Double Machine Learning for Static Panel Models with Fixed Effects

Paul S. Clarke[†] and Annalivia Polselli[‡]

[†]Institute for Social and Economic Research, University of Essex, Colchester CO4 3SQ, UK. E-mail: pclarke@essex.ac.uk

[‡]Institute for Analytics and Data Science, University of Essex, Colchester CO4 3ZL, UK. E-mail: annalivia.polselli@essex.ac.uk

Summary Recent advances in causal inference have seen the development of methods which make use of the predictive power of machine learning algorithms. In this paper, we develop novel double machine learning (DML) procedures for panel data in which these algorithms are used to approximate high-dimensional and nonlinear nuisance functions of the covariates. Our new procedures are extensions of the well-known correlated random effects, within-group and first-difference estimators from linear to nonlinear panel models, specifically, Robinson (1988)'s partially linear regression model with fixed effects and unspecified nonlinear confounding. Our simulation study assesses the performance of these procedures using different machine learning algorithms. We use our procedures to re-estimate the impact of minimum wage on voting behaviour in the UK. From our results, we recommend the use of first-differencing because it imposes the fewest constraints on the distribution of the fixed effects, and an ensemble learning strategy to ensure optimum estimator accuracy.

Keywords: *CART*, homogeneous treatment effect, hyperparameter tuning, LASSO, random forest.

1. INTRODUCTION

Recent advances in the econometric literature on Machine Learning (ML) use the power of ML algorithms, widely used in data science for solving prediction problems, to enhance existing estimation procedures for treatment and other kinds of causal effect. Notable developments include novel ML algorithms for causal analysis such as Causal Trees by Athey and Imbens (2016), Causal Forests by Wager and Athey (2018) and Generalised Random Forests by Athey et al. (2019). However, the key development, as far as this paper is concerned, is Double/Debiased Machine Learning (DML) by Chernozhukov et al. (2018) wherein ML is used to learn *nuisance functions* with *ex ante* unknown functional forms, and the predicted values of these functions used to construct (orthogonalized) scores for the interest parameters from which consistent and asymptotically normal estimators can be obtained. DML is a very general estimation framework but there are limited examples of its application to panel data, notable examples of which include Chang (2020), Klosin and Vilgalys (2023), and Semenova et al. (2023).

In this paper, we develop and assess novel DML procedures for estimating treatment (or causal) effects from panel data with *fixed effects*. The procedures we propose are extensions of the correlated random effects (CRE), within-group (WG) and first-difference (FD) estimators commonly used for linear models to scenarios where the underlying model is non-linear. Specifically, these are based on an extension of the partially linear regression (PLR) model proposed by Robinson (1988) to panel data through the inclusion of time-varying predictors and unobserved individual heterogeneity (i.e. individual fixed effects).

Our methodological contribution is twofold and complementary to the recent work on causal panel data estimation using DML by Chang (2020) on difference-in-differences, Klosin and Vilgalys (2023) and Semenova et al. (2023) on high-dimensional treatment heterogeneity in model with fixed effects.

A first contribution is that the procedures we develop are not based on *ex ante* taking the nuisance functions or fixed effects to be accurately approximated by high-dimensional sparse functions. Thus, we do not focus solely on the Least Absolute Shrinkage and Selection Operator (LASSO) in the first stage of DML. This is in contrast to Klosin and Vilgalys (2023) and Semenova et al. (2023) who explicitly exploit sparsity to propose alternative procedures. While we do not rule out using LASSO (and indeed do so in both our simulation study and empirical application), our approach requires only minimal assumptions (like Lipschitz continuity) about the nuisance functions. This is in line with Chang (2020) to encourage researchers to adopt estimation strategies which involve the use of different ML algorithms. Such strategies could involve the use of *ensemble* learning such as the stacking of Breiman (1996), the super-learning (a weighted average the best predictions over different choices of ML algorithm) of van der Laan et al. (2007), or selecting the best-performing learner on an application-by-application basis to ensure inference is based on the most accurate predictions. This is important because, in practice, many ML algorithms are difficult to tune and have been shown to perform very differently for different datasets in different contexts.

Second, on a computational note, we use *block-k-fold* cross-fitting, where the entire time series of the sampled unit is allocated to one fold to allow for possible serial correlation within each unit as is common with panel data. In this, we are aligned with the concurrent work by Klosin and Vilgalys (2023) in not needing to rely on the *weak* dependence assumption to deal with dependent data as do Semenova et al. (2023). Conversely, Chang (2020) does not require any additional assumption or change in the splitting algorithms because they use only one observation per unit after treatment, ruling out the presence of any fixed effects.

Our method for panel data models with individual fixed effects is general and particularly relevant for applied researchers. We provide new estimation tools within the existing DML framework for use on panel data. In doing so, we broaden the reach of DML to a large family of empirical problems for which the time dimension must be properly accounted for. Our focus, in the subsequent development, on the homogeneous treatment effect case is also because such models are widely used by applied researchers, but we also show how our procedures can be extended to heterogeneous treatment effects provided that the analyst is prepared to specify a (finite-dimensional) parametric model for this heterogeneity. More widely, we encourage researchers to use the procedures we propose in place of existing ones, or to test the robustness of their results - based on, say, linear models - to non-linearity.

We carry out a simulation study to assess our DML procedures and find large accuracy gains from using DML with flexible learners when the data generating process is highly non-linear (specifically, one involving a non-linear discontinuous function of the regressors); and, while ordinary least squares (OLS) *can* outperform ML when the data generating process is linear (as expected) and for some non-linear processes (specifically, for a smooth process that excludes interactions), it is not robust because the analyst never knows the form of the data generating process *ex ante*. Finally, we apply our new procedures to observational panel data by re-analyzing part of the study by Fazio and Reggiani (2023) on the effect of the introduction National Minimum Wage (NMW) in the United Kingdom (UK) on voting for conservative parties. We find evidence that the first-difference procedures are the most robust and perform well, and confirm the conclusion from our simulation study that ensemble learning strategies are crucial to obtaining robust results.

The remainder of the paper is structured as follows. Section 2 provides an overview of the literature and places our novel contribution within it. Section 3 motivates the partially linear panel regression model and the causal assumptions which must hold to ensure the target parameter can be interpreted as a causal effect. Section 4 introduces the two approaches we take to handle the fixed-effects problem. Section 5 formally introduces the DML estimation procedures. Section 6 briefly discusses the Monte Carlo simulation results. Finally, Section 7 illustrates an empirical application of the procedure and we make concluding remarks in Section 8.

2. RELATED LITERATURE

There is a growing body of econometrics literature on using ML for causal inference. One strand focuses on building or modifying existing learners to consistently estimate and make inferences about causal effects (e.g., Athey and Imbens, 2016; Wager and Athey, 2018; Athey et al., 2019; Lechner and Mareckova, 2022; Di Francesco, 2025). Another strand focuses on incorporating ML into traditional statistical estimators (e.g., least squares, generalised method of moments, maximum likelihood) to estimate causal effects more accurately (e.g., Belloni et al., 2014, 2016; Chernozhukov et al., 2018; Chang, 2020; Chernozhukov et al., 2022; Bia et al., 2023). This paper falls into the second strand.

Much of the ML literature in econometrics, including that for causal estimation, is built around penalized regression methods like LASSO, where the onus is on the analyst to specify a sufficiently rich data dictionary for the problem at hand. For example, Belloni et al. (2016) proposed two-step post-cluster LASSO procedures for panel data with additive individual-specific heterogeneity in which potential control variables are selected using LASSO at stage one followed by the estimation of the homogeneous treatment effect at stage two.

Two recent papers propose estimation procedures based on LASSO for non-linear panel models. Both focus on different causal targets to ours: in the context of continuous treatments, Klosin and Vilgalys (2023) proposed an estimator for the average partial effect based on a PLR model suitable for static panels, where a first-difference data transformation is used to handle omitted fixed time-invariant confounding; and Semenova et al. (2023) propose a procedure for CATE estimation based on a PLR model for dynamic panels. Both of these procedures rely crucially on the key nuisance functions being well approximated by high-dimensional sparse functions: without assuming this *ex-ante*, their procedures do not follow. We propose alternative procedures that make no such *ex-ante* assumptions, and set these up within the DML framework developed by Chernozhukov et al. (2018) who considered the PLR model at length but did not provide any panel examples.

The generality of our approach allows us to move beyond LASSO to other types of ML algorithm. LASSO is powerful but it requires the analyst to pre-specify a sufficiently rich data dictionary, and the memory required to store these in panels with many waves could be prohibitively large such that alternative learners are preferable. However, by staying general, we rule out solutions to the fixed-effects problem like that proposed by Semenova et al. (2023) who incorporate the fixed effects into the LASSO data dictionary

under a weak sparsity assumption. While recent work by Kolesár et al. (2023) argues that, in the absence of prior knowledge to the contrary, this assumption is unlikely to hold, Semenova et al. (2023) avoid this issue by decomposing the overall fixed effect into a fixed part, involving the individual-specific means of the (fixed) predictors, and a fixed-effect residual which *can* be incorporated into the dictionary under a weak sparsity assumption.

