Machine learning the first stage in 2SLS: Practical guidance from bias decomposition and simulation

Connor Lennon Edward Rubin* Glen R. Waddell

May 19, 2025

Abstract

Machine learning (ML) primarily evolved to solve "prediction problems." The first stage of two-stage least squares (2SLS) is a prediction problem, suggesting potential gains from ML first-stage assistance. However, little guidance exists on when ML helps 2SLS—or when it hurts. We investigate the implications of inserting ML into 2SLS, decomposing the bias into three informative components. Mechanically, ML-in-2SLS procedures face issues common to prediction *and* causal-inference settings—and their interaction. Through simulation, we show linear ML methods (*e.g.*, post-Lasso) work well, while nonlinear methods (*e.g.*, random forests, neural nets) generate substantial bias in second-stage estimates—potentially *exceeding* the bias of endogenous OLS.

JEL Codes: C26, C53, C18, C51, C52, C45 Keywords: machine learning, instrumental variables, two-stage least squares

*Lennon: Tacoma Power (clennon@uoregon.edu). Rubin (corresponding author): University of Oregon, Department of Economics, 1285 University of Oregon, 175 Prince Lucien Campbell Hall, Eugene, Oregon 97403-1205 (edwardr@uoregon.edu). Waddell: University of Oregon, Department of Economics and Research Fellow, IZA Bonn (waddell@uoregon.edu). We thank Joshua Angrist, Daniel Chen, Jiafeng Chen, Jonathan Davis, Thomas Dee, Max Farrel, Brigham Frandsen, Benjamin Hansen, Greg Lewis, John Loeser, Grant McDermott, Douglas Miller, James Stock, Jeffrey Wooldridge, and Eric Zou for helpful discussion. Lennon was affiliated with the University of Oregon during most of the project's work. The project's scripts and data will be posted to a publicly available Github repository. The project is exempt from IRB, as it does not involve human subjects. The authors have nothing to disclose.

1 Introduction

Machine learning (ML) methods appear in diverse empirical econometric applications. Despite this enthusiasm, the literature has little to say about the practical implications of combining ML methods with two-stage least squares (2SLS). In this paper we provide practical guidance on the potential risks and benefits of inserting standard ML methods into the first stage of 2SLS— drawing insights from a decomposition of potential sources of bias and a series of Monte Carlo simulations. The decomposition shows that ML-in-2SLS overlaps with the canonical *forbidden regression* but also reveals additional potential sources of bias. Our simulation results illustrate this potential is legitimate: some non-linear ML approaches generate larger second-stage bias than the original endogenous regression.

The motivation behind integrating machine learning in two-stage least squares is clear: to the extent incorporating "better" first-stage predictions is possible, researchers can obtain more precise second-stage estimates.¹ Because most ML methods are built explicitly for prediction—they typically outperform ordinary-least squares (OLS) at this task—using ML for first-stage predictions seems quite natural. Mullainathan and Spiess (2017) likewise highlight the reasoning that might lead a practitioner to inject ML into the 2SLS framework: "Machine learning... revolves around prediction" and "belongs in the part of the toolbox marked \hat{y} rather than in the more familiar $\hat{\beta}$ compartment." The authors immediately recognize that "the first stage of a linear instrumental variables regression is effectively prediction."² The risks of adopting out-of-the-box ML methods for 2SLS-type applications are less clear.

In this paper, we document the risks ML-in-2SLS poses for practitioners of applied econometrics. How does curating and generating first-stage predictions with ML affect the downstream, second-stage causal estimates of two-stage least squares?³ Using a simple decomposition, we discuss several phenomena that can bias ML-based 2SLS away from its target parameters. Some of these phenomena are implications of the *forbidden regression* (Angrist and Krueger, 2001; Angrist and Pischke, 2009; Wooldridge, 2010), which naïve implementations of ML in 2SLS are likely to lead

¹ The integration of 2SLS and ML has already appeared in applied work across many fields. Early examples have used linear models in the fist stage—predictors that typically perform well in our decomposition and simulation. For example, estimating labor market impacts of imprisonment (Mueller-Smith, 2015), the effects of racial-composition shocks during the Great Migration (Derenoncourt, 2022), the effect of expropriation on growth (Chen and Yeh, 2020), the "true" size of China's GDP growth (Chen et al., 2019), the inter-generational transmission of health (Bevis and Villa, 2020), and the heterogeneous impacts of family size and parental labor supply (Biewen and Kugler, 2021).

