3.594) | 12.460 (3.640) |
| Course effects | 0.015 (0.007) | 0.018 (0.007) | 0.016 (0.007) | 11.914 (0.510) | 12.104 (0.511) | 11.942 (0.512) |
| Sample size | 6057 | 6057 | 6057 | 6057 | 6057 | 6057 |
| R^2 | 0.948 | 0.948 | 0.948 | 0.947 | 0.948 | 0.948 |
| Adjusted R^2 | 0.897 | 0.897 | 0.897 | 0.895 | 0.895 | 0.895 |

Table 3: The results for specification (31) are shown under Model 1, for specification (32) are shown under Model 2, and for specification (33) are shown under Model 3. Standard errors in parantheses.

Second, we document via linear panel regressions similar to those above, the existence of subject-specific effects and the lack of invariance of the results to different normalizations of the outcome variable. More precisely, we run the following regression:

$$\tilde{Y}_{it} = c_{0i} + c_{1t}\bar{R}_{(-i)t} + c_{2t}\bar{M}_{(-i)t} + \lambda_t + U_{it}, \quad (34)$$

where the effects of teaching practices vary across the two subjects, and where the outcome variable is the standardized test score \tilde{Y}_{it} . We then run the same specifica-

¹⁴We repeat the exercise associated to specification 31 with both differenced instruments, i.e. $\bar{R}_{(-i)t=\text{math}} - \bar{R}_{(-i)t=\text{science}}$, $\bar{M}_{(-i)t=\text{math}} - \bar{M}_{(-i)t=\text{science}}$, and instruments in levels, i.e. $\bar{R}_{(-i)t=\text{math}}$, $\bar{R}_{(-i)t=\text{science}}$, $\bar{M}_{(-i)t=\text{math}}$, $\bar{M}_{(-i)t=\text{science}}$. Our results are virtually unchanged.

tion with the raw scores Y_{it} . The results of these regressions can be found in Table 4.

| | Standardized scores | Raw scores |
|----------------------------|---------------------|----------------|
| Traditional teaching index | | |
| mathematics | 0.397 (0.063) | 16.927 (4.906) |
| science | 0.217 (0.070) | 17.148 (5.459) |
| Modern teaching index | | |
| mathematics | 0.075 (0.058) | 14.438 (4.535) |
| science | 0.071 (0.060) | 9.798 (4.653) |
| Course effects | 0.129 (0.051) | 14.265 (4.018) |
| Sample size | 6057 | 6057 |
| Adjusted R^2 | 0.897 | 0.895 |

Table 4: The results for specification (34) with the standardized test scores and with the raw test scores. Standard errors in parantheses.

We find that, when using the raw scores, the effect of teaching practices on overall test scores varies across subjects and that both traditional teaching and modern teaching have a positive and statistically significant effect. Using standardized test scores instead obtains that modern teaching practices do not have a significant effect on overall test scores. Our results then suggest that there may be heterogeneous effects of teaching practices across subjects.

Taken together, we interpret the results of Tables 3 and 4 as suggestive evidence that the effects of teaching practices are sensitive to the type of transformation applied to the test scores and that they are heterogeneous across subjects – which justifies $h_t(\cdot)$ in our outcome equation specification in (27), and that we can use class-level indices that include student i ’s response as covariates and class-level indices that exclude it as instrumental variables.

We show below our results from estimating the model in (27) and (28). We use the rule-of-thumb bandwidth parameters and a regularization parameter of 10^{-5} for both functions. We also use all four instrumental variables, $\bar{R}_{(-i)t}$, $\bar{M}_{(-i)t}$, $t \in \{m, s\}$. The estimate for β is 1.79, while the two functions h_{math} and h_{science} are shown in Figures 5 and 6, respectively.

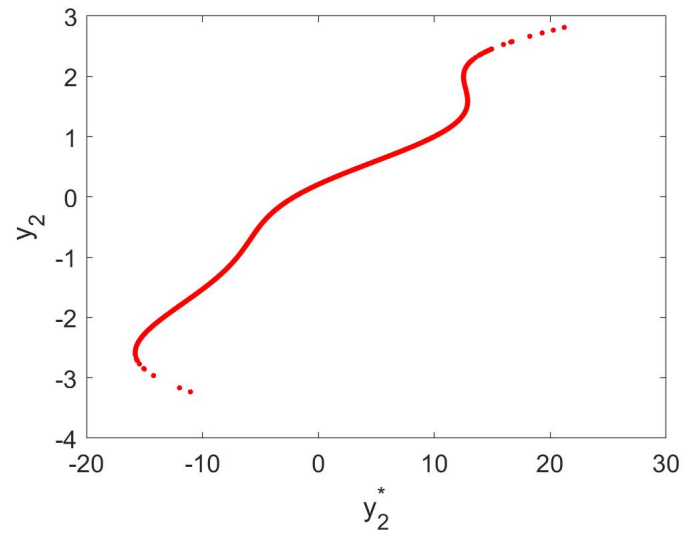


Figure 5: Estimated transformation of mathematics test scores.