By not formulating the learning problem in terms of LASSO, we must consider alternative approaches so, instead, propose three estimation procedures that obviate entirely the requirement to model the fixed effects. The first two are based on transforming the outcome to induce a reduced-form model that does not depend on the fixed effects. Because we are agnostic, our first-difference procedure differs from that proposed by Klosin and Vilgalys (2023) in ways we explain further on. We also consider a procedure based on a correlated random effects model similar to that used by Wooldridge and Zhu (2020) for limited dependent variable panel models to convert the fixed effects into random effects.

Finally, we note that our work also fits into the literature that leverages the power of ML for causal analysis and policy evaluation. The value added of doubly-robust procedures has been explored in applied works by, for example, Knaus (2022), Bach et al. (2023), Strittmatter (2023), Baiardi and Naghi (2024a,b). Because panel data are widely used in applied analyses, our proposed procedures for panel data models have the potential to attract the interest of applied researchers from various fields broadening the applicability of DML.

3. THE PARTIALLY LINEAR PANEL REGRESSION MODEL

3.1. Notation

Suppose the panel study collected information on each of N individuals at each of the t time periods, or waves. Let $\{(Y_{it}, D_{it}, \mathbf{X}_{it}) : t = 1, \ldots, T\}_{i=1}^{N}$ be N independent and identically distributed (*iid*) random vectors for individual i across all T waves, where $Y_{it} \in \mathcal{Y}$ is the outcome variable, $D_{it} \in \mathcal{D}$ a continuous or binary treatment variable (or intervention), and $\mathbf{X}_{it} = (X_{it,1}, \ldots, X_{it,p})' \in \mathcal{X}$ a $p \times 1$ vector of control (pre-determined) variables, usually including a constant term, able to capture time-varying confounding induced by non-random treatment selection.¹ We denote the realizations of these random variables by $\{(y_{it}, d_{it}, \mathbf{x}_{it})\}$, respectively. For continuous $D_{it} \in \mathcal{D} \subset \mathbb{R}$, if $d_{it} \geq 0$ then a dose-response relationship is presumed to hold with $d_{it} = 0$ indicating null treatment; otherwise, D_{it} is taken to be centered around its mean μ_D such that $D_{it} \equiv D_{it} - \mu_D$. For binary $D_{it} \in \{0,1\}, d_{it} = 0$ is taken to indicate the absence and $d_{it} = 1$ the presence of treatment.

The non-linear model we propose is a simple extension of the partially linear model of Robinson (1988) to panel data. We then use a potential outcomes causal framework (e.g. Rubin (1974)) to set out the assumptions under which its parameters can be interpreted causally. As pointed out by Lechner (2015), this allows us to interpret the resulting partially linear panel model as reduced-form without relying on specific parametric model for the data generating process. To this end, further define the set $Y_{it}(.) = \{Y_{it}(d) :$ $d \in \mathcal{D}\}$ containing all potential outcomes for individual *i* at wave *t*, where $Y_{it}(d)$ is the realization of the outcome for individual *i* at wave *t* were the treatment level set

¹Throughout, we use \mathbf{v}' to indicate the matrix transpose of arbitrary vector \mathbf{v} and take vectors to be column vectors unless explicitly stated to the contrary.

to d, with one potential outcome for every possible value the treatment could take. The realizations of the wave t potential outcomes are taken to occur before treatment selection at wave t, and are linked to the observed outcome by the consistency assumption that $Y_{it}(d_{it}) = Y_{it}$ with the others latent counterfactuals.² In the interval preceding wave t, it is also presumed that the realization of time-varying predictor X_{it} precedes that of (Y_{it}, D_{it}) .

Finally, define the sets ξ_i of time-invariant heterogeneity terms influencing (Y_{it}, D_{it}) , and $L_{t-1}(W_i) = \{W_{i1}, \ldots, W_{it-1}\}$ the lags of a generic random variable W_{it} available at wave t such that $L_0(W_i) \equiv \emptyset$.

3.2. Model and Assumptions

We extend the partially linear model proposed by Robinson (1988) to panel data by introducing fixed effect α_i^* to give the partially linear panel regression (PLPR) model

$$Y_{it} = D_{it}\theta_0 + g_1(\boldsymbol{X}_{it}) + \alpha_i^* + U_{it}, \qquad (3.1)$$

where g_1 is a non-linear nuisance function of \mathbf{X}_{it} and $E[U_{it}|D_{it}, \mathbf{X}_{it}, \alpha_i^*] = 0$ but $E[\alpha_i^* | D_{it}, \mathbf{X}_{it}] \neq 0$. The target parameter θ_0 is the average partial effect of continuous D_{it} such that $d\theta_0 = E[Y_{it}(d) - Y_{it}(0)]$, or the average treatment effect (ATE) $E[Y_{it}(1) - Y_{it}(0)]$ for binary treatments.

Following Chernozhukov et al. (2018), we focus on the partialled-out PLPR (PO-PLPR) model

$$Y_{it} = V_{it}\theta_0 + l_1(\boldsymbol{X}_{it}) + \alpha_i + U_{it}, \qquad (3.2)$$

$$V_{it} = D_{it} - m_1(\boldsymbol{X}_{it}) - \gamma_i, \qquad (3.3)$$

where l_1 and m_1 are nuisance functions, α_i is a fixed effect, $E[U_{it} | V_{it}, X_{it}, \alpha_i] = 0$, and V_{it} the residual of a non-linear additive noise treatment model depending on fixed effect γ_i and satisfying $E[V_{it}|X_{it}, \gamma_i] = 0$. The focus will be on the PO-PLPR rather than the PLPR model because, as shown below, l_1 is a conditional expectation over Y_{it} whereas g_1 is a conditional expectation over counterfactual $Y_{it}(0)$ that must generally be learnt iteratively. Moreover, the orthogonalized score on which DML estimation of θ_0 is based involves V_{it} whether it is derived under PLPR or PO-PLPR (see Section 5).

We now detail the underlying assumptions required for θ_0 to be feasibly estimable with a causal interpretation. Note that some of these assumptions are stated in terms of stochastic conditional independence \perp to simplify the development. Recalling that ξ_i represents the omitted time-invariant confounding variables, the first two assumptions can be stated as follows:

Assumption 3.1. (No feedback to predictors) $X_{it} \perp L_{t-1}(Y_i, D_i) \mid L_{t-1}(X_i), \xi_i$.

Assumption 3.2. (Static panel) $Y_{it}, D_{it} \perp L_{t-1}(Y_i, X_i, D_i) \mid X_{it}, \xi_i$.

Assumptions 3.1-3.2 together ensure that the panel is static and that any observed lag dependence is due to non-causal autocorrelation. Specifically, these assumptions ensure

²The stable unit treatment value (SUTVA) assumption, that $Y_{it}(d)$ does not depend on the treatment assignments of any other individual, is implicitly taken to hold.

that the joint distribution of outcomes and treatments given the time-varying predictors and omitted time-invariant influences satisfies

$$p\{Y_{i1}, D_{i1}, \dots, Y_{iT}, D_{iT} | L_T(\boldsymbol{X}_i), \xi_i\} = \prod_{t=1}^T p(Y_{it}, D_{it} | \boldsymbol{X}_{it}, \xi_i),$$

where p(.) denotes a density function for the (conditional or joint) distribution indicated by its arguments. Moreover, there is no need to model the distribution of the time-varying predictors, and the *initial conditions problem*, which would arise were the panel study to have started after the joint process began, is avoided.

To give θ_0 a causal interpretation, we require that the data generating process for the potential outcomes and treatment satisfies the following conditions:

Assumption 3.3. (Selection on observables and omitted time-invariant variables) $Y_{it}(.) \perp D_{it} \mid \mathbf{X}_{it}, \xi_i$.

Assumption 3.4. (Homogeneity and linearity of the treatment effect) $E[Y_{it}(d) - Y_{it}(0)|\mathbf{X}_{it}, \xi_i] = d\theta_0.$

Assumption 3.3 states that treatment selection at wave t is (strongly) ignorable given X_{it} and latent ξ_i . Assumption 3.4 requires that $Y_{it}(d) - Y_{it}(0)$ is constant or varies between individuals independently of X_{it} and ξ_i . Alternatively, under the weakly ignorable assumption that only $Y_{it}(0) \perp D_{it} \mid X_{it}, \xi_i$, target parameter $d\theta_0 = E[Y_{it}(d) - Y_{it}(0) \mid D_{it} = d]$. However, we assume throughout that the strong version holds.

