# Nonlinear Correlated Random Effects Models with Endogeneity and Unbalanced Panels

Michael D. Bates*, Leslie E. Papke†, Jeffrey M. Wooldridge ‡

September 13, 2022

## Abstract

We present simple procedures for estimating nonlinear panel data models in the presence of unobserved heterogeneity and possible endogeneity with respect to time-varying unobervables. We combine a correlated random effects approach with a control function approach while accounting for missing time periods for some units. We examine the performance of the approach in comparisons with standard estimators using Monte Carlo simulation. We apply the methods to estimating the effects of school spending on student pass rates on a standardized math exam. We find that a 10 percent increase in spending leads to an approximately two percentage point increase in math pass rates.

---

*Department of Economics, University of California at Riverside, Riverside, CA 92521 (email: mbates@ucr.edu).
†Department of Economics, Michigan State University, East Lansing, MI 48824-1038 (email: papke@msu.edu).
‡Department of Economics, Michigan State University, East Lansing, MI 48824-1038 (email: wooldri1@msu.edu).

# 1 Introduction

Unbalanced panel data – where some units do not have a complete set of observations in some time periods, are prevalent in empirical work. Researchers have documented unbalancedness stemming from both intermittent non-response and early attrition in panel surveys used in research in labor, public, and development economics. Aughinbaugh (2004) and Falaris and Peters (1998) note yearly non-response rates of three to four percent in the 1979 National Longitudinal Survey of Youth. Fitzgerald (2011) notes that only one-third of the children from the 1968 round of the Michigan Panel Survey of Income Dynamics remain in the data by 2007. Alderman et al. (1999) find significant annual attrition rates in household surveys from seven different developing countries that range from two to 20 percent.

As is well known, unbalanced panel data in a linear model context can be handled by fixed effects estimation provided the selection is based on observed variables or unobserved, time-constant heterogeneity; see, for example, Wooldridge (2019). When explanatory variables are endogenous with respect to time-varying unobservables, Joshi and Wooldridge (2019) show how linear fixed effects and control function methods can be applied to unbalanced panels for estimation and specification testing. But as pointed out in Wooldridge (2019), unbalanced panels cause significantly more difficulties in nonlinear panel data models. Wooldridge (2019) proposes a correlated random effects (CRE) approach to allow the heterogeneity to be correlated with time-constant functions of selection indicators for general nonlinear panel data models.

The CRE approach allows explanatory variables also to be correlated with time-constant unobservables – so-called "unobserved heterogeneity." In some cases, one might be concerned that a key explanatory variable is correlated with unobserved time-varying variables. In the panel data literature, this is called a failure of the "strict exogeneity" assumption. Failure of strict exogeneity is often due to omitted time-varying variables, or feedback from shocks to future outcomes of the explanatory variables. Simultaneity and (time-varying) measurement error can also cause failure of strict exogeneity.

In this paper we extend Wooldridge (2019) to nonlinear estimation of unbalanced panel data where the covariates may be endogenous with respect to time-constant heterogeneity as well as time-varying unobservables. Our work can also be viewed as extending Joshi and Wooldridge (2019), who consider linear models with unbalanced panels, to a nonlinear context. Our key assumption is that

the missingness of data is not correlated with idiosyncratic shocks. When explanatory variables are allowed to be endogenous with respect to idiosyncratic shocks, we require time-varying instrumental variables that are exogenous with respect to those shocks.

The approach we take is to combine the CRE approach for unbalanced panels – which we refer to as "CREU" for shorthand – with the control function approach when strict exogeneity fails. We consider different strategies for allowing correlation between unobserved heterogeneity and the selection indicators. We are specifically interested in comparing the CREU approach for nonlinear fractional response models, implemented using pooled quasi-maximum likelihood estimation (QMLE), with standard fixed effects estimation strategies for linear unobserved effects models. We find that the CREU approaches perform comparably to the CRE approach that ignores the unbalanced nature of the panel and linear fixed effects estimation in uncovering average partial effects (APEs). The CREU approach provides efficiency gains in estimating APEs and, because the fractional response model is nonlinear, allows us to study partial effects at different values of the key explanatory variables.

We illustrate our approach with an empirical application in the economics of education literature: estimating the effects of school spending on school pass rates of fourth graders on the Michigan state mathematics standardized exam. Papke (2005) used unbalanced school-level data and linear models estimated by fixed effects and instrumental variables. Given the bounded nature of the pass rate, a linear model may not be the best way to estimate average effects or effects at different points in the spending distribution. Papke and Wooldridge (2008) showed how to adapt fractional response models to a panel data setting, but they assumed a balanced panel and applied a combined correlated random effects/control function approach to balanced district-level data. We reexamine the results from Papke (2005) while accounting for the unbalancedness of the school-level data and the bounded nature of pass-rates. While we find some evidence of correlation between school spending and unbalancedness, our results largely uphold the evidence presented in Papke (2005): a 10% increase in spending leads to an approximately two percentage point increase in math pass rates, though our inference is sensitive to specification and level of clustering. These results are also similar to those in Papke (2008) using the balanced district-level data.

We organize the remainder of the paper as follows. In Section 2 we present the model and estimation methods, considering first the case where all explanatory variables are exogenous with

respect to the time-varying unobservables. We then derive a method that combines an extended version of the Mundlak (1978) device and a control function method to allow some explanatory variables to be correlated with time-varying unobservables. We present our simulation evidence in Section 3. In our application in Section 4, we demonstrate estimation with and without requiring school spending to be strictly exogenous with respect to idiosyncratic shocks in determining the effects of spending on fourth grade math pass rates. Section 5 concludes.

## 2 Model and Estimation

We begin with a population from which we draw a random sample of $N$ cross-sectional units. For each random draw $i$ from the cross section, there are potentially $T$ time observations, $t = 1, ..., T$, containing an outcome, $y_{it}$, and a vector of observed covariates, $\mathbf{x}_{it}$. Except for specific functional form and distributional assumptions, the approach proposed here applies to nonlinear models in general, but we focus on the case where $y_{it}$ is a fractional response that may take values at the endpoints in $[0, 1]$. Along with the $\mathbf{x}_{it}$, we expect unobserved heterogeneity, $c_i$, to play a role in determining $y_{it}$. In nonexperimental settings, it is likely that $c_i$ is correlated with at least some components of $\mathbf{x}_{it}$. We use a correlated random effects strategy to allow all elements of $\mathbf{x}_{it}$ that vary somewhat across $i$ and $t$ to be correlated with $c_i$. When one or more elements of $\mathbf{x}_{it}$ is correlated with underlying idiosyncratic shocks to $y_{it}$ – to be made precise shortly – we will assume the availability of some time-varying instrumental variables. Then, $\mathbf{z}_{it}$ will denote the vector of all variables strictly exogenous with respect to shocks. We still allow all elements of $\mathbf{z}_{it}$ to be correlated with $c_i$.

To account for the unbalanced nature of the panel data, we introduce a selection indicator – also known as a "complete cases" indicator, $s_{it}$. This indicator is one if we observe the outcome, all covariates, and any instrumental variables for unit $i$ in time $t$. It is important in what follows that only the complete cases are used, as using incomplete cases generally requires more assumptions and more complications. The default of estimation methods in econometrics packages is to use a data point only if all necessary variables are observed, and that is what the definition of $s_{it}$ captures. Therefore, $s_{it} = 1$ means we use observation $(i, t)$ in the estimation and $s_{it} = 0$ means we do not. The series of selection indicators for unit $i$ is $\{s_{i1}, ..., s_{iT}\}$.

## 2.1 Strict Exogeneity

We begin with the case where the explanatory variables are strictly exogenous conditional on the heterogeneity. The population model, written for a random draw $i$, is

$$E(y_{it}|\mathbf{x}_i, c_i) = E(y_{it}|\mathbf{x}_{it}, c_i) = \Phi(\mathbf{x}_{it}\boldsymbol{\beta} + c_i), t = 1, ..., T, \tag{1}$$

where $\mathbf{x}_i = (\mathbf{x}_{i1}, ..., \mathbf{x}_{iT})$ is the entire history of the covariates and $\Phi(\cdot)$ is the standard normal cumulative distribution function. Our use of $\Phi$ rather than some other cumulative distribution function leads to simple procedures in the presence of unobserved heterogeneity and easy calculation of average partial effects. It is also convenient when we have endogenous explanatory variables.

To account for sample selection, let $\mathbf{s}_i = (s_{i1}, ..., s_{iT})$ be the entire history of selection. We assume that, conditional on $\mathbf{x}_i$ and the unobserved heterogeneity, selection is strictly exogenous in the following sense:

$$E(y_{it}|\mathbf{x}_i, \mathbf{s}_i, c_i) = E(y_{it}|\mathbf{x}_i, c_i), \ \ t = 1, ..., T \tag{2}$$

This assumption allows selection to be arbitrarily correlated with both the explanatory variables and unobserved heterogeneity – because we are conditioning on them – but rules out correlation between selection and unobserved idiosyncratic fluctuations in the outcome.