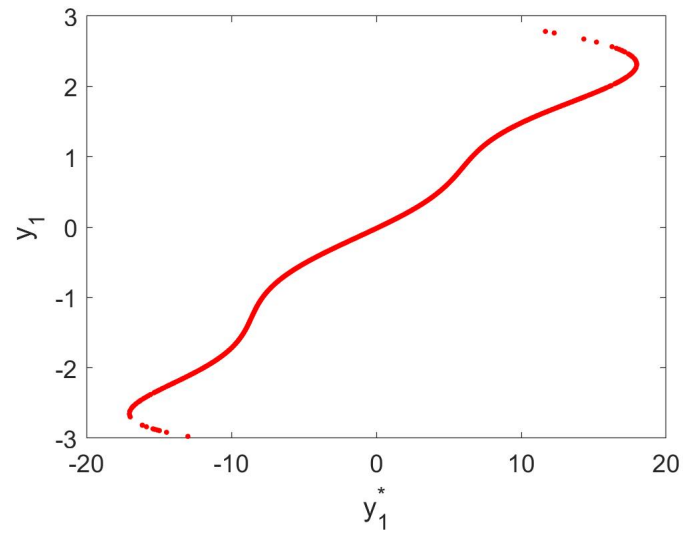


Figure 6: Estimated transformation of science test scores.

The estimated APEs for the effects of increasing traditional methods to 100% in (29) are $\hat{\delta}_{R,\text{math}} = 12.04$ and $\hat{\delta}_{R,\text{science}} = 13.34$, and those of increasing modern methods to 100% in (30) are $\hat{\delta}_{M,\text{math}} = 21.09$ and $\hat{\delta}_{M,\text{science}} = 23.95$. These results suggest that increasing traditional teaching methods has a similar effect on mathematics and science, and that the effect is much smaller than that of increasing modern teaching methods. The results remain unchanged when adjusting for the standard deviation of the test scores (74.5 points for mathematics and 79.8 for science).

References

- ABREVAYA, J. (1999): “Leapfrog estimation of a fixed-effects model with unknown transformation of the dependent variable,” *Journal of Econometrics*, 93, 203–228.
- AGUIRREGABIRIA, V. AND J. CARRO (2021): “Identification of Average Marginal Effects in Fixed Effects Dynamic Discrete Choice Models,” <https://arxiv.org/abs/2107.06141>.
- AGUIRREGABIRIA, V., J. GU, AND Y. LUO (2021): “Sufficient Statistics for Unobserved Heterogeneity in Structural Dynamic Logit Models,” *Journal of Econometrics*, 223, 280–311.
- ALTONJI, J. G. AND R. L. MATZKIN (2005): “Cross Section and Panel Data Estimators for Nonseparable Models with Endogenous Regressors,” *Econometrica*, 73, 1053–1102.
- ANDREWS, D. W. (2017): “Examples of L2-complete and boundedly-complete distributions,” *Journal of Econometrics*, 199, 213–220.
- ARELLANO, M. AND S. BOND (1991): “Some Tests of Specification for Panel Data: Monte Carlo Evidence and an Application to Employment Equations,” *The Review of Economic Studies*, 58, 277–297.
- ARELLANO, M. AND S. BONHOMME (2011): “Nonlinear Panel Data Analysis,” *Annual Review of Economics*, 3, 395–424.
- ARELLANO, M. AND B. HONORÉ (2001): “Panel Data Models: Some Recent Developments,” in *Handbook of Econometrics*, ed. by J. J. Heckman and E. Leamer, Elsevier, vol. 5, 3229–3296.
- BABII, A. AND J.-P. FLORENS (2020): “Is completeness necessary? Estimation in nonidentified linear models,” Working paper.
- BIETENBECK, J. (2014): “Teaching Practices and Cognitive Skills,” *Labour Economics*, 20, 143–153.