The final assumption is that the combined effect of the confounding variables is additively separable into that of the observed confounding variables and that of the omitted time-invariant confounding variables.

Assumption 3.5. (Additive Separability)

(a)
$$E[Y_{it}(0) \mid \mathbf{X}_{it}, \xi_i] = g_1(\mathbf{X}_{it}) + \alpha_i^*$$
 where $\alpha_i^* = \alpha^*(\xi_i)$,
(b) $E[D_{it} \mid \mathbf{X}_{it}, \xi_i] = m_1(\mathbf{X}_{it}) + \gamma_i$ where $\gamma_i = \gamma(\xi_i)$.

Assumptions 3.3 and 3.4 imply that $E[Y_{it} \mid D_{it}, \mathbf{X}_{it}, \xi_i] = E[Y_{it}(0) \mid \mathbf{X}_{it}, \xi_i] + D_{it}\theta_0$. Hence, Assumption 3.5(a) leads instantly to PLPR model (3.1). Moreover, Assumptions 3.3 and 3.4 also imply that $E[Y_{it}(0) \mid \mathbf{X}_{it}, \xi_i] = E[Y_{it} \mid \mathbf{X}_{it}, \xi_i] - E[D_{it} \mid \mathbf{X}_{it}, \xi]\theta_0$, from which it follows from both Assumptions 3.5(a) and 3.5(b) that PO-PLPR model (3.2)-(3.3) holds because $E[Y_{it} \mid \mathbf{X}_{it}, \xi_i] = l_1(\mathbf{X}_{it}) + \alpha_i$, where $l_1(\mathbf{X}_{it}) = g_1(\mathbf{X}_{it}) + m_1(\mathbf{X}_{it})\theta_0$ and $\alpha_i = \alpha_i^* + \gamma_i\theta_0$.

In the description of each approach, below, we focus on procedures for the homogenous treatment effects case. However, all of these procedures can be extended to heterogeneous treatment effects provided the analyst is prepared to specify a (finite-dimensional) parametric heterogeneity model. A summary of how this is done can be found in Section 4.3.

4. FIXED-EFFECTS ESTIMATION

We propose two approaches for estimating θ_0 from PO-PLPR model (3.2)-(3.3). The presence of the time-invariant fixed effects means that only $l_1(\mathbf{X}_{it}) + E[\alpha_i | \mathbf{X}_{it}]$ and

 $m_1(\mathbf{X}_{it}) + E[\gamma_i|\mathbf{X}_{it}]$ can be learnt from the observed data directly through $E[Y_{it} | \mathbf{X}_{it}]$ and $E[D_{it} | \mathbf{X}_{it}]$, respectively, so the first challenge is to remove the fixed effects just as it is in the linear case. As such, the approaches we propose are adaptations of existing techniques for *linear* panel data models. The first is based on correlated random effects (CRE) and the second on two varieties of data transformation: first-difference (FD) and within-group (WG).³

Crucially, to understand our contribution, it is important to note that the use of CRE or data transformations does not in itself lead to a consistent estimator of θ_0 . There are further challenges, created by the non-linearity of the nuisance functions l_1 and m_1 , to be overcome. Specifically, we do not solve these challanges by a priori taking l_1 and m_1 to be accurately approximated by high-dimensional sparse functions of the form $l_1(\mathbf{X}_{it}) \approx \mathbf{B}'_{it} \boldsymbol{\psi}$ and $m_1(\mathbf{X}_{it}) \approx \mathbf{B}'_{it} \boldsymbol{\phi}$, where vector $\mathbf{B}_{it} = \mathbf{B}(\mathbf{X}_{it})$ is an analyst-specified high-dimensional data dictionary of linear and non-linear terms, $\boldsymbol{\psi}$ and $\boldsymbol{\phi}$ are vectors of parameters with dim $(\boldsymbol{\psi}) >> ||\boldsymbol{\psi}||_0$ and dim $(\boldsymbol{\phi}) >> ||\boldsymbol{\phi}||_0$, and $||\mathbf{v}||_0$ is the number of non-zero elements in arbitrary vector \mathbf{v} . The FD and CRE procedures respectively developed by Klosin and Vilgalys (2023) and Semenova et al. (2023) rely explicitly on such representations. However, our procedures are suitable for DML based on ensemble learning strategies.

4.1. Correlated Random Effects

Correlated random effects (CRE) models are extensions of fixed effect models in which the fixed effect is replaced by a model for its dependence on the predictor variables. CRE models can depend on any function of $\{\boldsymbol{X}_{it} : t = 1, ..., T\}$ (e.g. Wooldridge and Zhu (2020)) but practice tends to focus on *exchangeable* functions and especially functions of $\overline{\boldsymbol{X}}_i = T^{-1} \sum_{t=1}^T \boldsymbol{X}_{it}$ following Mundlak (1978) who first demonstrated the equivalence of CRE and fixed effects estimators when the effects of \boldsymbol{X}_{it} in the fixed-effects model and the effects of $\overline{\boldsymbol{X}}_i$ in the model for the fixed effect are all linear.

In general, the fixed-effect model must be correctly specified, but a major advantage of the Mundlak device for linear panel models, where the effects of X_{it} are linear, is that least-squares estimators of the coefficients remain consistent even if the linear model for $E[\alpha_i|\overline{X}_i]$ is mis-specified, provided that α_i and \overline{X}_i are correlated. However, this robustness is lost for PO-PLPR model (3.2)-(3.3) because l_1 and m_1 are non-linear. This is also the case for the linear panel model with random coefficients for heterogeneous treatment effects of Wooldridge (2019) and the non-linear probit model of Wooldridge and Zhu (2020).

In the spirit of Wooldridge and Zhu (2020), we derive a CRE model based on PO-PLPR model (3.2)-(3.3). Suppose that the data generating process satisfies Assumptions 3.1-3.5 and the fixed effects follow non-linear additive noise models $\alpha_i = \omega_{\alpha}(\overline{X}_i) + a_i$ and $\gamma_i = \omega_{\gamma}(\overline{X}_i) + c_i$, where a_i and c_i are random effects satisfying $(a_i, c_i) \perp L_T(X_i)$ and

³Under the random effects assumption, where $E[\alpha_i | \mathbf{X}_{it}] = E[\alpha_i] = 0$ and $E[\gamma_i | \mathbf{X}_{it}] = E[\gamma_i] = 0$, l_1 and m_1 can be straightforwardly learnt from the observed data, but we assert that this is unlikely to hold in practice and, hence, causal inference would not be credible. However, were the analyst prepared to make the random effects assumption, we note that Sela and Simonoff (2012) developed an algorithm for using tree-based learners to estimate non-causal partially linear regression models.

 $a_i \perp c_i$. Then the CRE model with nuisance parameters \tilde{l}_1 and \tilde{m}_1 is

$$Y_{it} = V_{it}\theta_0 + l_1(\boldsymbol{X}_{it}, \overline{\boldsymbol{X}}_i) + a_i + U_{it}$$

$$(4.1)$$

$$V_{it} = D_{it} - \widetilde{m}_1(\boldsymbol{X}_{it}, \boldsymbol{X}_i) - c_i, \qquad (4.2)$$

where $E[U_{it} | V_{it}, \mathbf{X}_{it}, \overline{\mathbf{X}}_i, a_i] = E[V_{it} | \mathbf{X}_{it}, \overline{\mathbf{X}}_i, c_i] = 0$, $\tilde{l}_1(\mathbf{X}_{it}, \overline{\mathbf{X}}_i) = l_1(\mathbf{X}_{it}) + \omega_{\alpha}(\overline{\mathbf{X}}_i)$ and $\tilde{m}_1(\mathbf{X}_{it}, \overline{\mathbf{X}}_i) = m_1(\mathbf{X}_{it}) + \omega_{\gamma}(\overline{\mathbf{X}}_i)$. The random effects simply capture autocorrelation between observations on the same individual. Crucially, \tilde{l}_1 and \tilde{m}_1 can be learnt directly from $E[Y_{it} | \mathbf{X}_{it}, \overline{\mathbf{X}}_i]$ and $E[D_{it} | \mathbf{X}_{it}, \overline{\mathbf{X}}_i]$, respectively.⁴

4.2. Data Transformations

The second approach we consider follows more conventional techniques for panel data by transforming the data to remove entirely the fixed effects from the analysis.

Let W_{it} be a generic random variable and Q a panel data transformation operator such that $Q(W_{it}) = Q_t(W_{i1}, \ldots, W_{iT})$ is a function of the random variable at wave t and the remaining realizations for individual i. We consider two such transformations: the WG (within-group) or time-demeaning transformation, i.e., $Q(W_{it}) = W_{it} - \overline{W}_i$, where $\overline{W}_i = T^{-1} \sum_{t=1}^{T} W_{it}$; and the FD (first-difference) transformation, i.e., $Q(W_{it}) = W_{it} - W_{it-1}$ for $t = 2, \ldots, T$.