Following Wooldridge (2019), we use a correlated random effects approach to specify a model for the following conditional distribution:

$$D(c_i|\{(s_{it}\mathbf{x}_{it}, s_{it}) : t = 1, ..., T\}), \tag{3}$$

where multiplying the covariates by the selection indicator reflects our usage of complete cases only. Generally, Wooldridge (2019) suggests modeling (3) as fairly simple time-constant functions, say $\mathbf{w}_i$, of $\{(s_{it}\mathbf{x}_{it}, s_{it}) : t = 1, ..., T\}$ that effectively act as sufficient statistics in the relationship between the covariates and selection. It is natural to extend the Mundlak (1978) device to the unbalanced case by using the time averages

$$\bar{\mathbf{x}}_i = T_i^{-1} \sum_{t=1}^{T_i} \mathbf{x}_{it},$$

where $T_i = \sum_{t=1}^{T} s_{it}$ is the number of complete cases for unit $i$. (If $T_i = 0$, then there are no complete time periods for unit $i$, and such units are not used in the estimation). To handle correlation between $c_i$ and selection, we use a flexible mean specification where the intercept and slopes can depend on the number of complete cases, as given by the indicators $1[T_i = r]$, which are one if and only if unit $i$ has $r$ complete cases. Then,

$$E(c_i|\mathbf{w}_i) = \sum_{r=1}^{T} \psi_r 1[T_i = r] + \sum_{r=1}^{T} \left(1[T_i = r] \cdot \bar{\mathbf{x}}_i\right) \boldsymbol{\xi}_r. \tag{4}$$

If we also assume $D(c_i|\mathbf{w}_i)$ is a normal distribution, then we have

$$E(y_{it}|\mathbf{x}_{it}, \mathbf{w}_i, s_{it} = 1) = \Phi\left(\frac{\mathbf{x}_{it}\boldsymbol{\beta} + \sum_{r=1}^{T} \psi_r 1[T_i = r] + \sum_{r=1}^{T} 1[T_i = r] \cdot \bar{\mathbf{x}}_i \boldsymbol{\xi}_r}{+Var(c_i|\mathbf{w}_i)\}^{\frac{1}{2}}}\right), \tag{5}$$

because a mixture of independent normal distributions is normal. Equation (5) extends Papke and Wooldridge (2008), who assumed $Var(c_i|\mathbf{w}_i)$ is constant, to the case of unbalanced panels. Rather than assume $Var(c_i|\mathbf{w}_i)$ is constant, it is natural to allow, at a minimum, the variance of $c_i$ to vary with the number of complete cases. A simple way to do this is

$$Var(c_i|\mathbf{w}_i) = \exp\left(\tau + \sum_{r=1}^{T-1} 1[T_i = r]\boldsymbol{\omega}_r\right), \tag{6}$$

where $\exp(\tau)$ is the variance for the complete-cases base group ($T_i = T$) and each $\omega_r$ captures the deviation from the base group.

Combining (5) and (6), we have

$$E(y_{it}|\mathbf{x}_{it}, \mathbf{w}_i, s_{it} = 1) = \Phi\left(\frac{\mathbf{x}_{it}\boldsymbol{\beta} + \sum_{r=1}^{T} \psi_r 1[T_i = r] + \sum_{r=1}^{T} 1[T_i = r] \cdot \bar{\mathbf{x}}_i \boldsymbol{\xi}_r}{\left\{1 + \exp\left(\tau + \sum_{r=1}^{T-1} 1[T_i = r]\boldsymbol{\omega}_r\right)\right\}^{\frac{1}{2}}}\right). \tag{7}$$

For all $r \geq 2$ these scaled coefficients are identified as long as there is some time variation in all elements of $\mathbf{x}_{it}$ and no perfect collinearity among the elements of $\mathbf{x}_{it}$.

Given the expression (7) for the conditional mean, we can follow Papke and Wooldridge (2008) in estimating the parameters using a pooled quasi-maximum likelihood approach with the log-likelihood being chosen to be that for the Bernoulli distribution. Given the functional form in (7), the pooled quasi-log-likelihood function is equivalent to that from a particular heteroskedastic probit

model, where the heteroskedasticity function is $1 + \exp\left(\tau + \sum_{r=1}^{T-1} 1[T_i = r]\boldsymbol{\omega}_r\right)$. As a practical matter, we can drop the "1+" term because we allow an intercept $\tau$ inside the exponential function. The resulting parameters in the "mean" function, $\mathbf{x}_{it}\boldsymbol{\beta} + \sum_{r=1}^{T} \psi_r 1[T_i = r] + \sum_{r=1}^{T} 1[T_i = r] \cdot \bar{\mathbf{x}}_i\boldsymbol{\xi}_r$, get rescaled, but this does not affect estimating the magnitudes of the effects. It is easy to use software, such as Stata, that has a command for estimating fractional response models with heteroskedasticity. In obtaining proper standard errors and inference, we obtain a cluster-robust variance-covariance matrix estimator that accounts for both serial correlation and the fact that the variance $Var(y_{it}|\mathbf{x}_{it}, \mathbf{w}_i, s_{it} = 1)$ does not have the same form as when $y_{it}$ is a binary variable.

Estimating the average partial effects – the quantities typically of interest – requires some care if data are missing on the $\mathbf{x}_{it}$. At a minimum, we can plug in reasonable values of the covariates and average across the functions of $(\mathbf{x}_i, \mathbf{s}_i)$ that act as proxies for the heterogeneity. This leads to

$$\widehat{APE}_j(\mathbf{x}_t) = \hat{\beta}_j \left[N^{-1} \sum_{i=1}^{N} \phi\left(\frac{\mathbf{x}_t\hat{\boldsymbol{\beta}} + \sum_{r=1}^{T} \psi_r 1[T_i = r] + \sum_{r=1}^{T} (1[T_i = r] \cdot \bar{\mathbf{x}}_i)\hat{\boldsymbol{\xi}}_r}{1 + \exp\left(\hat{\tau} + \sum_{r=1}^{T-1} 1[T_i = r]\hat{\boldsymbol{\omega}}_r\right)}\right)\right].$$

It is harder to obtain an effect averaged across the distribution of $\mathbf{x}_{it}$ because data may be missing as a systematic function of $\mathbf{x}_{it}$. The simplest approach is to average the APEs across the selected observations:

$$\widehat{APE}_j = \hat{\beta}_j \left[N^{-1} \sum_{i=1}^{N} T_i^{-1} \sum_{t=1}^{T} s_{it}\phi\left(\frac{\mathbf{x}_{it}\hat{\boldsymbol{\beta}} + \sum_{r=1}^{T} \psi_r 1[T_i = r] + \sum_{r=1}^{T} (1[T_i = r] \cdot \bar{\mathbf{x}}_i)\hat{\boldsymbol{\xi}}_r}{1 + \exp\left(\hat{\tau} + \sum_{r=1}^{T-1} 1[T_i = r]\hat{\boldsymbol{\omega}}_r\right)}\right)\right].$$

As an extension, $\bar{\mathbf{x}}_i$ can be added to the variance function along with interactions between the dummies $1[T_i = r]$ and $\bar{\mathbf{x}}_i$.

## 2.2 Endogenous Explanatory Variables

In many applications, researchers are hesitant to assume strict exogeneity of covariates. In our application, we worry that deviations in school spending may be linked with unobserved fluctuations in student performance. This may come from unobserved demands of cohorts or accountability pressure, as depicted in Chiang (2009).

Here we present a straightforward approach to handle endogeneity of an explanatory variable $y_{it2}$ in nonlinear panel data models in the presence of unobserved heterogeneity and panel imbalance. We first assume the presence of instrumental variables $\mathbf{z}_{it2}$ that are both relevant to $y_{it2}$ and otherwise exogenous. We more precisely state these assumptions below. We allow there to be additional exogenous covariates denoted as $\mathbf{z}_{it2}$ with $\mathbf{z}_{it} = (\mathbf{z}_{it1}, \mathbf{z}_{it2})$ denoting the complete vector of pertinent exogenous variables. We follow Papke and Wooldridge (2008) in modeling the conditional mean as

$$
\begin{aligned}
E\left(y_{it1} | y_{it2}, \mathbf{z}_i, \mathbf{s}_i, c_{i1}, v_{it1}\right) =&E\left(y_{it1} | y_{it2}, \mathbf{z}_i, c_{i1}, v_{it1}\right) = E\left(y_{it1} | y_{it2}, \mathbf{z}_{it2}, c_{i1}, v_{it1}\right) \\
=&\Phi\left(\beta_1 y_{it2} + \mathbf{z}_{it1}\delta_1 + c_{i1} + v_{it1}\right),
\end{aligned}
\tag{8}
$$

where $c_{i1}$ is time-invariant unobserved heterogeneity across units and $v_{it1}$ is an omitted factor that varies over both units and time. Once we have already conditioned on the explanatory variables and the source of endogeneity, the conditional mean is unaffected by conditioning on $\mathbf{z}_{it2}$. Thus, $\mathbf{z}_{it2}$ is excluded from equation (8). Additionally, note that we continue to assume that selection is ignorable conditional on the observed variables and the unobservables, $c_{i1}$ and $v_{it1}$.

In equation (8) the variable $y_{it2}$ may now be endogenous with respect to $v_{it1}$ as well as with respect to $c_{i1}$. We handle the latter endogeneity similarly to the strict exogeneity case by using a correlated random effects approach to specify a model for $c_{i1}$, following the unbalanced case in Wooldridge (2019). In particular, to account for the unbalanced panel, we allow the coefficients on the time averages to change with the number of time periods observed for each $i$, in addition to allowing separate intercepts for each $T_i$:

$$
c_{i1} = \sum_{r=1}^{T} \psi_{r1} 1[T_i = r] + \sum_{r=1}^{T} \left(1[T_i = r] \cdot \bar{\mathbf{z}}_i\right) \boldsymbol{\xi}_{r1} + a_{i1}, \ a_{i1} | \mathbf{z}_i \sim Normal(0, \sigma_{a1}^2),
\tag{9}
$$

where $\bar{\mathbf{z}}_i = T_i^{-1} \sum_{r=1}^{T_i} s_{it} \mathbf{z}_{it}$ is the time average over the complete cases and $a_{i1}$ is an error term that we assume to be independent of $(\mathbf{z}_i, \mathbf{s}_i)$. As before, conditional normality leads to a relatively straightforward analysis.