- BIRKE, M., S. VAN BELLEGEM, AND I. VAN KEILEGOM (2017): “Semi-parametric estimation in a single-index model with endogenous variables,” *Scandinavian Journal of Statistics*, 44, 168–191.
- BLUNDELL, R. AND S. BOND (1998): “Initial conditions and moment restrictions in dynamic panel data models,” *Journal of econometrics*, 87, 115–143.
- BONHOMME, S. AND U. SAUDER (2011): “Recovering distributions in difference-in-differences models: A comparison of selective and comprehensive schooling,” *The Review of Economics and Statistics*, 93, 479–494.
- BOTOSARU, I. AND C. MURIS (2017): “Binarization for panel models with fixed effects,” Cemmap working paper CWP31/17.
- BOTOSARU, I., C. MURIS, AND K. PENDAKUR (2021): “Identification of Time-Varying Transformation Models with Fixed Effects, with an Application to Unobserved Heterogeneity in Resource Shares,” *Journal of Econometrics*, forthcoming.
- BUN, M. AND V. SARAFIDIS (2015): “Dynamic panel data models,” *The Oxford handbook of panel data*, 76–110.
- CARRASCO, M. AND J. FLORENS (2011): “A Spectral Method for Deconvolving a Density,” *Econometric Theory*, 27, 546–581.
- CENTORRINO, S., F. FÉVE, AND J.-P. FLORENS (2017): “Additive Nonparametric Instrumental Regressions: A Guide to Implementation,” *Journal of Econometric Methods*, 6, 1–25.
- CHEN, X., V. CHERNOZHUKOV, S. LEE, AND W. NEWHEY (2014): “Local Identification of Nonparametric and Semiparametric Models,” *Econometrica*, 2, 785–809.
- CHEN, X., O. LINTON, AND I. VAN KEILEGOM (2003): “Estimation of Semiparametric Models when the Criterion Function Is Not Smooth,” *Econometrica*, 71, 1591–1608.

- CHERNOZHUKOV, V., I. FERNÁNDEZ-VAL, AND A. GALICHON (2010): “Quantiles and Probability Curves without Crossing,” *Econometrica*, 78, 1093–1125.
- CHERNOZHUKOV, V., I. FERNÁNDEZ-VAL, J. HAHN, AND W. NEWEY (2013): “Average and Quantiles Effects in Nonseparable Panel Models,” *Econometrica*, 81, 535–580.
- CHIAPPORI, P.-A., I. KOMUNJER, AND D. KRISTENSEN (2015): “Nonparametric identification and estimation of transformation models,” *Journal of Econometrics*, 188, 22–39.
- CUNHA, F. AND J. HECKMAN (2008): “Formulating, Identifying and Estimating the Technology of Cognitive and Noncognitive Skill Formation,” *Journal of Human Resources*, 43, 738–782.
- DAROLLES, S., Y. FAN, J.-P. FLORENS, AND E. RENAULT (2011): “Nonparametric instrumental regression,” *Econometrica*, 79, 1541–1565.
- DAVEZIES, L., X. D’HAULTFOEUILLE, AND L. LAAGE (2021): “Identification and Estimation of Average Marginal Effects in Fixed Effects Logit Models,” <https://arxiv.org/abs/2105.00879>.
- D’HAULTFOEUILLE, X. (2010): “A new instrumental method for dealing with endogenous selection,” *Journal of Econometrics*, 154, 1–15.
- (2011): “On the completeness condition in nonparametric instrumental problems,” *Econometric Theory*, 27, 460–471.
- DOBRONYI, C., J. GU, AND K. IL KIM (2021): “Identification of Dynamic Panel Logit Models with Fixed Effects,” <https://arxiv.org/abs/2104.04590>.
- FÉVE, F. AND J.-P. FLORENS (2010): “The practice of nonparametric estimation by solving inverse problems: the example of transformation models,” *The Econometrics Journal*, 13, S1–S27.

- (2014): “Non parametric analysis of panel data models with endogenous variables,” *Journal of Econometrics*, 181, 151–164.
- FLORENS, J.-P., J. JOHANNES, AND S. VAN BELLEGEM (2012): “Instrumental regression in partially linear models,” *The Econometrics Journal*, 15, 304–324.
- FLORENS, J.-P. AND S. SOKULLU (2017): “Nonparametric estimation of semiparametric transformation models,” *Econometric Theory*, 33, 839–873.
- FREYBERGER, J. (2018): “Non-parametric panel data models with interactive fixed effects,” *Review of Economic Studies*, 85, 1824–1851.
- HONORÉ, B. E. AND E. KYRIAZIDOU (2000): “Panel Data Discrete Choice Models with Lagged Dependent Variables,” *Econometrica*, 68, 839–874.
- HONORÉ, B. E., C. MURIS, AND M. WEIDNER (2021): “Dynamic Ordered Panel Logit Models,” *arXiv preprint arXiv:2107.03253*.
- HONORÉ, B. E. AND M. WEIDNER (2020): “Moment Conditions for Dynamic Panel Logit Models with Fixed Effects,” *arXiv:2005.05942 [econ]*.
- HOROWITZ, J. (2011): “Applied Nonparametric Instrumental Variables Estimation,” *Econometrica*, 79, 347–394.
- HOROWITZ, J. L. (2009): *Semiparametric and Nonparametric Methods in Econometrics*, Springer.
- HU, Y. AND J.-L. SHIU (2018): “Nonparametric identification using instrumental variables: sufficient conditions for completeness,” *Econometric Theory*, 34, 659–693.
- LIU, L., A. POIRIER, AND J.-L. SHIU (2021): “Identification and Estimation of Average Partial Effects in Semiparametric Binary Response Panel Models,” <https://arxiv.org/abs/2105.12891>.