The reduced-form model for $Q(Y_{it})$ and $Q(V_{it})$ under PO-PLPR model (3.2)-(3.3) is

$$Q(Y_{it}) = Q(V_{it})\theta_0 + Q(l_1(X_{it})) + Q(U_{it})$$
(4.3)

$$Q(V_{it}) = Q(D_{it}) - Q(m_1(X_{it})),$$
(4.4)

which does not depend on fixed effects α_i and γ_i because $Q(\alpha_i) = Q(\gamma_i) = 0$.

The challenge of learning the transformed nuisance functions $Q(l_1(\mathbf{X}_{it}))$ and $Q(m_1(\mathbf{X}_{it}))$ is complicated by the non-linearity of l_1 and m_1 because, as set out previously, we do not a priori assume that the nuisance functions admit high-dimensional sparse representations of l_1 and m_1 where $Q(l_1(\mathbf{X}_{it}) \approx Q(\mathbf{B}'_{it})\boldsymbol{\psi})$ and $Q(m_1(\mathbf{X}_{it})) \approx Q(\mathbf{B}'_{it})\boldsymbol{\phi}$ could be learnt directly as per the FD procedure proposed by Klosin and Vilgalys (2023).

4.3. Heterogeneity

If Assumption 3.4 is implausible or the inferential target of the analysis is heterogeneity of the treatment effects itself then further modelling is required. Completely relaxing Assumption 3.4 results in $\theta_0 d = E[Y_{it}(d) - Y_{it}(0) | \mathbf{X}_{it}, \xi_i]$, that is, conditional average partial effects or conditional average treatment effects (CATE) dependent on all of the confounding variables. Focussing on the binary treatment case for simplicity, extending Assumption 3.5(a) so that $E[Y_{it}(1) | \mathbf{X}_{it}, \xi_i]$ is also additively separable would lead to an additively separable CATE $\theta_0 = h_0(\mathbf{X}_{it}) + \delta_i$, where $\delta_i = \delta(\xi_i)$ is what Wooldridge (2019) refers to as a random coefficient.

CRE model (4.1)-(4.2) can be specified in terms of CATEs if Assumption 3.4 fails.

⁴Semenova et al. (2023) also use a CRE model in their formulation for dynamic models. In our static panel case with only fixed X_{it} predictors, this would be equivalent to specifying $l_1(X_{it}) + \alpha_i = \psi_0 \overline{X}_i + \alpha_i$ with $(\alpha_1, \ldots, \alpha_N)$ assumed to be a *weakly sparse* vector following a Mundlak-style model which can be estimated using LASSO. Alternative models for weakly sparse fixed effects have been developed by Kock (2016) and Kock and Tang (2019) which can be estimated using *de-sparsified* LASSO.

Temporarily discarding vector notation, if $E[Y_{it}(1) - Y_{it}(0) | X_{it}, \xi_i] = h_0(X_{it})$ then our procedure allows parametric model $h_0(X_{it}) = S_{it}\theta_0$, where $S_{it} = S(X_{it})$ captures treatment effect heterogeneity though interaction $S_{it}V_{it}$ replacing V_{it} in (4.1). If the heterogeneity also depends on ξ_i then, following Wooldridge (2019), the analyst can proceed by additionally specifying a further linear model for $E[\delta_i | \overline{X}_i]$ so that S_{it} is extended to $S_{it} = S(X_{it}, \overline{X}_i)$ with no interactions between X_{it} and \overline{X}_i included. As noted by Wooldridge (2019), this model must be correctly specified to induce a consistent estimator, despite having the appearance of the Mundlak device.

Finally, for the transformation models, the failure of Assumption 3.4 also requires the analyst to specify parametric model $h_0(X_{it}) = S_{it}\theta_0$. Then (4.3) becomes $Q(Y_{it}) = Q(V_{it}S_{it})\theta_0 + Q(l_1(\mathbf{X}_{it})) + Q(U_{it})$. In contrast to the CRE, where the analyst must assume $\delta_i = 0$ or specify the model $E[\delta_i | \overline{\mathbf{X}}_i]$, random coefficient δ_i cancels out.

Our approach contrasts with those of Klosin and Vilgalys (2023) and Semenova et al. (2023), both of which are based on learning the high-dimensional sparse representation of h_1 using a suitably specified data dictionary. The former is a FD procedure that removes δ_i from consideration provided the CATE is additively separable as above, whereas the latter implicitly assumes that $\delta_i = 0$.

5. DML PROCEDURES FOR PANEL DATA

We now describe our DML procedures for the fixed-effects estimators presented in Section 4. The objective is to make inferences on the target parameter θ_0 given a suitable predictions of nuisance functions η_0 based on the observed data $W_i = \{W_{it} : t = 1, \ldots, T\}$, where $W_{it} = \{Y_{it}, D_{it}, \mathbf{X}_{it}, \overline{\mathbf{X}}_i\}$ for CRE, and $W_{it} = \{Q(Y_{it}), Q(D_{it}), Q(\mathbf{X}_{it})\}$ for the transformation approach (noting that $W_{i1} \equiv \emptyset$ if Q is the FD transformation).

5.1. Learning the Nuisance Parameters

The first stage of DML involves using a suitably chosen learner (or combination of learners) to predict flexibly the unknown nuisance parameters. The nuisance functions $\eta_{0i} = \eta_0(\mathbf{X}_{i1}, \ldots, \mathbf{X}_{iT})$ vary between individuals and so the task is to learn η_0 . As discussed previously, the presence of the fixed effects and non-linear functional forms pose challenges which need to be addressed. Below we describe different procedures for learning η_0 .

5.1.1. Correlated Random Effects The estimation of θ_0 based on model (4.1)-(4.2) requires learning $\tilde{l}_1(\mathbf{X}_{it}, \overline{\mathbf{X}}_i)$ from the data $\{Y_{it}, \mathbf{X}_{it}, \overline{\mathbf{X}}_i : t = 1, \ldots, T\}_{i=1}^N$, and obtaining predicted residual \hat{V}_{it} to plug into Equation (4.1). Simply learning $\tilde{m}_1(\mathbf{X}_{it}, \overline{\mathbf{X}}_i)$ from the data $\{D_{it}, \mathbf{X}_{it}, \overline{\mathbf{X}}_i : t = 1, \ldots, T\}_{i=1}^N$ and using $\hat{V}_{it} = D_{it} - \hat{m}_1(\mathbf{X}_{it}, \overline{\mathbf{X}}_i)$ would ignore individual random effect c_i and so introduce bias. Hence, we propose the following approach:

STEP 1. Learn $\widetilde{m}_1(.)$ from $\{D_{it}, \mathbf{X}_{it}, \overline{\mathbf{X}}_i : t = 1, ..., T\}_{i=1}^N$ with prediction $\widehat{m}_{1i} = \widetilde{m}_1(\widehat{\mathbf{X}_{it}, \overline{\mathbf{X}}_i}).$ STEP 2. Calculate $\widehat{\overline{m}}_i = T^{-1} \sum_{t=1}^T \widehat{m}_{1i}.$ STEP 3. Calculate $\widehat{m}_1^*(\mathbf{X}_{it}, \overline{\mathbf{X}}_i, \overline{D}_i) = \widehat{m}_{1i} + \overline{D}_i - \widehat{\overline{m}}_i,$

where $m_1^*(\boldsymbol{X}_{it}, \overline{\boldsymbol{X}}_i, \overline{\boldsymbol{D}}_i) = E[D_{it} \mid \boldsymbol{X}_{it}, \overline{\boldsymbol{X}}_i] + c_i$. To capture c_i , it is essential to include the individual-specific treatment mean \overline{D}_i to ensure V_{it} is included in the prediction of m_1^* . For example, if the conditional distribution $D_{i1}, \ldots, D_{iT} \mid \mathbf{X}_{i1}, \ldots, \mathbf{X}_{iT}$ is multivariate normal with $E[D_{it}|X_{it}, \overline{X}_i] = \widetilde{m}_1(X_{it}, \overline{X}_i)$ and a homoskedastic variance-covariance matrix with random effects structure (induced by marginalizing over c_i) then $E(D_{it} \mid$ $X_{it}, \overline{X}_i, \overline{D}_i) = \widetilde{m}_1(X_{it}, \overline{X}_i) + c_i = m_1^*$ as required, but it is generally necessary to follow the steps above. More details are given in Section S1.1 of the Online Supplementary Information.