Substituting equation (9) into equation (8) gives

$$E(y_{it1}|y_{it2}, \mathbf{z}_i, r_{it1}, s_{it} = 1) =$$

$$\Phi\left(\beta_1 y_{it2} + \mathbf{z}_{it1}\delta_1 + \sum_{r=1}^{T} \psi_{r1}1[T_i = r] + \sum_{r=1}^{T} (1[T_i = r] \cdot \bar{\mathbf{z}}_i) \boldsymbol{\xi}_{r1} + r_{it1}\right), \tag{10}$$

where $r_{it1} = a_{i1} + v_{it1}$ is a composite error term. Researchers may also wish to follow Lin and Wooldridge (2019) by including $\bar{y}_{2i} = T^{-1}\sum_{r=1}^{T} y_{2ir}$ to clearly separate the endogeneity due to $c_{i1}$ from the endogeneity due to $v_{it1}$. We omit it here to coincide with previous approaches in our application.

Secondly, we must deal with the endogeneity of $y_{it2}$. Following Mundlak (1978), we linearly model $y_{it2}$ as a function of the exogenous explanatory variables, excluded instruments, and their time averages. As selection may be correlated with $y_{it2}$, we generally include indicators for the number of time observations and interactions with time averages here as well. We present this first-stage equation below:

$$y_{it2} = \mathbf{z}_{it}\pi_2 + \sum_{r=1}^{T} \psi_{r2}1[T_i = r] + \sum_{r=1}^{T} 1[T_i = r] \cdot \bar{\mathbf{z}}_i\boldsymbol{\xi}_{r2} + v_{it2}, \tag{11}$$

where $v_{it2}$ represents time-varying unobserved elements of $y_{it2}$ and we have included a full set of dummies and omitted an intercept. In equation (11) the endogeneity in $y_{it2}$ is due to the correlation between $r_{it1}$ and $v_{it2}$. Following Rivers and Vuong (1988) and Papke and Wooldridge (2008), we model $r_{it1}$ as linear in $v_{it2}$ and conditionally normal:

$$r_{it1} = \eta_1 v_{it2} + e_{it1}, \ e_{it1} \,|\, (\mathbf{z}_i, \mathbf{s}_i, v_{it2}) \sim Normal(0, \sigma_e^2).$$

Note that $e_{it1}$ is also independent of $y_{it2}$. As in the balanced case in Papke and Wooldridge (2008), given the assumptions, we can replace $r_{it1} = \eta_1 v_{it2} + e_{it1}$ and then integrate out $e_{it1}$ using the properties of the normal distribution. The resulting coefficients are scaled by $(1 + \sigma_e^2)^{-\frac{1}{2}}$:

$$E(y_{it1}|y_{it2}, \mathbf{z}_i, v_{it2}, s_{it} = 1) =$$

$$\Phi\left(\beta_{1e} y_{it2} + \mathbf{z}_{it1}\delta_1 + \sum_{r=1}^{T} \psi_{r1e}1[T_i = r] + \sum_{r=1}^{T} (1[T_i = r] \cdot \bar{\mathbf{z}}_i) \boldsymbol{\xi}_{r1e} + \eta_{1e} v_{it2}\right), \tag{12}$$

9

where subscript $e$ denotes the scaling of the coefficients. The average partial effects – where we necessarily average over the selected sample – depend on the scaled coefficients, as discussed in Papke and Wooldridge (2008) in the balanced panel case. Therefore, in what follows, we drop the $e$ subscript from the parameters.

We follow a two-step procedure to estimate equation (12). In the first step, we estimate (11) by regressing our endogenous explanatory variable, $y_{it2}$, on the exogenous variables, $\mathbf{z}_{it}$, that include the instruments and time indicators, indicators for the number of time observations per unit, and interactions between those indicators and time averages of the exogenous variables. We save the residuals from that regression, $\hat{v}_{it2}$, for the complete cases. In step two, we substitute these residuals for $v_{it2}$ and estimate equation (12) using the complete cases and pooled probit QMLE of $y_{iy1}$ on $\mathbf{z}_{it1}$; the indicators, $1[T_i = r]$; all interactions, $1[T_i = r] \cdot \bar{\mathbf{z}}_i$; and the first-stage residuals, $\hat{v}_{it2}$.

Due to the estimation of $\hat{v}_{it2}$ in the first step, the standard errors in the second stage should be adjusted. Bootstrapping the entire procedure by resampling individual units with replacement is one way to account for the first stage estimation. We adopt this approach in our empirical application.

We are mainly interested in the APE of the endogenous explanatory variable, $y_{it2}$. To obtain it, we first follow Blundell and Powell (2003) in defining an average structural function (ASF) for $y_{it1}$, such that

$$ASF(y_{2t}, \mathbf{z}_{1t}) = \Phi \left( \beta_e y_{t2} + \mathbf{z}_{t1} \boldsymbol{\delta}_e \right) \tag{13}$$

is the conditional mean with $\bar{\mathbf{z}}_i$ averaged out to account for the heterogeneity and $v_{it2}$ averaged out to account for the contemporaneous endogeneity.

In practice, we uncover the APE by averaging out $y_{it2}, \mathbf{z}_{it1}, \sum_{r=1}^{T} 1[T_i = r], \sum_{r=1}^{T} 1[T_i = r] \cdot \bar{\mathbf{z}}_i$, and $\hat{v}_{it2}$ across the sample for a given $t$, and difference or differentiate. For instance, in the sample, we estimate the APE of $y_{it2}$ as

$$\hat{\beta}_e \cdot \left( N^{-1} \sum_{i=1}^{N} \phi[\hat{\beta}_e y_{it2} + \mathbf{z}_{it2} \hat{\boldsymbol{\delta}}_e + \sum_{r=1}^{T} \hat{\psi}_{re1} 1[T_i = r] + \sum_{r=1}^{T} 1[T_i = r] \cdot \bar{\mathbf{z}}_i \hat{\boldsymbol{\xi}}_{re1} + \hat{\eta}_e v_{it2} \right), \tag{14}$$

where the "ˆ" denotes that the coefficients have been estimated by pooled probit QMLE. Again, applying a clustered bootstrap is a convenient way to obtain valid standard errors.

# 3 Simulation Evidence

We conduct a simulation study to investigate the performance of approaches that handle the panel unbalancedness against standard estimators that do not. In particular, we are interested in the bias that correlated unbalancedness may produce in estimated APEs and the relative efficiency of the estimators. For comparison, we first use POLS and FE as approximations of the APE from linear models. We then consider standard nonlinear approaches; namely, pooled fractional response probit QMLE (PFR), and PFR where we model the correlated random effects using the time averages of covariates (CRE).

In using CRE in the linear case, once $\bar{\mathbf{x}}_i$ has been included, interacting the indicators, $1[T_i = r]$, with $\bar{\mathbf{x}}_i$ does not change the estimates; they are the usual fixed effects estimates. This follows from Wooldridge (2019). To handle the imbalance of the panel in the linear model, we mimic fixed effects estimation using time averages of all explanatory variables and add time-observation indicators interacted with time averages of the explanatory variables (FEU). In nonlinear formulations, we add to CRE time-observations indicators (CREU). We next add to CREU interactions between time-observation indicators and covariate time averages (CREU1). We then add to CREU1 interactions between time-observation indicators and the covariates themselves (CREU2). Finally, we also include triple interactions between time-observation indicators, the covariates, and their time averages (CREU3).

## 3.1 Data Generating Process

We generate the data using a slight generalization of our single-stage model of interest. We consider two cases ($g = 0, 1$) and generate the outcome, $y$, according to the following:

$$y_{it} = \Phi[\alpha + (\beta_1 + u_{ig})x_{it1} + \beta_2 x_{it2} + c_i + v_{it}], \ v_{it} \sim Normal(0, 0.2), \tag{15}$$

where $\Phi$ is the standard normal CDF and $\beta_1 = \beta_2 = 1$. We first draw school-year variables $x_{it1} \sim Normal(0, \ 0.2)$ and $x_{it2} \sim Binomial(1, \ 0.3)$ over 500 "schools" across 5 "years," and generate time averages of these variables by school building. The standard deviation of $x_{it1}$ is set to approximate the standard deviation of our spending variable in the empirical application. The generalization in this simulation is that unobserved school heterogeneity takes the form of fixed effects, $c_i$, and

random coefficients, $u_{ig}$.

The two cases differ depending on the construction of the unobserved heterogeneity and which form of heterogeneity is correlated with selection. In both cases we generate the unobserved fixed effects according to the following equation:

$$c_i = \sqrt{T}\bar{x}_{i1}\gamma_1 + \eta_i, \ \eta_i \sim Normal(0, 0.14). \tag{16}$$

However, in the first case, $g = 0$, the school-level, random slope of $x_1$, $u_{i0}$, is defined as the time average of an independently distributed normal random variable, and thus, is not correlated with $x_1$ and selection into the panel. Consequently, it is not a far departure from the standard model introduced above. In the second case, $g = 1$, we extend the data generating process to include a correlated random coefficient on $x_1$. This random slope, $u_{i1}$, is correlated with $x_1$ and selection into the panel. Specifically, the random slopes take the following form:

$$u_{ig} = \begin{cases} u_{i0} = T^{-1}\sum_{t=1}^{T} e_{it0}, \ e_{it0} \sim Normal(0, 0.14) \text{ in simulation one,} \\ u_{i1} = \sqrt{T} \times \bar{x}_{i1}\gamma_2 + \gamma_3 c_i + e_{i1}, \ e_{i1} \sim Normal(0, \ 0.14) \text{ in simulation two,} \end{cases} \tag{17}$$

where $T$ represents the five possible time-observations, $\gamma_1$ and $\gamma_2$ are each set to 0.7, and $\gamma_3$ is set to 0.2, and $\eta_i$, $e_{i0}$, and $e_{i1}$ are each drawn from independent, mean-zero, normal distributions.