- MURIS, C., P. RAPOSO, AND S. VANDOROS (2020): “A dynamic ordered logit model with fixed effects,” *arXiv preprint arXiv:2008.05517*.
- NEWHEY, W. K. AND J. L. POWELL (2003): “Instrumental variable estimation of nonparametric models,” *Econometrica*, 71, 1565–1578.
- NEWHEY, W. K., J. L. POWELL, AND F. VELLA (1999): “Nonparametric estimation of triangular simultaneous equations models,” *Econometrica*, 67, 565–603.
- PAKEL, C. AND M. WEIDNER (2021): “Bounds on Average Effects in Discrete Choice Panel Data Models,” Working paper.
- VAN DER VAART, A. AND J. WELLNER (1996): *Weak convergence and Empirical processes*, Springer.
- VANHEMS, A. AND I. VAN KEILEGOM (2019): “Estimation of a semiparametric transformation model in the presence of endogeneity,” *Econometric Theory*, 35, 73–110.

A Additional results

A.1 Measurable separability

In this section, we establish low-level assumptions for measurable separability, which is a sufficient assumption for Assumption 6(ii) in the main text.

Assumption 14. *The random variables $W \equiv (Y_1, Y_2, \Delta X)$ are such that for any $\delta_2 - \delta_1 \in L_W^2$ and any $b \in \mathbb{R}^k$, if*

$$\delta_2(Y_2) - \delta_1(Y_1) - \Delta X b = 0 \text{ a.s. } F_{Y_1, Y_2, \Delta X},$$

then there exist constants $c_t \in \mathbb{R}$, $t = 1, 2$, such that

$$\delta_t(Y_t) = c_t \text{ a.s. } F_{Y_t}, \quad t = 1, 2.$$

Assumption 14 is a high-level assumption that rules out a linear relationship between Y_1, Y_2 , and ΔX . The assumption is a slightly weaker version of the measurable separability assumption made in the NPIV literature. The assumption fails if there exists an *additive* functional relationship between Y_1, Y_2 , and ΔX , see, e.g., Newey et al. (1999).¹⁵ Identification may still occur in the presence of a non-additive functional relationship between the three random variables. The Lemma below establishes sufficient low-level assumptions for Assumption 14.

Lemma 1. *Let the following assumptions hold: (L1) h_t is continuously differentiable for all t ; (L2) the support of X_t contains an open set and is continuous on that set; (L3) U_t is continuous for all t and is serially independent. Then for any h_1, h_2, β satisfying Assumptions 1 to 5 and Assumptions L1, L2, L3, for any random variables Y_t, X_t, Z following the model above, and for any $\delta_2 - \delta_1 \in L_W^2$ and for any $b \in \mathbb{R}^k$, if*

$$\delta_2(Y_2) - \delta_1(Y_1) - \Delta X b = 0 \text{ a.s. } F_{Y_1, Y_2, \Delta X}, \quad (35)$$

then there exist constants $c_t \in \mathbb{R}$ such that $\delta_t(Y_t) = c_t$ a.s. F_{Y_t} , $t = 1, 2$.

Proof. The conclusion of the Lemma follows by contradiction. That is, assuming both (35) and $\delta_1(Y_1) \neq c_1$ a.s. or $\delta_2(Y_2) \neq c_2$ a.s. for all $c_1, c_2 \in \mathbb{R}$, leads to a contradiction.

First, solve for α from the outcome equation for Y_1 and plug the resulting expression in the outcome equation for Y_2 to obtain:

$$Y_2 = h_2(h_1^{-1}(Y_1) + (X_2 - X_1)\beta + U_2 - U_1).$$

¹⁵For example, Newey et al. (1999) write that there exists a functional relationship between two random variables W_1 and W_2 provided that there exist functions $H(W_1, W_2)$ and a set \mathcal{W} such that $P(\mathcal{W}) > 0$ and

$$\begin{aligned} P(H(W_1, W_2) = 0) &= 1 \\ P(H(W_1, \bar{W}_2) = 0) &< 1 \end{aligned}$$

for all fixed $\bar{W}_2 \in \mathcal{W}$. In fact, Assumption 14 is implied by two measurable separability assumptions: one between (Y_1, Y_2) and ΔX , and another between Y_1 and Y_2 .

Consider then (35):

$$\delta_2 \left(h_2 \left(h_1^{-1}(y_1) + (x_2 - x_1) \beta - u_1 + u_2 \right) \right) - \delta_1(y_1) \equiv (x_2 - x_1) b, \quad (36)$$

for all $x_t \in \mathcal{X}_t, y_t \in \mathcal{Y}_t, u_t \in \mathcal{U}_t, t = 1, 2$.