5.1.2. Transformation approaches First, we consider learning based on the transformed data alone. Generally, $Q(l_1(\mathbf{X}_{it})) \neq l_1(Q(\mathbf{X}_{it}))$ and $Q(m_1(\mathbf{X}_{it})) \neq m_1(Q(\mathbf{X}_{it}))$ unless l_1 and m_1 are linear. This means a consistent estimator cannot straightforwardly be constructed from the transformed data. However, we propose the following procedures: (a) Approximate procedure. Approximate model (4.3)-(4.4) by

 $O(Y_{ii}) \approx O(V_{ii})\theta_0 + l_1(O(\mathbf{X}_{ii})) + O(U_{ii})$

$$Q(Y_{it}) \approx Q(V_{it})\theta_0 + l_1(Q(\boldsymbol{X}_{it})) + Q(U_{it})$$
(5.1)

$$Q(V_{it}) \approx Q(D_{it}) - m_1(Q(\boldsymbol{X}_{it})), \qquad (5.2)$$

where l_1 and m_1 can be learnt from the transformed data $\{Q(Y_{it}), Q(X_{it}) : t = 1, \dots, T\}_{i=1}^N$ and $\{Q(D_{it}), Q(X_{it}) : t = 1, \dots, T\}_{i=1}^{N}$, respectively. This approach can produce good approximations of some non-linear nuisance functions but it is not robust and we were able to choose functions where it performed poorly.

(b) Exact procedure. Consider first the FD transformation $Q(Y_{it}) = Y_{it} - Y_{it-1}$. Then under Assumptions 3.1-3.5, $E[Y_{it} - Y_{it-1} | \mathbf{X}_{it-1}, \mathbf{X}_{it}] = \Delta l_1(\mathbf{X}_{it-1}, \mathbf{X}_{it})$ and $E[D_{it} - D_{it-1} | \mathbf{X}_{it-1}, \mathbf{X}_{it}] = \Delta m_1(\mathbf{X}_{it-1}, \mathbf{X}_{it})$, so that $\Delta l_1(\mathbf{X}_{it-1}, \mathbf{X}_{it})$ can be learnt using $\{Y_{it} - Y_{it-1}, \mathbf{X}_{it-1}, \mathbf{X}_{it} : t = 2, \dots, T\}_{i=1}^N$ and $\Delta m_1(\mathbf{X}_{it-1}, \mathbf{X}_{it})$ can be learnt using $\{D_{it} - D_{it-1}, \mathbf{X}_{it-1}, \mathbf{X}_{it} : t = 2, \dots, T\}_{i=1}^N$. However, for WG transformation $E[Y_{it} - \overline{Y}_i | L_T(\mathbf{X}_i)] = l_1(\mathbf{X}_{it}) - T^{-1} \sum_{s=1}^T l_1(\mathbf{X}_{is})$ and $E[D_{it} - \overline{D}_i|L_T(\mathbf{X}_i)] = m_1(\mathbf{X}_{it}) - T^{-1} \sum_{s=1}^T l_1(\mathbf{X}_{is})$ and $E[D_{it} - \overline{D}_i|L_T(\mathbf{X}_i)] = m_1(\mathbf{X}_{it}) - T^{-1} \sum_{s=1}^T l_1(\mathbf{X}_{is})$ $T^{-1}\sum_{s=1}^{T} m_1(\boldsymbol{X}_{is})$, the dimension of the learning problem is $O(T^2)$ times greater than that of the FD transformation and so unfeasible for non-trivial T.

(c) Hybrid procedure. Suppose the conditions set out in Section 4.1 for model (4.1)-(4.2) hold. Then $l_1(\mathbf{X}_{it}, \overline{\mathbf{X}}_i) = l_1(\mathbf{X}_{it}) + \omega_{\alpha}(\overline{\mathbf{X}}_i)$ and $\widetilde{m}_1(\mathbf{X}_{it}, \overline{\mathbf{X}}_i) = m_1(\mathbf{X}_{it}) + \omega_{\gamma}(\overline{\mathbf{X}}_i)$ learnt from the sample data satisfy $Q(\tilde{l}_1(\boldsymbol{X}_{it}, \overline{\boldsymbol{X}}_i)) = Q(l_1(\boldsymbol{X}_{it}))$ and $Q(\tilde{m}_1(\boldsymbol{X}_{it}, \overline{\boldsymbol{X}}_i)) =$ $Q(m_1(X_{it}))$. The hybrid approach can be used with both WG and FD transformations.

We elaborate on the justification for these procedures in Sections S1.2 and S1.3 of the Online Supplementary Information.

5.2. Neyman Orthogonal Score Function

The second stage of DML requires specifying an orthogonal score function for the structural parameter of interest to be solved after the nuisance functions have been plugged in. Following Chernozhukov et al. (2018), we construct a generic Neyman orthogonal score function for panel data that accounts for (a) the presence of the unobserved individual heterogeneity, and (b) non-linearity of the nuisance functions. Its properties allow us to obtain valid inferences using DML algorithms, provided the ML algorithms converge at rate $N^{1/4}$. We wish to make inference on the target parameter θ_0 given a suitable estimate of the nuisance parameter η_0 from the data. We repurpose the definition of W_{it} to define the entire dataset as $W = \bigcup_{i=1}^{N} W_i$, where $W_i = \{W_{it} : t = 1, \dots, T\}$, $W_{it} = w\{Y_{it}, D_{it}, X_{it}\}$ and w is a function or transformation of the data (possibly the identity) chosen by the analyst to implement one of the estimation approaches discussed in Section 4.

For individual i, let $\mathbf{r}_i = (r_{i1}, ..., r_{iT})'$ be the T residuals of the PO-PLPR model (either (4.1)-(4.2) or (5.1)-(5.2)), $X_i = X(X_{i1}, ..., X_{iT})$ the appropriately chosen set of predictor variables, and $\Sigma(X_i) = E[\mathbf{r}_i \mathbf{r}'_i \mid X_i]$ the (potentially) heteroskedastic residual variance-covariance matrix. Note that r_{it} differs from U_{it} in that it potentially incorporates the model random effect and other sources of autocorrelation.

Then the Neyman orthogonal score has the form

$$\psi^{\perp}(W_i;\theta_0,\boldsymbol{\eta}_0) = \mathbf{V}_i^{\perp} \Sigma_0^{-1}(X_i) \mathbf{r}_i, \qquad (5.3)$$

where row vector $\mathbf{V}_i^{\perp} = (V_{i1}^{\perp}, \dots, V_{iT}^{\perp})$ contains the orthogonalized regressors chosen to ensure Neyman orthogonality.

- (a) For CRE: Under model (4.1)-(4.2), $r_{it} = a_i + U_{it} = Y_{it} V_{it}\theta_0 \tilde{l}_1(\boldsymbol{X}_{it}, \overline{\boldsymbol{X}}_i),$ $V_{it} = D_{it} \tilde{m}_1(\boldsymbol{X}_{it}, \overline{\boldsymbol{X}}_i) c_i \text{ and } V_{it}^{\perp} = V_{it}.$ (b) For FD: Under model (4.3)-(4.4), $r_{it} = Q(u_{it}) = Q(Y_{it}) Q(V_{it})\theta_0 Q(l_1(\boldsymbol{X}_{it}))$
- and $V_{it}^{\perp} = Q(V_{it}).$

A derivation of (5.3) and its generalization to estimating conditional average treatment effects (CATEs) in the presence of heterogeneous treatment effects is provided in Section S2 of the Online Supplement.

5.3. Estimation and Inference About the Target Parameter

Suppose that suitable ML algorithms are available to learn the nuisance functions and obtain estimate $\widehat{\eta}$ on a pathway converging to η_0 at rate $N^{1/4}$. We now set out the DML2 (Chernozhukov et al., 2018, Def. 3.2) algorithm we use to make post-regularized inference about θ_0 based on Neyman-orthogonal score (5.3).

DML uses cross-fitting to control the overfitting or regularization bias introduced by plugging in $\hat{\eta}$ to (5.3). A key feature of our procedure is that we have adapted DML to use *block-k-fold* cross-fitting for panel data. This involves treating the time-series for each sample unit, W_i , as k-fold-sampling units when partitioning the sample into K equi-sized folds. Because W_i satisfies the independent and identically distributed assumption, whereas W_{it} does not, block cross-fitting accounts for the within-individual autocorrelation structure and so avoids having to make further post-fitting adjustments for autocorrelation to the estimmated standard errors as did Semenova et al. (2023).