We model selection depending on the unobserved effect, $c$, in simulation one and on the unobserved correlated random slope, $u_1$, in simulation two. In both cases, the selection of each time-observation is drawn from a binomial distribution with probability $p_{ig}$ defined below.

$$p_{ig} = \begin{cases} \Phi(a_{it} + c_i) \text{ in simulation one,} \\ \Phi(a_{it} + u_{i1}) \text{in simulation 2,} \end{cases} \tag{18}$$

where $a_{it}$ is an independent normal distributed random variable with a mean of 0.75 and a standard deviation of 0.2.

## 3.2 Simulation Results

We present the resulting correlations from simulation one on the left and simulation two on the right of Table 1. In the first case, only the unobserved fixed heterogeneity is correlated with time-observation selection and with $x_1$. The resulting average number of time-observations across the 500 replications is 3.83, with a correlation of 0.299 between the number of time-observations and the unobserved fixed effect. In contrast, the correlation between the number of time-observations and the random slope is 0.001. The correlation between $x_1$ and $c$ is 0.315, while the correlation between $x_1$ and $u$ is 0.0006.

Due to the positive correlation between $x_1$ and $c$, we may expect POLS and PFR to exhibit an upward bias for $\beta_1$. Indeed, we see exactly this in Table 2. The first row of Table 2 provides the "true" APEs of $x_1$ when averaged over the population (as if the panel were balanced), the sample (where the number of time-observations is non-randomly unbalanced), and then disaggregated by each number of time-observations that the schools are present in the data. Over the population, the APE of $x_1$ is 0.2979, and averaged over the sample, the APE of $x_1$ is 0.2953. Both POLS and PFR overstate this effect by 0.087. The bias is easily statistically significant, as the standard deviations of the POLS and PFR estimates over the 500 replications are 0.0111 and 0.0102 respectively, .

Beyond POLS and PFR, all estimated APEs are quite close to the true APE, and none are more than a third of a standard deviation of the estimated APEs away from the true APE over the population. Still, the FE estimated APE is the furthest from the truth, with an estimated APE of 0.2939, 1.3 percent lower than the true effect. Adding time-observation indicators to the time-means in FE estimation in FEU increases the estimated APE to 0.2947, 1.1 percent lower than the true APE.

Even without accounting for the unbalancedness of the panel, with an estimated APE of 0.2947, the CRE estimates are remarkably close to those estimated by FEU. Neither adding indicators for the number of time-observations in CREU nor including time-observation indicators interacted with time averages in CREU1 alter the estimated APEs to the fourth decimal place. Further, the standard deviations of the APE estimates remain remarkably similar among CRE (0.0097), CREU (0.0098), and CREU1 (0.0098).

We examine additional specifications by adding interactions between covariates and time aver-

ages of covariates to the covariates used in CREU1 estimation (labeled CREU2), and second, by including the triple interactions between the covariates, their time averages, and the number of time-observations (labeled CREU3). Including these additional interactions makes sense for a model that includes correlated random slopes, as shown in equation (15). Here, we model the heterogeneous school-level slopes by interacting the time averages of covariates with each covariate just as we model the fixed unobserved heterogeneity by inserting the time averages additively. Incorporating the triple interaction between the covariates, the time averages, and indicators for the number of time-observations in CREU3 addresses the potential that selection of time-observations is related to random slopes.

Both estimators provide very similar APE estimates to those from CRE. The CREU2 approach yields an estimated APE of 0.2949, while CREU3 estimates the APE of $x_1$ at 0.2951. The estimates become slightly less precise as we add covariates—the standard deviation of the CREU3 APEs increases to 0.0092.

We provide both the average standard error as well as the standard deviation of estimates across repetitions to show how well the estimated standard errors reflect the precision of each estimator. The standard errors of most estimators perform well, with the ratio of average standard errors to Monte Carlo standard deviations ranging from 0.93 to 1.02.

We turn next to the case where the selection of time-observations depends on the random coefficient of $x_1$, which is $u_1$. The average number of time-observation across the 500 replications is 3.82. The resulting correlation between the number of time-observations and the unobserved fixed effect is 0.2065, and the correlation between the number of time-observations and the random slopes is 0.3300. The correlation between $x_1$ and $c_0$ is 0.3152, while the correlation between $x_1$ and $u_1$ is 0.3404.

The results of the simulation with selection based upon correlated random slopes appear in Table 3. The "true" APE of $x_1$, when averaged over the population (with the panel balanced) is 0.2902. The correlation between selection and the heterogeneous slopes is apparent looking at the top row across the true APEs averaged across schools with one, two, three, four, or five time-observations appearing in the data. The relationship is monotonically positive with the APE among those with only one time-observation being 0.2523, whereas the APE among those with five time-observations is 0.3034. The true APE over the unbalanced sample is about 1 percent higher than the true APE

over the population. This positive correlation may be expected in many contexts where those with favorable numbers may be more likely to report their data. In our application, schools that report their data more frequently tend to be higher-performing.

In the presence of correlation between the heterogeneous slopes and selection of time-observations, all estimators overstate the average effect of the endogenous regressor. All estimates are greater than the true APE among the unbalanced sample. POLS and PFR overstate this effect most. POLS estimates the APE to be 0.3823 (standard deviation 0.0101). PFR probit estimates the APE to be 0.3827 (standard deviation 0.0094). Again, the true value lies far outside the 95% confidence interval of both estimators.

FE and CRE estimates lie significantly closer to the true estimates at 0.2954 and 0.2945, respectively. The CRE estimates (with a standard deviation of 0.0083) are more precise than the FE estimates (with a standard deviation of 0.0094). Adding indicators for the number of time-observations moves the FE estimates closer to the true APEs, though not statistically significantly so. The FEU estimate of the APE of $x_1$ is 0.2941 (standard deviation 0.0104).

Regarding the nonlinear estimators, adding indicators for the number of time-observations in CREU only marginally affects the estimates of the APE (0.2946) nor its precision (standard deviation of 0.0083). Adding interactions between time averages and time-observation indicators marginally increases the estimate of the APE to 0.2947 (standard deviation of 0.0095), though again not statistically significant. Only adding the triple interactions between the covariates, their time averages, and time-observation indicators in CREU3 makes a somewhat larger impact. With an estimated APE of 0.2937, CREU3 again provides the estimates closest to the true APE, though it is less precise with a standard deviation across simulations of 0.0096. Still, all CREU estimates fall within 0.3 percent of the estimated APE using the standard CRE approach.

Across specifications, the standard errors perform similarly to the standard deviations across repetitions. The ratio of mean standard errors to the standard deviation of the APEs across replications is 0.89 to 1.1 in this second simulation. CREU provides the most conservative standard errors relative to the standard deviation of estimates, and the ratio for CREU3 indicates that the standard errors perhaps overstate the estimator's precision. The nonlinear approaches mostly appear to be more efficient than the approaches using a linear specification. The mean standard errors are approximately 10 percent smaller using one of the fractional response probit specifications as opposed

to an analogous linear specification. The relative precision of the nonlinear estimators makes sense given the nonlinearity of the estimated effects. Panel C of Table 3 shows the estimated partial effects at the tenth, thirtieth, fiftieth, seventieth, and ninetieth decile of $x_1$ using the CREU1 approach. The estimated partial effect at the tenth percentile is 11 percent larger than the partial effect at the median and 31 percent larger than the estimated partial effect at the ninetieth percentile.

To summarize, under both formulations of selection of time-observations into the sample, all estimators that account for unobserved heterogeneity do comparably well in avoiding bias. Nonlinear estimators have the additional advantage of greater efficiency and the ability to detect nonlinear effects, particularly at the tails of the support. In the next section, we apply the methods to study the effect of school spending on student achievement.

## 4    Empirical Application

In revisiting Papke (2005), we initially treat spending as strictly exogenous and apply single-stage estimation of both linear and nonlinear models. Then, following Papke (2005), Chaudhary (2009), and Roy (2011), we use the 1994 centralization of school financing that occurred in Michigan under Proposal A to provide plausibly exogenous variation in school expenditures.[1] We use this policy to apply instrumental variables to spending and demonstrate these methods with an endogenous regressor. As in Papke (2005), we conduct our analysis at the school-building level over the time period of 1993-1998.

For building-level expenditure data, we use the average per pupil expenditures taken from the Michigan School Reports indexed for inflation using the Consumer Price Index normalized to 1997 dollars. As Papke (2005) notes, spending in previous years may impact students' fourth-grade math scores as might contemporaneous spending. Thus, we measure spending as the log of average real expenditures over the current and previous year. Our data contain 7,242 building-year observations from the 1,771 elementary schools in the state over the five year period when funding equalization was most dramatic. We focus on the effects of spending on *math4,* which measures the fraction of 4th grade students who pass the mathematics section of the Michigan Education Assessment Program (MEAP). During this time period, the mathematics and reading MEAP tests were only offered in grades 4 and 7. We focus on mathematics, as the reading MEAP exams changed format

---

[1]Papke (2005), Chaudhary (2009), and Roy (2011) provide fuller discussion of this school finance reform.

in the 1994/1995 school year and coincides with the policy change.

We note significant imbalance in the school-building-by-year panels, making this an appropriate application of the methods discussed in Wooldridge (2019). Table 4 demonstrates this unbalancedness. As the bottom row of Table 4 shows, 37 percent of schools are missing at least one of the five possible observations, and 23 percent are missing at least two observations. Further, there is significant variation in fourth-grade math pass rates, size, student composition, and spending across schools that appear in the data for each number of years, suggesting that the unbalancedness may be consequential for estimating the APE of spending. In addressing this unbalancedness, we move beyond Papke and Wooldridge (2008), who conduct analysis at the district level due to this issue.