First, note that since X_2 and U_2 are correlated, we can think of X_2 as a function of U_2 , e.g., $X_2 = \gamma(U_2) + \eta_2$, $\eta_2 = X_2 - \gamma(U_2)$. Second, note that h_t being differentiable guarantees that δ_t is also differentiable. Then differentiating (36) with respect to u_2 obtains

$$\frac{\partial \delta_2}{\partial h_2} \left(\frac{\partial h_2}{\partial x_2} \frac{\partial \gamma_2}{\partial u_2} \beta + \frac{\partial h_2}{\partial u_2} \right) = \frac{\partial \gamma_2}{\partial u_2} b, \quad (37)$$

where we used Assumptions L1, L2, and L3, and that X_2 is correlated with U_2 .

However, since $\delta_2(Y_2) \neq c_2$ a.s. it follows that

$$\frac{\partial \delta_2}{\partial h_2} \left(\frac{\partial h_2}{\partial x_2} \frac{\partial \gamma_2}{\partial u_2} \beta + \frac{\partial h_2}{\partial u_2} \right) \neq 0. \quad (38)$$

Combining (37) and (38), it must be that for all $b \in \mathbb{R}^k$,

$$\frac{\partial \gamma_2}{\partial u_2} b \neq 0.$$

Since X_2 is correlated with U_2 , $\frac{\partial \gamma_2}{\partial u_2} \neq 0$. Hence it follows that $b \neq 0$, which is not true since $b \in \mathbb{R}^k$.

Similarly, we can show that assuming (35) and $\delta_1(Y_1) \neq c_1$ for all $c_1 \in \mathbb{R}$ leads to a contradiction. \square

A.2 Normal equations

In this section we show that the system of three normal equations based on the three linear operators that we use for identification is identical to the system of two equations based on a bilinear operator that we use for estimation.

The normal equations using the three operators K_x , K_{y_1} , and K_{y_2} are:

$$K_{y_1}^* r = K_{y_1}^* K_{y_2} h_2^{-1} - K_{y_1}^* K_{y_1} h_1^{-1} - K_{y_1}^* K_x \beta, \quad (39)$$

$$K_{y_2}^* r = K_{y_2}^* K_{y_2} h_2^{-1} - K_{y_2}^* K_{y_1} h_1^{-1} - K_{y_2}^* K_x \beta, \quad (40)$$

$$K_x^* r = K_x^* K_{y_2} h_2^{-1} - K_x^* K_{y_1} h_1^{-1} - K_x^* K_x \beta. \quad (41)$$

Notice that (41) can be written as

$$K_x^* r = K_x^* (K_{y_2} h_2^{-1} - K_{y_1} h_1^{-1}) - K_x^* K_x \beta = K_x^* K_y (h_1^{-1}, h_2^{-1}) - K_x^* \beta,$$

where we used the definition of K_y . The expression above is (17) in the main text.

Consider now (39) and (40), and rewrite them as

$$\begin{aligned} K_{y_1}^* K_y (h_1^{-1}, h_2^{-1}) &= K_{y_1}^* r + K_{y_1}^* K_x \beta, \\ K_{y_2}^* K_y (h_1^{-1}, h_2^{-1}) &= K_{y_2}^* r + K_{y_2}^* K_x \beta. \end{aligned}$$

Imposing Assumption 2(ii), multiplying the second equation above by -1 , and using the definition of K_y^* , obtains equation (16) in the main text.

B Proofs

B.1 Proof of Theorem 1

Let (h_1, h_2, β) be the true value of the model parameters, and let $(g_1, g_2, B) \in L_{\mathcal{Y}_1}^2 \otimes L_{\mathcal{Y}_2}^2 \otimes \mathbb{R}^k$ be alternative values such that

$$(g_1, g_2, B) \neq (h_1, h_2, \beta)$$

and such that they satisfy the same assumptions as (h_1, h_2, β) , i.e. Assumptions 1, 2, 3, and 4. In particular, for any $z \in \mathcal{Z}$:

$$E(g_2^{-1}(Y_2) - g_1^{-1}(Y_1) - \Delta X B \mid Z = z) = E(\Delta X_0 \mid Z = z). \quad (42)$$

Equating (6) and (42), and re-arranging yields

$$E(\delta_2(Y_2) - \delta_1(Y_1) - \Delta X b | Z = z) = 0,$$

where

$$\delta_t(Y_t) \equiv h_t^{-1}(Y_t) - g_t^{-1}(Y_t), t = 1, 2, \quad (43)$$

$$b \equiv \beta - B. \quad (44)$$

Assumption 5 obtains that

$$\delta_2(Y_2) = 0, \delta_1(Y_1) = 0, \Delta X b = 0 \text{ a.s. } F_{Y_1, Y_2, \Delta X}, \quad (45)$$

which is a contradiction. Then since $h_t^{-1}, t = 1, 2$ have been identified, the pre-images of $h_t, t = 1, 2$, and, hence $h_t, t = 1, 2$, are identified.