Letting $\mathcal{W} = \{1, \ldots, N\}$ denote the indices of the sample units, our DML estimation procedure with *block-k-fold* cross-fitting can be set out as follows:

Algorithm 5.1. (DML estimator)

STEP 1. (SAMPLE SPLITTING AND CROSS-FITTING) Randomly partition the crosssectional units in the estimating sample into K folds of the same size. Denote the units in fold k = 1, ..., K by $\mathcal{W}_k \subset \mathcal{W}$ and let \mathcal{W}_k^c be its complement such that $N_k \equiv |\mathcal{W}_k| = N/K, |\mathcal{W}_k^c| = N - N_k$ and, because the folds are mutually exclusive and exhaustive, $\mathcal{W}_k \cap \mathcal{W}_j = \mathcal{W}_k \cap \mathcal{W}_k^c = \emptyset$ and $\mathcal{W}_k \cup \mathcal{W}_k^c = \mathcal{W}_1 \cup \ldots \cup \mathcal{W}_K = \mathcal{W}$. For K > 2, the larger complementary sample \mathcal{W}_k^c is used to learn the potentially complex nuisance parameters η , and \mathcal{W}_k for the relatively simple task of estimating the target parameter θ_0 .

STEP 2. (NUISANCE LEARNING) For fold k, learn η_0 from the data $\{W_i : i \in \mathcal{W}_k^c\}$ using one of the approaches from Section 5.1. Then use the learnt prediction rule $\hat{\eta}_k$ to construct $\psi^{\perp}(W_k; \theta, \hat{\eta}_k)$ for each fold.

STEP 3. (ESTIMATION) The DML estimator $\hat{\theta}^{DML}$ is obtained by solving

× 7

$$\frac{1}{K}\sum_{k=1}^{K}\frac{1}{N_{k}}\sum_{i\in\mathcal{W}_{k}}\boldsymbol{\psi}^{\perp}(W_{i};\widehat{\boldsymbol{\theta}}^{DML},\widehat{\boldsymbol{\eta}}_{k})=\boldsymbol{0},$$
(5.4)

which has closed-form solution

$$\widehat{\theta}^{DML} = \left(\sum_{k=1}^{K} \sum_{i \in \mathcal{W}_k} \widehat{\boldsymbol{V}}'_i \widehat{\boldsymbol{V}}_i\right)^{-1} \sum_{k=1}^{K} \sum_{i \in \mathcal{W}_k} \widehat{\boldsymbol{V}}'_i \widehat{\boldsymbol{U}}_i,$$
(5.5)

where $\widehat{\mathbf{V}}_i = (\widehat{V}_{i1}, \ldots, \widehat{V}_{iT})'$ and $\widehat{\mathbf{U}}_i = (\widehat{U}_{i1}, \ldots, \widehat{U}_{iT})'$. Note that the closed-form expression is obtained by setting $\Sigma_0 = I_T$ in equation (5.3) and capturing autocorrelation using cluster-robust versions of the following variance estimators. The estimator of the variance of $\widehat{\theta}^{DML}$ is

$$\widehat{\sigma}^2 = \widehat{J}^{-1} \left\{ \frac{1}{K} \sum_{k=1}^K \frac{1}{N_k} \sum_{i \in \mathcal{W}_k} \psi^{\perp}(W_i; \widehat{\theta}^{DML}, \widehat{\eta}_k) \psi^{\perp}(W_i; \widehat{\theta}^{DML}, \widehat{\eta}_k)' \right\} \widehat{J}_k^{-1}$$
(5.6)

where
$$\widehat{J} = K^{-1} \sum_{k=1}^{K} N_k^{-1} \sum_{i \in \mathcal{W}_k} \widehat{V}'_i \widehat{V}_i^{.5}$$

STEP 4. (ITERATION) Repeat STEPS 1-3 for each of the k-folds and average out the results.

Because score (5.3) is linear in θ , our proposition is that $\hat{\theta}^{DML}$ converges to normality according to the vectorized equivalent of Chernozhukov et al. (2018, Theorem 4.1) under the regularity conditions set out by Chernozhukov et al. (2018, Assumption 4.1).

6. SIMULATION STUDY

The primary focus of our Monte Carlo simulation study is to assess the performance of the DML procedures we propose in terms of bias and precision for different ML algorithms. We generated data under DGPs satisfying the assumptions required for the PO-PLPR model to hold for three different pairs of nuisance functions: linear, non-linear (smoothly continuous with no interactions), and non-linear (discontinuous with interactions). These are described along with a detailed discussion of the results in Online Supplement S3. The DGP satisfies the usual sparsity constraints because only two of the thirty predictors included in the analysis have non-zero effects.

Table 1 provides a summary of the main results under the most complex non-linear DGP for four of the procedures introduced above (Panels A-C). For each procedure, standard OLS estimates are compared with DML using four different learners for the

⁵We additionally make a finite-sample correction $(\hat{\theta}^{DML} - K^{-1}\sum_{k}\hat{\theta}_{k})^{2}$ to $\hat{\sigma}^{2}$, where $\hat{\theta}_{k} = (\sum_{i \in \mathcal{W}_{k}} \hat{V}'_{i}\hat{V}_{i})^{-1} \sum_{i \in \mathcal{W}_{k}} \hat{V}'_{i}\hat{U}_{i}$ is the fold-specific estimate as in Chernozhukov et al. (2018, p. C30).

non-linear nuisance functions. Across all three procedures, OLS can exhibit large bias and DML-LASSO typically outperforms the other learners in terms of both bias and precision, with the exception of the WG estimator (Panel C). In practice, LASSO requires the analyst to specify a sufficiently rich dictionary of non-linear terms; had only linear terms been included, for example, its performance would have been closer to OLS. Gradient boosting outperforms the other tree-based approaches (CART and RF) the performance of which in terms of bias and precision (SE/SD) is far inferior. A closer investigation, described in Online Supplement S3.2, found the sampling distributions of tree-based $\hat{\theta}_{DML}$ to be highly non-normal (see Figure S.1) so that inferences based on first-order asymptotic results are unreliable. Because trees are sensitive to hyperparameter choice (tree depth, etc.), we experimented with an alternative strategy for hyperparameter tuning (described in Online Supplement S3.2) that led to normal sampling distributions but larger biases. An extensive discussion on hyperparameter tuning can be found in Section S4.2 of the Online Supplement.

The superior performance of LASSO is indicated by the Root Mean Square Errors (RMSEs) of the estimator, model, and nuisance parameters. These results support the use of an *ensemble* learning strategy, in which the best-performing learner is chosen because, for our example, this would have selected LASSO and ensured reliable inference.

There is little to choose between the CRE, FD (Exact) and WG (Approximate) procedures. This is not unexpected because the DGP in the simulation study satisfies the additional assumptions for modelling the unobserved heterogeneity with Mundlak-like CRE, but the results are slightly different because these procedures involve learning different nuisance functions. More generally, the FD (Exact) estimator is the most robust procedure because it does not rely on any additional assumptions. In comparison, the WG (Approximation) method performed poorly no matter which learner was chosen because the approximation error induced by using the linearly transformed data (the time-demanded transformation) is large for the non-linear DGP.

7. EMPIRICAL APPLICATION

We reanalyse Fazio and Reggiani (2023)'s study of the impact of the National Minimum Wage (NMW) on voting behaviour in the UK. The data used in the original investigation come from the British Household Panel Survey (BHPS) comprising 4,927 working individuals (aged 18-64) between waves 1 and 16 (from 1991 until 2007).⁶ The treatment is measured by the question 'Were you paid the minimum wage' asked to those responding at Wave 9 in 1999. The authors use least squares to estimate the average treatment effect on the treated (ATT) between the NMW policy on various outcomes, including voting for conservative political parties which is our outcome of interest. Fazio and Reggiani (2023) compared the results of four regression specifications with different sets of covariates and fixed effects to capture all potential confounders, showing that the estimates did not vary across these, hence, concluding that workers who were paid the NMW rate were more likely to support conservative parties.

We revisit Specification (2) of Table 5 from the original paper that includes all control

⁶BHPS is a longitudinal survey study for British households that run from 1991 until 2009. The online replication package provides the instructions for data access and the codes to run the analysis in this section. DML estimation is conducted in R using the XTDML package in the replication package (or accessible in its latest version at https://github.com/POLSEAN/XTDML at the time of writing) which is built on DoubleML by Bach et al. (2024).