## 4.1 Treating spending as strictly exogenous

The population linear model estimated in Papke (2005) can be written as the following:

$$
\begin{aligned}
math4_{it} =& \theta_t + \beta_1 log(avgrexp_{it}) + \beta_2 lunch_{it} + \beta_3 lunch_{it}^2 \\
& + \beta_4 log(enroll_{it}) + \beta_5 log(enroll_{it})^2 + c_i + e_{it}
\end{aligned}
\tag{19}
$$

We control for year indicators and quadratics of both the percent of free and reduced price lunch students and the log of student enrollment. This model may be estimated using pooled ordinary least squares (POLS). However, the estimated coefficients would be inconsistent if the unobserved heterogeneity is correlated with any of the explanatory variables. Consequently, researchers may opt to use fixed effects (FE) estimation or equivalently include time averages of all explanatory variables. The estimated coefficients may provide good approximations to the APEs when the actual model is nonlinear, but there is no general result that says so.

As noted previously, we do not observe each time-observation for each school building. Let $s_{it}$ represent an indicator for whether school $i$ appears in the data in year $t$. Thus, we can characterize the linear unobserved effects model with unbalanced data by multiplying equation (19) through by the selection indicator, $s_{it}$. Equation (20) represents this linear model in the presence of unbalancedness.

$$
s_{it} math4_{it} = s_{it}\theta_t + s_{it}\mathbf{x}_{it}\boldsymbol{\beta_a} + s_{it}c_i + s_{it}e_{it},
\tag{20}
$$

where $\mathbf{x}_{it}$ includes $log(avgrexp_{it})$, $lunch_{it}$, $lunch_{it}^2$, $log(enroll_{it})$, and $log(enroll_{it})^2$. In order for

fixed effects estimation to be consistent in the presence of such unbalancedness, we must assume strict exogeneity of the covariates *and* selection, conditional on the unobserved heterogeneity. Put more formally,

$$E(u_{it}|\mathbf{x}_i, \ \mathbf{s}_i, \ c_i) = 0, \tag{21}$$

where $\mathbf{x}_i = (\mathbf{x}_{i1}, \mathbf{x}_{i2}, ..., \mathbf{x}_{iT})$ and $\mathbf{s}_i = (s_{i1}, s_{i2}, ..., s_{iT})$. As an example, an idiosycratic low pass rate in year $t-1$ affecting selection in year $t$ would violate this condition.

In this setting our dependent variable, $math4$, is bounded between zero and one. While linear estimation may do well to approximate the average effect of spending on pass rates, researchers may opt to estimate a nonlinear model with a more plausible functional form to reduce bias in the APEs and possibly improve precision of the estimates. Using a nonlinear functional form has the added benefit of allowing estimation of partial effects at different points along the distribution of $\mathbf{x}_{it}$. Papke and Wooldridge (2008) estimate a fractional response probit unobserved effects model, where the unobserved heterogeneity is modeled using time averages of each covariate as in Chamberlain (1980) and Mundlak (1978). Accordingly, their correlated random effects (CRE) estimation equation is shown below:

$$E[math4_{it}|x_{i1}, x_{i2}, ..., x_{iT}] = \Phi(\psi_t + \mathbf{x}_{it}\boldsymbol{\beta} + \overline{\mathbf{x}}_i\boldsymbol{\xi}) \tag{22}$$

where $\overline{\mathbf{x}}_i$ includes the time averages of each covariate, $\Phi$ represents the normal CDF, and $\psi_t$ allows for year-specific intercepts.[2]

We use indicators for each number of times a particular school appears in the data, $\mathbf{T}_i = T_{1i}, T_{2i}, ...., T_{5i}$, as sufficient statistics for the dependence between the unobserved, school-level heterogeneity and the selection of time-observations into our data. Thus, in accounting for the unbalancedness of the school-level data, we initially estimate the following:

$$E[math4_{it}|\mathbf{x}_{it}, \overline{\mathbf{x}}_i, \mathbf{T}_i] = \Phi(\psi_t + \mathbf{x}_{it}\boldsymbol{\beta} + \mathbf{T}_i\boldsymbol{\gamma} + \overline{\mathbf{x}}_i\boldsymbol{\xi}) \tag{23}$$

We term the above approach Correlated Random Effects for Unbalancedness (CREU).

We also include a specification in which we interact time averages with indicators for the number of time-observations in an approach we label CREU1. The corresponding estimating equation

---

[2]As discussed above, the coefficients remain scaled by the variance of the unobserved heterogeneity.

appears below:

$$E[math4_{it}|\mathbf{x}_{it}, \overline{\mathbf{x}}_i, \mathbf{T}_i] = \Phi(\psi_t + \mathbf{x}_{it}\boldsymbol{\beta} + \mathbf{T}_i\boldsymbol{\gamma} + \overline{\mathbf{x}}_i\boldsymbol{\xi} + (\mathbf{T}_i \otimes \overline{\mathbf{x}}_i\boldsymbol{\delta})). \tag{24}$$

Table 5 provides estimates of both the linear and nonlinear models above using each methodology. Moving from left to right across the table, the first two columns report estimates from POLS and FE linear regressions. The third and fourth columns report estimates from the analogous pooled fractional response (PFR) and CRE estimation. Columns five and six report the estimates from correlated random effects estimation accounting for unbalancedness without (CREU) and with (CREU1) interactions between time averages and time-observation indicators.

It is unclear how we should compute standard errors in this application. The observations are structured as school buildings appearing (or sometimes not appearing) over time, making school-level clustering an obvious choice. However, the policy variation used to instrument for spending in the next section occurs at the district level, a situation similar to the simple situation studied in Abadie et al. (2022), where a policy intervention is correlated within clusters. The Abadie et al. (2022) paper does not cover the panel data setting or instrumental variables, but it seems reasonable to conclude here that clustering at the district level is either correct or somewhat conservative. Consequently, we present school-building-clustered standard errors in parentheses and district-clustered standard errors in brackets. Unless stated otherwise, we use the generally more conservative, district-clustered standard errors for inference in our discussion.

Across all estimators, the estimated effect of expenditures on fourth-graders' achievement in math is positive. However, accounting for unobserved heterogeneity leads to a decrease in the estimated effect of expenditures on achievement. From the first row of column one using POLS, a 10% increase in average spending leads to an 0.84 percentage point (p-value $< 0.001$) increase in fourth-grade math pass rates. From column 2, the fixed-effects estimated effect of the same increase in spending is 0.72 percentage points (p-value $= 0.034$). Using fully-robust, district-clustered standard errors, the estimated effect is statistically significantly positive using either linear approach.

PFR probit estimation yields very similar results to those from POLS. Using this nonlinear approach, a 10% increase in average spending leads to an 0.87 percentage point (p-value $< 0.001$) increase in fourth-grade math pass rates. However, including time averages in the PFR probit causes

the estimated average effect of spending to fall to 0.49 percentage points (p-value = 0.136) under the CRE approach.

As shown in the fifth and sixth columns, accounting for panel imbalance does little to change the CRE estimates. The estimated average effect of spending remains between 0.48 and 0.5 percentage points with p-values between 0.13 and 0.15. The similarity of these CREU estimated coefficients to those estimated ignoring the selection of time-observations suggests that the unbalancedness of the panel is not driving the estimates.

Table 6 provides the results from Wald tests between nested models. We cluster the data at the school-building level for inference in the top two rows, while we cluster the data at the district level to construct the chi-squared test statistics in the bottom two rows. We test the constraints that the nine time averages have zero effect on 4th-grade math pass rates in the linear model in the first column and the nonlinear model in the second column. In both cases with either clustering, we reject the hypothesis that the coefficient estimates on the time averages are zero.

In the third column, we test the coefficient estimates on the four indicators for the number of time-observations. We reject the null hypothesis that all four coefficients are zero at the 5% confidence level when we cluster the data at the school-level. However, clustering at the district-level, we fail to reject the null with a chi-squared test statistic of only 6.2 corresponding to a p-value of 0.18. However, when testing the 30 interactions between indicators for the number of time-observations and the time averages, we reject the null hypothesis that the coefficient on each is zero.[3] In summary, there is significant evidence of unobserved heterogeneity across schools. Further, while the bulk of the evidence points to significant relationships between panel imbalance and fourth-grade math pass rates, the estimated effects of spending on those pass rates remains robust to controlling for panel imbalance.

## 4.2 Allowing spending to be endogenous

The FE and CRE results are robust to unobserved heterogeneity that may be correlated with spending and math pass rates. However, these results are susceptible to the criticism that idiosyncratic unobserved shocks may impact spending and math pass rates. We use instrumental variables

---

[3]Note that due to collinearity, interactions between the indicator for 5 time-observations and time averages for indicators for years 1995, 1996, 1997, and 1998 are omitted as is the interaction between the indicator for four time-observations and the time average of the indicator for year 1997 and 1998.

to address the possible endogeneity of spending and violations of the strict exogeneity assumption.