Here, we show by contradiction that Assumptions 6(i) and 6(ii) imply Assumption 5. Suppose that Assumptions 6(i) and 6(ii) hold and that Assumption 5 does not. Let δ_1, δ_2, b be such that a.s.

$$K_{y_1} \delta_1 = 0, K_{y_2} \delta_2 = 0, \text{ and } K_x b = 0$$

so that, by injectivity of the operators, $\delta_1 = \delta_2 = b = 0$. Then

$$K(\delta_1, \delta_2, b) = K_{y_2} \delta_2 - K_{y_1} \delta_1 - K_x b = 0 \text{ a.s.}$$

and $(\delta_1, \delta_2, b) \neq (0, 0, 0)$ since Assumption 5 does not hold. This leads to a contradiction.

B.2 Proof of Theorem 2

Proof. The proof follows from Florens and Sokullu (2017). Here we provide a sketch. First, note that

$$\hat{H}^{\gamma^n} - H = A + B + C,$$

where

$$A \equiv (\gamma_n I + \hat{K}_y^*(I - \hat{P}_x)\hat{K}_y)^{-1} \hat{K}_y^*(I - \hat{P}_x)\hat{r} - (\gamma_n I + \hat{K}_y^*(I - \hat{P}_x)\hat{K}_y)^{-1} \hat{K}_y^*(I - \hat{P}_x)\hat{K}_y H, \quad (46)$$

$$B \equiv (\gamma_n I + \hat{K}_y^*(I - \hat{P}_x)\hat{K}_y)^{-1} \hat{K}_y^*(I - \hat{P}_x)\hat{K}_y H - (\gamma_n I + K_y^*(I - P_x)K_y)^{-1} K_y^*(I - P_x)K_y H, \quad (47)$$

$$C \equiv (\gamma_n I - K_y^*(I - P_x)K_y)^{-1} K_y^*(I - P_x)K_y H - H, \quad (48)$$

where A captures the estimation error on the right hand side of the equation, B shows the error coming from estimation of the operators, and C captures the regularisation bias. Following Florens and Sokullu (2017), A can be shown to be $O_p\left(\frac{1}{\gamma_n^2}\left(\frac{1}{n} + b_n^{2s}\right)\right)$, while B and C are $O_p\left(\frac{1}{\gamma_n^2}\left(\frac{1}{nb_n^{q+1}} + b_n^{2s}\right)\gamma_n^\nu\right)$ and $O_p(\gamma_n^\nu)$, respectively.

Second, $\sqrt{n}(\hat{\beta} - \beta)$ can be decomposed as:

$$\begin{aligned} \sqrt{n}(\hat{\beta} - \beta) = & \hat{M}_\gamma^{-1} \left\{ \underbrace{\sqrt{n}[K_x^*(I - P_y)\hat{E}(U_2 - U_1|Z)]}_I \right. \\ & - \underbrace{\sqrt{n}[K_x^*(I - P_y) - \hat{K}_x^*(I - \hat{P}_y^\gamma)]\hat{E}(U_2 - U_1|Z)}_{II} \\ & \left. + \underbrace{\sqrt{n}[\hat{K}_x^*(I - \hat{P}_y^\gamma)\hat{K}_y(h_1^{-1}, h_2^{-1})]}_{III} \right\} \end{aligned}$$

where $\hat{P}_y^\gamma = \hat{K}_y(\gamma I + \hat{K}_y^*\hat{K}_y)^{-1}\hat{K}_y^*$. The proof then proceeds showing the following which lead to the final result:

$$\|\hat{M}_\gamma^{-1} - M^{-1}\| \rightarrow o_p(1),$$

where

$$M = K_x^* K_y (K_y^* K_y)^{-1} K_y^* K_x - K_x^* K_x,$$

and

$$\begin{aligned}
\|II\| &\rightarrow O_p(1), \\
\|III\| &\rightarrow O_p(1), \\
\hat{M}_\gamma^{-1} \left\{ \sqrt{n} [K_x^* (I - P_y) \hat{E}(U_2 - U_1 | Z)] \right\} &\rightarrow \mathcal{N} \left(0, \sigma^2 M^{-1} \left(\sum_{j/\psi_j \in \mathcal{R}(K_y)^\perp} E(\Delta X \psi_j) E(\Delta X \psi_j)' M^{-1} \right) \right).
\end{aligned}$$

□

B.3 Proof of Corollary 1

Following Darolles et al. (2011), we first show that rate of convergence of \hat{H}^{γ_n} can be shown to be equal to $n^{-\frac{\nu}{2+\nu}}$. And then we show that $\nu = 2/3$, this rate is equal to $n^{-1/4}$. Consider the convergence rate of \hat{H}^{γ_n} given in Theorem 1. The proof based on making the middle term negligible. Assume that $b_n^{2s} \sim \frac{1}{n}$, together with assumption $nb_n^{q+1} \rightarrow \infty$, this implies that $s \geq \frac{q+1}{2}$ and then the middle term is $O_p \left(\frac{\gamma_n^{\nu-2}}{nb_n^{q+1}} \right)$.