P. S. Clarke and A. Polselli

Table 1: Average MC simulation results, nonlinear and discontinuous DG
--

	$\operatorname{Bias}(\widehat{\theta})$	$\mathrm{RMSE}(\widehat{\theta})$	$\mathrm{SE}(\widehat{\theta})/\mathrm{SD}(\widehat{\theta})$	Model RMSE	RMSE_l	RMSE_m			
Panel A: CRE approach									
OLS	0.993	0.993	0.999						
DML-Lasso	0.009	0.014	1.235	5.818	1.981	1.432			
DML-CART	-0.087	0.199	0.084	21.316	7.169	5.910			
DML-RF	0.149	0.151	1.320	6.773	2.427	1.779			
DML-Boosting	-0.007	0.033	0.871	7.432	2.523	1.860			
Panel B: Exact approach with FD transformation									
OLS	0.993	0.993	0.951						
DML-Lasso	0.005	0.008	1.050	4.302	1.605	1.432			
DML-CART	0.291	0.385	0.049	46.564	18.889	12.596			
DML-RF	0.752	0.765	0.071	14.180	8.824	5.840			
DML-Boosting	-0.014	0.059	0.461	9.815	3.514	2.612			
Panel C: Approximation approach for WG transformation									
OLS	0.993	0.993	0.999						
DML-Lasso	0.977	0.977	1.120	4.347	9.625	6.450			
DML-CART	0.754	0.764	0.050	20.049	11.363	7.552			
DML-RF	0.972	0.972	0.739	5.070	9.817	6.578			
DML-Boosting	0.918	0.918	0.426	9.837	10.017	6.711			
Note: The figures in the table are the average values over the total number of Monte Carlo replications									

Note: The figures in the table are the average values over the total number of Monte Carlo replications (R = 100). The true target parameter is $\theta = 0.50$; N = 4000 and T = 10. The quantities displayed correspond to: $Bias(\hat{\theta}, \theta) = R^{-1} \sum_{r=1}^{R} (\mathbb{E}(\hat{\theta}_{r}) - \theta)$; $RMSE(\hat{\theta}) = \sqrt{R^{-1} \sum_{r=1}^{R} (Var(\hat{\theta}_{r}) + Bias(\hat{\theta}_{r}, \theta)^{2})}$, where $Var(\hat{\theta}) = \mathbb{E}[(\hat{\theta} - \mathbb{E}(\hat{\theta}))^{2}]$; Model RMSE= $(RK)^{-1} \sum_{r=1}^{R} \sum_{k=1}^{K} \sqrt{(|W|)^{-1} \sum_{i \in W} (\widehat{U_{it}} - \widehat{V_{it}} \widehat{\theta}_{k})^{2}}$, where $\widehat{U_{it}}$ and $\widehat{V_{it}}$ are the residuals of the PO structural equations; the RMSE of the nuisance parameters are $RMSE_{l} = (RK)^{-1} \sum_{r=1}^{R} \sum_{k=1}^{K} \sqrt{(|W^{c}|)^{-1} \sum_{i \in W^{c}} (y_{it} - \widehat{l_{k}}(.))^{2}}$ and $RMSE_{m} = (RK)^{-1} \sum_{r=1}^{R} \sum_{k=1}^{K} \sqrt{(|W^{c}|)^{-1} \sum_{i \in W^{c}} (d_{it} - \widehat{m_{k}}(.))^{2}}$. DML-Lasso uses 525 raw variables; the rest of the learners 30 raw variables. We use the Neyman-orthogonal PO score and five-fold cross-fitting.

variables, regions and wave fixed effects, by using DML procedure to fit the following PO-PLPR model

$$Vote_{it} = V_{it}\theta_0 + l_1(\boldsymbol{X}_{it}) + \alpha_i + U_{it}$$

$$(7.1)$$

$$V_{it} = NMW_{it} - m_1(\boldsymbol{X}_{it}) - \gamma_i \tag{7.2}$$

where $Vote_{it}$ is an indicator corresponding to one if respondent *i* voted for conservative parties in wave *t*, and zero otherwise; NMW_{it} is the treatment variable which is equal to one if the respondent's hourly pay increased due to the introduction of the NMW in 1999 and the respondent is observed in wave 9 onwards; X_{it} are the confounding variables whose functional form is *ex ante* unknown; and θ is the target parameter. The original base control variables included information about the age, education, marital status, household size, and income of other household members for respondent *i* in wave *t*.

We add thirty additional variables to the original specification, providing information on demographic characteristics, socio-economic status, employment and work-related variables, and ideology of the respondents. The complete list of confounders is reported in Table S.5 of the Online Supplement S5. We also increase the sample size by includ-

Double Machine Learning for Static Panel Models with Fixed Effects Table 2: The Effect of National Minimum Wage on Voting Behaviour in the UK

	OLS (1)	DML-Lasso (2)	DML-CART (3)	DML-RF (4)	DML-Boosting (5)					
Dependent variable: Voting for conservative parties										
Panel A: CRE approach										
NMW	0.051^{***}	0.045^{**}	0.069^{*}	0.180	-0.355					
Model RMSE	(0.019)	(0.020)	1.030	(0.151) 1.009	(0.232) 1.051					
RMSE of learner m No. control variables	144	0.422 0.064 1.477	0.419 0.051 144	0.413 0.016 144	0.435 0.010 144					
Panel B: Exact approach with FD										
ΔNMW	0.029 (0.026)	$0.030 \\ (0.025)$	$0.026 \\ (0.026)$	0.018 (0.026)	0.011 (0.026)					
Model RMSE RMSE of learner l RMSE of learner m	149	0.644 0.283 0.051 1.476	0.643 0.282 0.050 142	$\begin{array}{c} 0.640 \\ 0.282 \\ 0.049 \\ 1.42 \end{array}$	0.624 0.286 0.053 142					
No. control variables	140	1,470	140	140	140					
Panel C: Approximation approach with WG										
NMW	0.051^{***} (0.019)	0.048^{**} (0.019)	0.048^{**} (0.019)	0.041^{**} (0.019)	0.047^{***} (0.017)					
Model RMSE RMSE of learner l RMSE of learner m		$\begin{array}{c} 0.532 \\ 0.214 \\ 0.065 \end{array}$	$\begin{array}{c} 0.530 \\ 0.213 \\ 0.064 \end{array}$	$0.528 \\ 0.212 \\ 0.065$	$0.525 \\ 0.216 \\ 0.070$					
No. control variables	71	738	71	71	71					

Note: The table displays our estimates based on Specification (2) of Table 5 in Fazio and Reggiani (2023) using a different sample and confounders. Figures in Column (1) are estimated through OLS while the remaining columns using DML with different learners. The reported number of control variables includes the original variables, the individual means in DML-CRE, the lags of included controls in DML-FD, and the extended dictionary with polynomials and interactions for Lasso. The number of observations (NT) is 59,745 in Panels A and C, and 49,823 in Panel B; the number of cross-sectional units is 9,922 in all panels. The DML estimation uses 5-fold block cross-fitting and partialling-out score. The hyperparameters of the base learners are tuned with grid search. Cluster-robust standard errors at the respondent level in parenthesis. Significance levels: * p < 0.10, ** p < 0.05, *** p < 0.01.

ing two additional waves for years 2008 and 2009 (waves 17 and 18, respectively). Our final sample includes 9,922 working individuals. The functional form of the confounding variables is learnt using four different base learners: LASSO with a dictionary of non-linear terms (i.e., polynomials of order three and interaction terms of each raw variable), CART, RF, and boosted trees. The hyperparameter tuning of the base learners used in the implementation is discussed in the Online Supplement S4.3.

The point estimates of NMW effect on voting conservative parties are reported in Table 2. Column (1) shows OLS estimates (i.e. using standard panel estimation based on linear models) while the remaining columns contain the results using DML with different learners. Panel A displays the CRE estimates, Panel B the FD (Exact) estimates, and Panel C the WG (Approximation) estimates. The regression equations include the individual means of all included control variables when the CRE is used, and one-period

lagged variables of the controls when FD exact is used. Cluster-robust standard errors are reported in parenthesis with clustering at the respondent level.

The results exhibit considerable differences in the estimated effects between learners and estimators but, we argue, are consistent with those from a DGP which is either close to linear or smoothly non-linear as in our simulation study (see Online Supplement S3). The first difference we consider is between the FD (Exact) and CRE estimators. The former should be the most robust because it does not rely on the Mundlak-type models for fixed effects and also allows for different sets of omitted variables fixed only between distinct wave pairs t-1 and t (which could be represented by adding parameters like $\alpha_{i(t-1,t)}$ to the models). The magnitudes of the FD (Exact) point estimates are (a) smaller for OLS and LASSO and (b) more stable across all learners than those based on CRE. The RMSEs for the outcome (l) and treatment (m) models for FD (Exact) are smaller than those for CRE, which indicates the learners are more effective at learning nuisance functions of X_{it} and \mathbf{x}_{it-1} than of X_{it} and \mathbf{x}_i and so more likely to have errors bounded by $N^{1/4}$. Moreover, the consistency across OLS and the learners would suggest a DGP that is linear or smoothly non-linear. Second, the WG (Approximate) estimates are also consistent across learners and larger than those obtained using FD (Exact). The differences between FD (Exact) and WG (Approximate) are also commensurate with the additional robustness of the former. The CRE estimates obtained using tree-based learners, compared with all other estimates, are both larger (in absolute value) and less precise, with the DML-Boosting estimate particularly imprecise. However, the RMSEs indicate the tree-based algorithms do a better job learning the nuisance functions in X_{it} and $\bar{\mathbf{x}}_i$ than LASSO, and an ensemble strategy would have selected the estimate based on the random forest (DML-RF) for comparison with OLS and so pointed to a slightly larger but less significant estimate than the best estimate obtained using WG (Approximate). We finally reiterate the preference of FD (Exact) in this application, and the importance of comparing different procedures and learners in reaching this conclusion.