Michigan's Proposal A equalized revenue to districts according to a non-smooth function determined by district spending in 1994. Consequently, we use the log of the foundation grants set forth by Proposal A ($lfound_{it}$) to instrument for spending ($lrexppp_{it}$). Further, we control for the log of real per-pupil expenditures in 1994 ($lrexppp_{i94}$) to capture this initial heterogeneity in spending. Thus we model cumulative spending according to equation (25) below.

$$
\begin{aligned}
log(avgrexp_{it}) =& \eta_t + \pi_1 lfound_{it} + \pi_2 lrexppp_{i94} + \pi_3 lunch_{it} + \pi_4 lunch_{it}^2 \\
& + \pi_5 log(enroll_{it}) + \pi_6 log(enroll_{it})^2 + c_i + v_{1it}
\end{aligned}
\tag{25}
$$

As a benchmark, we estimate equation (19) using pooled two-stage least squares (P2SLS) with the fitted values coming from POLS estimation of equation (25). We add the residuals from POLS estimation of equation (25), $\widehat{v}_{1it}$, to the pooled fractional probit model to accommodate the nonlinear functional form of fourth-grade math pass rates. This fractional response probit control function approach appears as PFR CF in Table 7.

In order for P2SLS or PFR CF to produces consistent estimates of the causal effect of spending on math pass rates, the instruments must be predictive of $log(avgrexp_{it})$, and the excluded instrument must be uncorrelated with both the unobserved heterogeneity, $c_i$, and the idiosyncratic error term, $e_{it}$, from equation (19). The F-statistic on $lfound_{it}$ from the first stage POLS regression is 180.95 demonstrating that the state foundation grants are indeed predictive of spending.

The second condition is untestable. However, we make a less restrictive exclusion restriction by directly addressing the unobserved heterogeneity, $c_i$. We treat this unobserved heterogeneity by modeling it as a function of the building-level time averages, and including them as regressors as in Chamberlain (1980) and Mundlak (1978). Equation (26) reflects the first stage of this approach.

$$
\begin{aligned}
log(avgrexp_{it}) =& \eta_t + \pi_1 lfound_{it} + \pi_2 lrexppp_{i94} + \pi_3 lunch_{it} + \pi_4 lunch_{it}^2 \\
& + \pi_5 log(enroll_{it}) + \pi_6 log(enroll_{it})^2 + \pi_7 \overline{lunch}_i + \pi_8 \overline{lunch_i^2} + \pi_9 \overline{log(enroll_i)} \\
& + \pi_{10} \overline{log(enroll_i)^2} + \pi_{11} \overline{y96}_i + \pi_{12} \overline{y97}_i + \pi_{13} \overline{y98}_i + v_{1it}.
\end{aligned}
\tag{26}
$$

Note that due to the unbalancedness of the data, we also include time averages of the year indicators – $\overline{y96}_i$, $\overline{y97}_i$, and $\overline{y98}_i$.

In linear specifications, the fitted values of $log(avgrexp_{it})$ from POLS estimation of equation (26) are used for estimating the second stage. This procedure is akin to fixed effects instrumental variables (FEIV) except that we use the base expenditures ($lrexppp_{i94}$) to proxy for time-invariant spending as opposed to the time average of spending.

In nonlinear specifications, we incorporate the estimated residuals, $\widehat{v}_{1it}$, from the same first-stage regression into our CRE fractional response probit. This correlated random effects control function (CRE CF) approach allows us to handle the endogeneity of spending while accommodating the nonlinear functional form.

As with the single equation model, we handle the possibly endogenous unbalancedness of the panel by incorporating the number of time-observations into estimation of the model. We do this by first including indicators for the number of time-observations to the CRE CF in both estimation stages in what we term a correlated random effects unbalancedness control function (CREU CF) approach. Secondly, we incorporate interactions between time averages (and $lrexppp_{i94}$) and the number of time-observations to more fully account for the unbalancedness of the panel in what we term CREU1 CF.

We also perform several robustness checks. First, we incorporate the interactions between time averages of each covariate and indicators for the number of time-observations of each school in the linear fixed effects model. We term this approach FEIVU, which serves as the linear analogous methodology to the CREU1 CF approach. Second, we treat the time-constant heterogeneity in schools by demeaning each covariate in the first stage to generate the residuals used in each of the three control function approaches. Naturally, the differencing eliminates time-constant covariates, such as the number of time-observations and the level of expenditures in 1994, from the first stage estimation. These three approaches (CRE FECF, CREU FECF, CREU1 FECF) appear in the last three columns of Table 7.

Across all instrumental variables approaches the point estimates range from 0.169 to 0.25. These findings all lie within the P2SLS confidence interval reported in Papke (2005). From the first column of Table 7, the P2SLS estimated average partial effect (APE) of a 10% increase in $log(avgrexp_{it})$ is a 2.07 percentage point increase in fourth-grade math pass-rates (p-value = 0.013). Once we instrument for spending, modeling the unobserved heterogeneity does little to change the point estimates in the linear model. FEIV estimates of the effect of the same spending increase math pass

rates by 1.94 percentage points (p-value = 0.021). Accounting for the unbalancedness of the panel drops the point estimate of a 10% increase in spending to 1.69 percentage points (p-value = 0.069), though it remains economically significant.

The nonlinear estimators find slightly larger effects of spending than do the linear instrumental variables approaches. The instrumental variables linear approaches provide estimates of the APE of spending that range from 0.169 to 0.207. Still, the nonlinear estimated APEs lie well within the 95% confidence intervals of the linear estimates. The standard fractional response probit estimates a 10% increase in spending leads to a 2.4 percentage point increase in math pass rates on average. Including time averages using CRE CF decreases the estimated APE of spending to a 2.23 percentage point increase in fourth-grade math pass rates (p-value < 0.001). Adjusting for panel imbalance leads to smaller effects though these differences are far from statistically significant. The estimated APE of $log(avgrexp_{it})$ on pass rates is 0.5% (or 1% of a CREU CF standard error) smaller when including time-observation indicators in the CREU CF approach than when using the more standard CRE CF approach. Including interactions between indicators for the number of time-observations and the time averages of covariates makes a somewhat bigger difference, though the economic and statistical interpretation of the effects remains similar. The estimated APE of spending using CREU1 CF is 9.5% (or 19.6% of a CREU1 CF standard error) smaller than the CRE CF estimate.

Furthermore, there is little loss of efficiency when accounting for panel imbalance. Due to the estimation of the residuals in the first stage, we cluster-bootstrap the standard errors over 500 repetitions to account for the estimation error. When clustering at the building level, the standard errors are identical to three decimal places between CRE CF and CREU CF.[4] There is a larger difference when incorporating interactions between time averages and the number of time-observations. The building-clustered standard errors are approximately 16% larger when applying CREU1 CF as opposed to CREU CF or CRE CF.

The APEs from approaches using fixed effects estimation in the first stage range from 0.212 to 0.25, making them only slightly larger in magnitude than those previously discussed. The primary difference in these estimates is in their precision. Demeaning the foundation grants in the first-stage fixed-effects estimation rather than incorporating base-period expenditures in the first-stage

---

[4]When clustering at the district level the standard errors in columns 4 through 7 are about 30% larger than those from CRE CF.

increases the magnitude of the district-cluster-bootstrapped standard errors by roughly a 20 to 30%.

Beyond directly addressing the potential endogeneity of spending, the test statistic on $\widehat{v}_{1it}$ from these control function approaches conveniently provides evidence regarding the prevalence of endogeneity of $log(avgrexp_{it})$. The coefficient estimate on $\widehat{v}_{1it}$ is -0.194 using CREU CF with a standard error of 0.097 from 500 district-level-cluster-bootstrap replications. Across the primary specifications, using district-level clustering, the p-values range from 0.023 (using the PFR CF approach) to 0.096 (using the CREU1 CF approach).[5] These results provide some evidence against the hypothesis that spending is strictly exogenous.[6] The estimated APEs from the instrumental variables approaches in Table 7 are much larger than the estimated APEs in Table 5 assuming strict exogeneity of spending. For instance, the 95% confidence interval from the CREU CF approach ranges from 0.102 to 0.342, excluding the 0.049 APE of spending estimated using CREU.

Table 8 provides the results from Wald tests between nested two-stage models. Again, we cluster the data at the school-building level for inference in the top two rows, while we cluster the data at the district level to construct the chi-squared test statistics in the bottom two rows. We test the constraints that the seven time averages have zero effect on 4th-grade math pass rates in the linear model in the first column, in the standard nonlinear model in the third column, and in the nonlinear model with first-stage demeaning in the sixth column. Despite the closeness of the estimated APEs, in all three cases with either clustering, we reject the hypothesis that the coefficient estimates on the time averages are zero with p-values less than 0.001 in each case.

In the fourth and seventh columns of Table 8, we test the coefficient estimates on the four indicators for number of time-observations. Column 4 depicts test statistics from our primary CREU CF specification while column 7 depicts test statistics from CREU CF when we demean covariates in the first-stage estimation. With either approach to the first stage, and either level of clustering, we are unable to reject the hypothesis that the coefficient on the four time-observation indicators are zero at the 95% level. The relevant Chi-squared statistics range from just 2.4 to 8.2.

When testing the 27 interactions between indicators for the number of time-observations and

---

[5]Using building-level-clustered standard errors across all specifications, the t-statistics on the residuals range from -2.289 (using the CREU1 FECF approach) to -3.174 (using the FR CF approach).

[6]The district-level t-statistics are much smaller when using fixed effects estimation in the first stage. Comparing CRE CF to CRE FECF the standard errors on the estimated coefficients on both $log(avgrexp_{it})$ and $\widehat{v}_{1it}$ increase when using the fixed effects residuals. The pattern continues when treating the unbalancedness.

the time averages (and $lrexppp_{i94}$), however, the results reject the null hypothesis.[7] Both in linear and nonlinear approaches, regardless of the level of clustering, we reject the null hypothesis of zero coefficients on the interactions between indicators for the number of time-observations and the time averages with p-values consistently smaller than 0.001. While these Wald tests inform whether or not the unobserved heterogeneity and unbalancedness affect test scores conditional on the covariates and foundation grants, the stability of the estimated APEs of spending is reassuring. The effects of spending on fourth-grade math pass rates are not driven by the unbalancedness of the panel.