If the middle term is negligible, together with $b_n^{2s} \sim 1/n$, optimal γ_n is obtained by setting equal the first and the third term:

$$\frac{1}{\gamma_n^2 n} \sim \gamma_n^\nu,$$

which will lead to $\gamma_n \sim n^{-\frac{1}{2+\nu}}$. Going back to the middle term, one can then choose a bandwidth which satisfies:

$$\frac{1}{nb_n^{q+1}} = O_p \left(\frac{\gamma_n^\nu}{\gamma_n^\nu - 2} \right)$$

If we replace the γ_n with its optimal rate in the above equation, we obtain the first condition of the corollary. Then under $\gamma_n \sim n^{-\frac{1}{2+\nu}}$ and if $s \geq \frac{(q+1)(\nu+2)}{2\nu}$, the rate of convergence of \hat{H}^{γ_n} is given by:

$$\|\hat{H}^{\gamma_n} - H\|^2 = O_p(n^{-\nu/\nu+2}),$$

which is equal to $O_p(n^{-1/4})$ for $\nu = 2/3$.

B.4 Proof of Theorem 3

The proof consists in verifying the conditions in Theorem 2 in Chen et al. (2003). Conditions (2.1), (2.2), (2.4) are standard and hold. Condition (2.5) holds since the conditions of Lemma 1 in Chen et al. (2003) hold, which is sufficient for Condition (2.5), see Remark 2 in Chen et al. (2003). In particular, the class of functions $\{m(\delta_{k,t}, h_t, \beta_k) : h_t \in L^2(\mathbb{R}), \beta_k \in \mathbb{R}, \delta_{k,t} \in \mathbb{R}\}$ is P -Donsker, where P is the probability measure of Y_t given that h_t is strictly increasing and Lipschitz continuous, and given Donsker preservation results in van der Vaart and Wellner (1996). Conditions (2.3) and (2.6) are directly assumed at the time of writing.

C Extension

It is possible to analyze the more general model below. For any $z_t \in \mathcal{Z}_t$,

$$Y_{it} = h_t(\rho(X_{it}) + \alpha_i + U_{it}), \quad E(U_{it} - U_{it-1} | Z_{it} = z_t) = 0. \quad (49)$$

This model nests that of Fève and Florens (2014) when $h_t(s) = s$.

Assuming that the instrumental variable is time-invariant obtains for $t = 2$:

$$E(h_2^{-1}(Y_2) - h_1^{-1}(Y_1) + \rho(X_2) - \rho(X_1) | Z = z) = 0. \quad (50)$$

Via an observational equivalence argument as above with (g_1, g_2, R) that are observationally equivalent to (h_1, h_2, ρ) and, in particular, that satisfy

$$E(g_2^{-1}(Y_2) - g_1^{-1}(Y_1) + R(X_2) - R(X_1) | Z = z) = 0, \quad (51)$$

subtracting (51) from (50) obtains

$$E(\tilde{\delta}_2(Y_2) - \tilde{\delta}_1(Y_1) + r(X_2) - r(X_1) | Z = z) = 0, \quad (52)$$

where

$$\tilde{\delta}_t(Y_t) \equiv h_t^{-1}(Y_t) - g_t^{-1}(Y_t), \quad t = 1, 2, \quad (53)$$

and

$$r(X_t) = \rho(X_t) - R(X_t), \quad t = 1, 2. \quad (54)$$

As before, the identification argument involves completeness and measurable separability assumptions.

Assumption 15. (i) $E(h_t^{-1}(Y_t)|Z) \in L_Z^2$, $E(\rho(X_t)|Z) \in L_Z^2$, $t = 1, 2$; (ii) The random variables (Y_1, Y_2, X_1, X_2) are strongly identified by Z , i.e. for $\tilde{\delta}_t \in L_2(Y_t)$, $r \in L_2(X_t)$, $t = 1, 2$, defined in (53) and (54), respectively, if

$$E\left(\tilde{\delta}_2(Y_2) - \tilde{\delta}_1(Y_1) + r(X_2) - r(X_1) \middle| Z\right) = 0 \text{ a.s. } F_Z$$

then

$$\tilde{\delta}_2(Y_2) - \tilde{\delta}_1(Y_1) + r(X_2) - r(X_1) = 0 \text{ a.s. } F_{Y_1, Y_2, X_1, X_2};$$