8. CONCLUSION

DML is a powerful tool for leveraging the power of ML for robust estimation of treatment effects, or policy-intervention. We provide estimation tools for applying DML when panel data are available which practitioners can use in place of existing ones or in a complementary way to test the robustness of their results to non-linearity. Moreover, the nature of static panels means that our procedures extend naturally to unbalanced panels provided that the non-response at each wave is conditionally independent of Y_{it} given \mathbf{X}_{it} and ξ_i . Further work on relaxing this assumption for CRE could build on Wooldridge (2019).

In general, although the three approaches we presented can be considered complementary to each other, the suitability of each will depend on the specific empirical example and the assumptions that the analyst is prepared to make. However, following our empirical studies, our recommendation for practice is to employ the FD (exact) approach because the ML algorithms we consider proved effective at estimating the nuisance functions involving predictors from waves t and t-1, the random coefficient treatment effect heterogeneity of Wooldridge (2019) is automatically accounted for, and it places fewer assumptions on the the distribution of the fixed effects.

From our simulation study, we also recommend using multiple (base) learners for DML stage one in the context of *ensemble* learning. In our simulations, the LASSO with an extended dictionary of non-linear terms performed best across all scenarios. However,

despite Lipschitz continuity and weak sparsity conditions holding, we found tree-based learners to perform poorly in terms of bias (of the point and interval estimates) and in terms of standard deviation and normality, even after making extensive efforts to implement adaptive hyperparameter tuning. We do not claim this to be a general result (there are many successful applications of tree-based learners in the literature), but we strongly recommend that analysts use an *ensemble* strategy involving multiple learners (the tree-based learners could perform better than the others in some scenarios) as is standard practice in other disciplines where ML and other AI-based methods are widely used.

Finally, our method is limited when there is treatment effect heterogeneity but the analyst still wishes to target the ATE rather than the CATE. Further work on extending these procedures, based on adapting the *interactive model* of Chernozhukov et al. (2018, Sec. 5) to panel data, would therefore be of great value for practice.

ACKNOWLEDGEMENTS

The authors are grateful to Andreas Alfons, Anna Baiardi, Thomas Cornelissen, Riccardo Di Francesco, Maria Grith, Omar Hussein, Damian Machlanski, Spyros Samothrakis, David Zentler-Munro, Wendun Wang, and the participants to the Annual MiSoC Workshop, RSS International Conference 2023, RCEA-ICEEF 2024, IPDC 2024, the internal seminars in ISER and in the Econometric Institute at the Erasmus University Rotterdam for the helpful comments and suggestions. We also thank the editor Christoph Rothe and three anonymous referees, whose comments significantly improved the final version of this paper. This research was funded by the UK Economic and Social Research Council award ES/S012486/1 (MiSoC). The authors also acknowledge the use of the High Performance Computing Facility (Ceres) and its associated support services at the University of Essex in the completion of this work.

REFERENCES

- Athey, S. and G. Imbens (2016). Recursive partitioning for heterogeneous causal effects. Proceedings of the National Academy of Sciences 113(27), 7353–7360.
- Athey, S., J. Tibshirani, and S. Wager (2019). Generalized random forests. The Annals of Statistics 47(2), 1148 – 1178.
- Bach, P., V. Chernozhukov, M. S. Kurz, M. Spindler, and S. Klaassen (2024). DoubleML
 An object-oriented implementation of double machine learning in R. Journal of Statistical Software 108(3), 1–56.
- Bach, P., V. Chernozhukov, and M. Spindler (2023). Heterogeneity in the US gender wage gap. Journal of the Royal Statistical Society Series A: Statistics in Society 187(1), 209–230.
- Baiardi, A. and A. A. Naghi (2024a). The effect of plough agriculture on gender roles: A machine learning approach. *Journal of Applied Econometrics*.
- Baiardi, A. and A. A. Naghi (2024b). The value added of machine learning to causal inference: Evidence from revisited studies. *The Econometrics Journal*, utae004.
- Belloni, A., V. Chernozhukov, and C. Hansen (2014). High-dimensional methods and inference on structural and treatment effects. *Journal of Economic Perspectives* 28(2), 29–50.
- Belloni, A., V. Chernozhukov, C. Hansen, and D. Kozbur (2016). Inference in highdimensional panel models with an application to gun control. *Journal of Business and Economic Statistics* 34(4), 590–605.

- Bia, M., M. Huber, and L. Lafférs (2023). Double machine learning for sample selection models. Journal of Business and Economic Statistics, 1–12.
- Breiman, L. (1996). Stacked regressions. Machine learning 24, 49-64.
- Chang, N.-C. (2020). Double/debiased machine learning for difference-in-differences models. *The Econometrics Journal* 23(2), 177–191.
- Chernozhukov, V., D. Chetverikov, M. Demirer, E. Duflo, C. Hansen, W. Newey, and J. Robins (2018). Double/debiased machine learning for treatment and structural parameters. *The Econometrics Journal* 21(1), C1–C68.
- Chernozhukov, V., W. K. Newey, and R. Singh (2022). Automatic debiased machine learning of causal and structural effects. *Econometrica* 90(3), 967–1027.
- Di Francesco, R. (2025). Ordered correlation forest. Econometric Reviews, 1–17.
- Fazio, A. and T. Reggiani (2023). Minimum wage and tolerance for high incomes. European Economic Review 155, 104445.
- Klosin, S. and M. Vilgalys (2023). Estimating continuous treatment effects in panel data using machine learning with an agricultural application. *arXiv preprint arXiv:2207.08789*. Second version; last revised 13 Sep 2023.
- Knaus, M. C. (2022). Double machine learning-based programme evaluation under unconfoundedness. *The Econometrics Journal* 25(3), 602–627.
- Kock, A. B. (2016). Oracle inequalities, variable selection and uniform inference in high-dimensional correlated random effects panel data models. *Journal of Economet*rics 195(1), 71–85.
- Kock, A. B. and H. Tang (2019). Uniform inference in high-dimensional dynamic panel data models with approximately sparse fixed effects. *Econometric Theory* 35(2), 295– 359.
- Kolesár, M., U. K. Müller, and S. T. Roelsgaard (2023). The fragility of sparsity. arXiv preprint arXiv:2311.02299.
- Lechner, M. (2015). Treatment effects and panel data. In The Oxford Handbook of Panel Data (Online Edition). Oxford Academic.
- Lechner, M. and J. Mareckova (2022). Modified causal forest. arXiv preprint arXiv:2209.03744.
- Mundlak, Y. (1978). On the pooling of time series and cross section data. *Economet*rica 46(1), 69–85.
- Robinson, P. M. (1988). Root-n-consistent semiparametric regression. Econometrica: Journal of the Econometric Society 56(4), 931–954.
- Rubin, D. B. (1974). Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of Educational Psychology* 66(5), 688–701.
- Sela, R. J. and J. S. Simonoff (2012). Re-em trees: a data mining approach for longitudinal and clustered data. *Machine Learning* 86, 169–207.
- Semenova, V., M. Goldman, V. Chernozhukov, and M. Taddy (2023). Inference on heterogeneous treatment effects in high-dimensional dynamic panels under weak dependence. *Quantitative Economics* 14(2), 471–510.
- Strittmatter, A. (2023). What is the value added by using causal machine learning methods in a welfare experiment evaluation? *Labour Economics* 84, 102412.
- University of Essex, Institute for Social and Economic Research (2018). British Household Panel Survey: Waves 1-18, 1991-2009. [data collection]. 8th Edition. UK Data Service. SN: 5151, DOI: http://doi.org/10.5255/UKDA-SN-5151-2.
- van der Laan, M. J., E. C. Polley, and A. E. Hubbard (2007). Super learner. Berkeley Devision of Biostatistics Working Paper Series, 222, pp.1–20.

- Wager, S. and S. Athey (2018). Estimation and inference of heterogeneous treatment effects using random forests. *Journal of the American Statistical Association* 113(523), 1228–1242.
- Wooldridge, J. M. (2019). Correlated random effects models with unbalanced panels. Journal of Econometrics 211(1), 137–150.
- Wooldridge, J. M. and Y. Zhu (2020). Inference in approximately sparse correlated random effects probit models with panel data. *Journal of Business and Economic Statistics* 38(1), 1–18.