## 5    Discussion

This paper considers estimation of nonlinear panel data models when the panel is unbalanced in the presence of endogeneity. We allow the selection of time observations to be correlated with both unobserved heterogeneity as well as the explanatory variables. We take a correlated random effects approach and model the unobserved heterogeneity while controlling for selection of time observations. We incorporate a control function approach to handle endogeneity of explanatory variables. In estimating average partial effects we adopt quasi-maximum likelihood estimation such that consistency does not require knowledge of the specific distribution. Our approach is straightforward to implement with standard statistical software and may be used when the outcome is binary, fractional response, or otherwise bounded with known upper and lower bounds. As cases of unbalanced panels are common in many applied fields of economics as well as in other disciplines such as quantitative sociology and political science, there is wide potential for application across the social sciences. The approach is easily extended to other nonlinear panel data models, such as ordered probit and Tobit. Pooled (quasi-) MLE can be applied by combining the CRE and control function approach we propose here.

In estimating the effect of school spending on fourth-grade pass rates on state mathematics exams, we see significant unbalancedness in the underlying data. Estimation is additionally complicated by likely unobserved heterogeneity across schools and the potential of contemporaneous endogeneity of school spending. Indeed, we find evidence supporting the existence of both. Using instrumental variables to identify the effect of spending off of plausibly exogenous changes in the

---

[7]Note that due to collinearity, interactions between the indicator for five time-observations and time averages for indicators for years 1996, 1997, and 1998 are omitted as are the interactions between the indicator for four time-observations and the time averages of the indicators for years 1996 and 1997.

funding structure is consequential. Whereas we estimate that a 10 percent increase in spending leads to a 0.5 percentage point increase in pass rates when we ignore the potential of endogeneity in school spending decisions, we estimate the same spending change to increase pass rates by around 2 percentage points with a variety of instrumental variables approaches.

Once we address the contemporaneous endogeneity of school spending, we find our results to be quite robust. Despite our Wald tests rejecting the null of no unobserved fixed heterogeneity, our estimates remain relatively stable whether or not we include time averages of covariates. Further, though the fact that buildings with missing time observations on average have lower spending and lower pass rates than those with data present for each time period, our estimated effects remain stable regardless of our approach to address the unbalancedness of our panel. The robustness of these estimates provide further evidence of the positive effects of school spending on the students' academic achievement. This result has found additional support in recent work, such as Jackson et al. (2016); Hyman (2017); and Lafortune et al. (2018). Each find substantial positive effects of school spending on students' academic achievement – perhaps finally turning the prevailing narrative to more positively depicting the efficacy of expenditures on public schooling.

## Disclosures

## References

Abadie, A., S. Athey, G. W. Imbens, and J. Wooldridge (2022). When should you adjust standard errors for clustering? Technical report. https://arxiv.org/abs/1710.02926.

Alderman, H., J. R. Behrman, H.-P. Kohler, J. A. Maluccio, and C. S. Watkins (1999). *Attrition in longitudinal household survey data: some tests for three developing-country samples*. The World Bank.

Aughinbaugh, A. (2004). The impact of attrition on the children of the nlsy79. *Journal of human resources 39*(2), 536–563.

Blundell, R. and J. L. Powell (2003). Endogeneity in nonparametric and semiparametric regression models.

Chamberlain, G. (1980). Analysis of covariance with qualitative data. *Review of Economic Studies 47*, 225–238.

Chaudhary, L. (2009, February). Education inputs, student performance and school finance reform in Michigan. *Economics of Education Review 28*(1), 90–98.

Chiang, H. (2009). How accountability pressure on failing schools affects student achievement. *Journal of Public Economics 93*(9-10), 1045–1057.

Falaris, E. M. and H. E. Peters (1998). Survey attrition and schooling choices. *Journal of Human Resources*, 531–554.

Fitzgerald, J. M. (2011). Attrition in models of intergenerational links using the psid with extensions to health and to sibling models. *The BE journal of economic analysis & policy 11*(3).

Hyman, J. (2017). Does money matter in the long run? effects of school spending on educational attainment. *American Economic Journal: Economic Policy 9*(4), 256–80.

Jackson, C. K., R. C. Johnson, and C. Persico (2016). The effects of school spending on educational and economic outcomes: Evidence from school finance reforms. *The Quarterly Journal of Economics 131*(1), 157–218.

Joshi, R. and J. M. Wooldridge (2019). Correlated random effects models with endogenous explanatory variables and unbalanced panels. *Annals of Economics and Statistics* (134), 243–268.

Lafortune, J., J. Rothstein, and D. W. Schanzenbach (2018). School finance reform and the distribution of student achievement. *American Economic Journal: Applied Economics 10*(2), 1–26.

Lin, W. and J. M. Wooldridge (2019). Testing and correcting for endogeneity in nonlinear unobserved effects models. In *Panel data econometrics*, pp. 21–43. Elsevier.

Mundlak, Y. (1978). On the pooling of time series and cross-sectional data. *Econometrica 46*, 69–86.

Papke, L. E. (2005, June). The effects of spending on test pass rates: evidence from Michigan. *Journal of Public Economics 89*(5–6), 821–839.

Papke, L. E. (2008, July). The Effects of Changes in Michigan's School Finance System. *Public Finance Review 36*(4), 456–474.

Papke, L. E. and J. M. Wooldridge (2008, July). Panel data methods for fractional response variables with an application to test pass rates. *Journal of Econometrics 145*(1–2), 121–133.

Rivers, D. and Q. H. Vuong (1988). Limited information estimators and exogeneity tests for simultaneous probit models. *Journal of econometrics 39*(3), 347–366.

Roy, J. (2011, February). Impact of School Finance Reform on Resource Equalization and Academic Performance: Evidence from Michigan. *Education Finance and Policy 6*(2), 137–167.

Wooldridge, J. M. (2019). Correlated random effects models with unbalanced panels. *Journal of Econometrics 211*(1), 137–150.

Table 1: Average correlations across simulation repetitions for both correlated fixed effects and correlated random-coefficient data generating processes (DPGs)

| | Correlated Fixed Effect DGP | | | | | Correlated Random Coefficient DGP | | | |
|---|---|---|---|---|---|---|---|---|---|
| | $x_1$ | $c$ | $u_0$ | $T$ | | $x_1$ | $c$ | $u_1$ | $T$ |
| $x_1$ | 1 | | | | $x_1$ | 1 | | | |
| $c$ | 0.3152 | 1 | | | $c$ | 0.3152 | 1 | | |
| $u_0$ | 0.0006 | 0.0009 | 1 | | $u_1$ | 0.3404 | 0.6294 | 1 | |
| $T$ | 0.0941 | 0.2986 | 0.0013 | 1 | $T$ | 0.1119 | 0.2065 | 0.3300 | 1 |

Notes: Average correlations over 500 simulation repetitions. $x_1$ is the primary variable of interest. $c$ represents the unobserved fixed effects, and $u_0$ and $u_1$ are random slopes. $T$ represents the number of time-observations for a given "school." Correlated fixed effect DGP used in simulations appearing in Table 2. Correlated random-coefficient DGP used in simulations appearing in Table 3.

Table 2: Simulation evidence with selection based on unobserved fixed effects

| True mean APEs of x1 over: | | Population | Sample | T =1 | T =2 | T =3 | T =4 | T =5 |
|---|---|---|---|---|---|---|---|---|
| Mean APE | | 0.2979 | 0.2953 | 0.3252 | 0.3163 | 0.3068 | 0.2965 | 0.2851 |

| Estimates | POLS | FE | FEU | PFR | CRE | CREU | CREU1 | CREU2 | CREU3 |
|---|---|---|---|---|---|---|---|---|---|
| Mean APE | 0.3850 | 0.2939 | 0.2947 | 0.3850 | 0.2947 | 0.2947 | 0.2947 | 0.2949 | 0.2951 |
| Mean SE | 0.0112 | 0.0099 | 0.0099 | 0.0101 | 0.0086 | 0.0086 | 0.0086 | 0.0086 | 0.0086 |
| SD | 0.0111 | 0.0098 | 0.0098 | 0.0102 | 0.0087 | 0.0087 | 0.0086 | 0.0089 | 0.0092 |
| Mean SE/SD | 1.0096 | 1.0108 | 1.0055 | 0.9915 | 0.9950 | 0.9924 | 0.9919 | 0.9611 | 0.9298 |

| Partial effects (PEs) from CREU1 at selected deciles | 10 | 30 | 50 | 70 | 90 |
|---|---|---|---|---|---|
| PE at decile | 0.3314 | 0.3145 | 0.3001 | 0.2837 | 0.2576 |
| SD of PE at decile | 0.0104 | 0.0097 | 0.0091 | 0.0082 | 0.0069 |

Notes: Simulated over 500 repetition with 500 individual "buildings" and a mean of 3.83 out of 5 possible time-observations (T). All standard errors (SE) clustered at the building level. FE regressions are estimated as POLS including time averages. PFR includes the same covariates as POLS with a nonlinear functional form (normal GLM). CRE incorporates individual time averages into PFR. CREU incorporates indicators for the number of individual Tobs into CRE. CREU1 adds interactions between Tobs indicators and covariate time averages to CREU. CREU2 incorporates interactions between covariates and the time averages into CREU1. CREU3 adds triple interactions between covariates, their time averages, and indicators for Tobs.