(iii) The random variables Y_1, Y_2, X_1, X_2 are measurably separable in the sense that for $\tilde{\delta}_t(Y_t)$, $r(X_t)$, $t = 1, 2$, defined in (53) and (54), respectively, if

$$\tilde{\delta}_2(Y_2) - \tilde{\delta}_1(Y_1) + r(X_2) - r(X_1) = 0 \text{ a.s. } F_{Y_1, Y_2, X_1, X_2},$$

then there exist constants $\tilde{c}_t, d_t \in \mathbb{R}$, $t = 1, 2$, such that

$$\begin{aligned} \tilde{\delta}_t(Y_t) &= \tilde{c}_t \text{ a.s. } F_{Y_t}, \\ r(X_t) &= d_t \text{ a.s. } F_{X_t}. \end{aligned}$$

(iv) $E(h_t^{-1}(Y_t)) = 0$, $E(\rho(X_t)) = 0$, $t = 1, 2$.

Theorem 4. Let Assumptions 1 and 15 hold, and let Y_t, X_t, Z_t satisfy (49). Then h_t and ρ are identified.

Proof. Consider (52). By Assumption 15(ii), it follows that:

$$\tilde{\delta}_2(Y_2) - \tilde{\delta}_1(Y_1) + r(X_2) - r(X_1) = 0 \text{ a.s.}$$

By Assumption 15(iii) there exist constants $\tilde{c}_t, d_t \in \mathbb{R}$ such that $\tilde{\delta}_t(Y_t) = \tilde{c}_t$ a.s., $r(X_t) = d_t$ a.s. By Assumption 15(iv) these constants are all equal to zero. \square

D Illustration: A nonlinear dynamic panel model

The nonlinear panel model studied in this paper nests a nonlinear version of the canonical linear panel data model. Consider the outcome equations in (3) and (4). Setting $h_t(v) = v$ for $t \geq 1$ obtains the outcome equation of the standard dynamic panel model. The period-0 equation has its own ϕ_i that can capture i -specific terms regarding to the initial condition, the history of a given unit i , and a period-0 error term.

In the linear version of this model, estimation via differences is problematic because $\Delta Y_{i,t-1}$ is correlated with ΔU_{it} . Naturally, this problem carries over to the nonlinear generalization. We can use internal instruments to address this endogeneity issue. Internal instruments are available under a strict exogeneity assumption and restrictions on the serial correlation in U_{it} .

Assumption 16. *For each t , $E[U_{it}]$ does not depend on t , and*

$$U_{it} \perp \left(\tilde{X}_{i0}, \dots, \tilde{X}_{iT}, \alpha_i, \phi_i, U_{i1}, \dots, U_{i,t-1} \right).$$

This assumption is stronger than necessary: serial independence in the errors and the strict exogeneity condition on the regressors can be relaxed to a form of mean-independence. However, given the nonlinear nature of our model, it will be convenient to maintain statistical independence.

To place the nonlinear panel model in the notation of the general specification above, set

$$\begin{aligned} X_{it} &= \left(\tilde{X}_{it}, Y_{i,t-1} \right), \\ \bar{\beta} &= \left(\tilde{\beta}, \rho \right) \end{aligned}$$

and assume that \tilde{X}_{it} is non-empty. Then we can rewrite the nonlinear dynamic panel model as

$$Y_{it} = h_t(\alpha_i + X_{it}\bar{\beta} + U_{it}), i = 1, \dots, n, t = 1, \dots, T, \quad (55)$$

and, with three periods of data on Y_i and two periods on X_i , we can use as instruments $Z_i = (Y_{i0}, \tilde{X}_{i1}, \tilde{X}_{i2})$ for the difference $U_{i2} - U_{i1}$. We have the following result:

Theorem 5. *Suppose that $(Y_0, Y_1, Y_2, \tilde{X}_1, \tilde{X}_2)$ follow the nonlinear dynamic panel model above and that Assumptions 1, 3, 4, either 5 or 6, and 16 hold. Then $h_1, h_2, \tilde{\beta}$, and ρ are identified.*

Proof. This follows immediately from Theorem 1 once we verify that Assumption 16 implies Assumption 3, i.e. that

$$E(U_{i,2} - U_{i,1} | Y_0, X_1, X_2) = 0.$$

But this follows immediately from the fact that $(U_{i1}, U_{i2}) \perp (\phi_i, X_{i1}, X_{i2})$. \square

Note that it may be possible to relax the completeness assumption for this dynamic model along the lines of those in Fève and Florens (2014). It may also be possible to derive sufficient conditions for the completeness assumptions using arguments from d'Haultfoeuille (2010).