Table 3: Simulation evidence with selection based on correlated heterogeneous slopes

| True mean APEs of x1 over: | | Population | Sample | T =1 | T =2 | T =3 | T =4 | T =5 |
|---|---|---|---|---|---|---|---|---|
| Mean APE | | 0.2902 | 0.2933 | 0.2523 | 0.2670 | 0.2812 | 0.2931 | 0.3034 |

| Estimates | POLS | FE | FEU | PFR | CRE | CREU | CREU1 | CREU2 | CREU3 |
|---|---|---|---|---|---|---|---|---|---|
| Mean APE | 0.3823 | 0.2954 | 0.2941 | 0.3827 | 0.2945 | 0.2946 | 0.2947 | 0.2945 | 0.2937 |
| Mean SE | 0.0109 | 0.0099 | 0.0099 | 0.0099 | 0.0090 | 0.0090 | 0.0090 | 0.0089 | 0.0086 |
| SD | 0.0101 | 0.0094 | 0.0091 | 0.0094 | 0.0083 | 0.0083 | 0.0083 | 0.0085 | 0.0096 |
| Mean SE/SD | 1.0756 | 1.0494 | 1.0840 | 1.0584 | 1.0795 | 1.0798 | 1.0794 | 1.0458 | 0.8941 |

| Partial effects (PEs) from CREU1 at selected deciles | 10 | 30 | 50 | 70 | 90 |
|---|---|---|---|---|---|
| PE at decile | 0.3321 | 0.3147 | 0.3000 | 0.2833 | 0.2568 |
| SD of PE at decile | 0.0103 | 0.0095 | 0.0087 | 0.0078 | 0.0064 |

Notes: Simulated over 500 repetition with 500 individual "buildings" and a mean of 3.82 out of 5 possible time-observations (T). All standard errors (SE) clustered at the building level. FE regressions are estimated as POLS including time averages. PFR includes the same covariates as POLS with a nonlinear functional form (normal GLM). CRE incorporates individual time averages into PFR. CREU incorporates indicators for the number of individual Tobs into CRE. CREU1 adds interactions between Tobs indicators and covariate time averages to CREU. CREU2 incorporates interactions between covariates and the time averages into CREU1. CREU3 adds triple interactions between covariates, their time averages, and indicators for Tobs.

Table 4: Summary statistics by number of time-observations per school

| Number of time-observations per school | 1 | | 2 | | 3 | | 4 | | 5 | | Total | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Pass rate on fourth-grade math test | 0.72 | (0.23) | 0.60 | (0.21) | 0.64 | (0.20) | 0.57 | (0.23) | 0.65 | (0.19) | 0.64 | (0.20) |
| Average real expenditure per-pupil ($) | 3949 | (1181) | 4082 | (1051) | 3968 | (663) | 3875 | (520) | 3875 | (614) | 3897 | (626) |
| Percent FRL eligible | 0.19 | (0.2) | 0.49 | (0.25) | 0.4 | (0.26) | 0.57 | (0.26) | 0.31 | (0.22) | 0.37 | (0.25) |
| Number of enrolled students | 282 | (173) | 366 | (223) | 411 | (153) | 540 | (226) | 398 | (137) | 420 | (165) |
| Number of schools | 54 | | 42 | | 506 | | 259 | | 910 | | 1771 | |

Notes: Sample means with standard deviations appearing in parentheses.

Table 5: APE estimates assuming spending to be strictly exogenous

|  | Linear | | | Fractional Probit | | |
| VARIABLES | POLS | FE | PFR | CRE | CREU | CREU1 |
| --- | --- | --- | --- | --- | --- | --- |
| lavgrexpp | 0.084 | 0.072 | 0.087 | 0.049 | 0.049 | 0.048 |
|  | (0.017) | (0.026) | (0.018) | (0.025) | (0.025) | (0.025) |
|  | [0.023] | [0.034] | [0.025] | [0.033] | [0.033] | [0.033] |
| lunch | -0.447 | -0.067 | -0.439 | -0.071 | -0.070 | -0.070 |
|  | (0.012) | (0.044) | (0.012) | (0.043) | (0.043) | (0.042) |
|  | [0.026] | [0.047] | [0.026] | [0.046] | [0.046] | [0.045] |
| lenrol | -0.015 | -0.022 | -0.014 | -0.019 | -0.019 | -0.019 |
|  | (0.009) | (0.021) | (0.008) | (0.021) | (0.021) | (0.021) |
|  | [0.010] | [0.023] | [0.010] | [0.022] | [0.022] | [0.022] |
|  |  |  |  |  |  |  |
| Observations | 7,242 | 7,242 | 7,242 | 7,242 | 7,242 | 7,242 |

Notes: School-clustered standard errors appear in parentheses. District-clustered standard errors are in brackets. CREU estimation includes indicators for each number of time-observations. CREU1 includes indicators for time-observations as well as interactions between time-observations and time averages of covariates. All regressions include year indicators. Regressions including time averages also include time averages of year indicators.

Table 6: Wald testing between nested single-stage models

| Model Comparison | POLS vs FE | PFR vs CRE | CRE vs CREU | CREU vs CREU1 |
| --- | --- | --- | --- | --- |
| Panel A: School-level clustered standard errors | | | | |
| $\chi^2$ (constraints) | 13.7 (9) | 114.7 (9) | 9.7 (4) | 93.7 (30) |
| Prob $> \chi^2$ | <0.001 | <0.001 | 0.046 | <0.001 |
|  |  |  |  |  |
| Panel B: District-level clustered standard errors | | | | |
| $\chi^2$ (constraints) | 13.2 (9) | 107.2 (9) | 6.2 (4) | 145.1 (30) |
| Prob $> \chi^2$ | <0.001 | <0.001 | 0.184 | <0.001 |
|  |  |  |  |  |
| Variables tested |  |  |  |  |
| time averages | X | X |  |  |
| time-observation |  |  | X |  |
| time-observation interactions |  |  |  | X |

Table 7: APEs allowing spending to be contemporaneously endogenous

| VARIABLES | Two-stage linear model | | | Fractional resonse probit control functions | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | P2SLS | FEIV | FEIVU | PFR CF | CRE CF | CREU CF | CREU1 CF | CRE FECF | CREU FECF | CREU1 FECF |
| lavgrexpp | 0.207 | 0.191 | 0.168 | 0.240 | 0.223 | 0.222 | 0.201 | 0.250 | 0.248 | 0.212 |
| | (0.055) | (0.056) | (0.063) | (0.061) | (0.064) | (0.064) | (0.074) | (0.073) | (0.073) | (0.076) |
| | [0.083] | [0.083] | [0.092] | [0.093] | [0.093] | [0.095] | [0.107] | [0.127] | [0.126] | [0.129] |
| residuals | | | | -0.219 | -0.197 | -0.194 | -0.178 | -0.223 | -0.220 | -0.190 |
| | | | | (0.069) | (0.073) | (0.074) | (0.082) | (0.081) | (0.081) | (0.083) |
| | | | | [0.096] | [0.095] | [0.097] | [0.107] | [0.126] | [0.126] | [0.127] |
| lunch | -0.454 | -0.022 | -0.014 | -0.655 | -0.037 | -0.043 | -0.053 | -0.039 | -0.044 | -0.056 |
| | (0.015) | (0.062) | (0.061) | (0.014) | (0.062) | (0.062) | (0.060) | (0.065) | (0.065) | (0.063) |
| | [0.030] | [0.066] | [0.064] | [0.048] | [0.065] | [0.066] | [0.064] | [0.070] | [0.070] | [0.068] |
| lenrol | -0.005 | -0.037 | 0.025 | -0.197 | 0.018 | 0.035 | -0.009 | 0.010 | 0.023 | -0.032 |
| | (0.012) | (0.033) | (0.034) | (0.011) | (0.033) | (0.033) | (0.034) | (0.030) | (0.030) | (0.030) |
| | [0.015] | [0.042] | [0.042] | [0.015] | [0.042] | [0.045] | [0.046] | [0.036] | [0.036] | [0.035] |
| Observations | 4,853 | 4,853 | 4,853 | 4,853 | 4,853 | 4,853 | 4,853 | 4,853 | 4,853 | 4,853 |

Notes: School-clustered standard errors appear in parentheses. District-clustered standard errors are in brackets. Control function standard errors are from 500 cluster-bootstrap repetitions to handle residuals' estimation error. CREU estimation includes indicators for each number of time-observations. CREU1 includes indicators for time-observations as well as interactions between time-observations and time averages of covariates. All regressions include year indicators. Regressions including time averages also include time average of year indicators.

Table 8: Wald testing between nested two-stage models

| Model Comparison | IV vs FEIV | FEIV vs FEIVU | PFR CF vs CRE CF | CRE CF vs CREU CF | CREU CF vs CREU1 CF | PFR CF vs CRE FECF | CRE FECF vs CREU FECF | CREU CF vs CREU1 FECF |
|---|---|---|---|---|---|---|---|---|
| **Panel A: School level clustered standard errors** | | | | | | | | |
| $\chi^2$ (constraints) | 78.7 (7) | 524.4 (27) | 85.9 (7) | 8.2 (4) | 80.1 (27) | 69.4 (7) | 6.6 (4) | 80.1 (27) |
| Prob > $\chi^2$ | >0.001 | >0.001 | >0.001 | 0.083 | >0.001 | >0.001 | 0.159 | >0.001 |
| **Panel B: District level clustered standard errors** | | | | | | | | |
| $\chi^2$ (constraints) | 73.1 (7) | 220.2 (27) | 55.7 (7) | 2.5 (4) | 59.4 (27) | 50.5 (7) | 2.4 (4) | 58.7 (27) |
| Prob > $\chi^2$ | >0.001 | >0.001 | >0.001 | 0.644 | >0.001 | >0.001 | 0.668 | >0.001 |
| **Variables tested** | | | | | | | | |
| time averages | X | | X | | | X | | |
| time-observation indicators | | | | X | | | X | |
| time-observation interactions | | X | | | X | | | X |