² Similar intuition is offered in (Belloni et al., 2011, 2012, 2013; Chernozhukov et al., 2015, 2018; Singh et al., 2019; Angrist and Frandsen, 2022; Singh et al., 2020; Chen et al., 2020) inclusive of new artificial intelligence (AI) methods (Hartford et al., 2017; Bennett et al., 2019; Liu et al., 2020). In a recent working paper, Chen et al. (2020) also recognizes this motivation, suggesting that the traditional OLS-based implementation of 2SLS "leaves on the table some variation provided by the instruments that may improve precision of estimates." If one is willing to accept the fairly strong assumption that any function (nonlinear or linear) of valid instruments is a valid instrument, Chen et al. (2020) provides an interesting solution to some of the challenges involved with including ML methods in 2SLS. We do not make this assumption.

³ With valid instruments, applying OLS in the first stage of 2SLS produces predictions (\hat{x}) that are a linear combination of the exogenous instruments. Thus, \hat{x} is itself exogenous in the traditional 2SLS procedure. Predictions produced by nonlinear functions are not guaranteed to be orthogonal to their residuals, generating additional bias/inconsistency in second-stage estimates.

to—injecting the predictions of a nonlinear estimator into the first stage of 2SLS.⁴

Other issues are less common to "traditional" econometrics but become key, we argue, to understanding ML-based results. These include:

- **Recovering endogeneity:** If the prediction algorithm is *too* good, then the first-stage predictor may entirely recover the endogenous regressor (including both *good* and *bad* variation). With (*i*) a small set of valid instruments and (*ii*) a linear estimator (*e.g.*, OLS), this scenario is of less concern. As the number of potential instruments increases and the estimator becomes more nonlinear and flexible (a hallmark of many ML methods), we show that this concern becomes real.
- Exclusion restrictions: ML methods are not designed to choose exclusion restrictions. If one relies on ML methods to determine a nonlinear functional form, choose instruments, and select first-stage controls in a 2SLS framework, then one must ultimately assume that the algorithm is capable of settling on a valid exclusion restriction. This places a lot of trust in ML to do something it is not typically designed to do. As Angrist and Frandsen (2022) point out, nonlinear estimators generate nonlinear combinations of the original instruments and thereby require additional exclusion restrictions *beyond* the original exclusion restriction implied by the linear combination of the instruments. With highly flexible ML methods, the set of exclusion restrictions is nearly infinite. One must either assume that (*i*) the ML algorithm will choose the appropriate exclusion restrictions or (*ii*) all possible exclusion restrictions are valid (as the algorithm's choice set is infinite).
- Amplified bias: As we show below, the bias of second-stage estimates in 2SLS is inversely related to the variance of the first-stage predictions (\hat{x}). Most ML methods reduce variance in the predictions (to reduce out-of-sample prediction performance in the canonical *bias-variance tradeoff*). This variance-reduction strategy leads to inflated bias in second-stage applications—a consideration not typical to OLS-based 2SLS applications.

In our simulations, most ML-rooted solutions that use common ML procedures in the first stage of 2SLS fail to improve upon standard 2SLS (*i.e.*, using OLS in the first stage) and generate more bias. Two linear estimators are the exception: post-Lasso selection and principal component analysis (PCA). Post-Lasso and PCA perform at least as well as standard OLS-based 2SLS. Perhaps more importantly, we show that highly nonlinear tree-based methods (*e.g.*, random forests and boosted trees) can amplify bias, providing parameter estimates farther from *truth* than naïve OLS regressions that ignore endogeneity. Given sufficient training time, naïve implementations of

⁴ Another flavor of the *forbidden regression* involves applying different specifications of controls in the first and second stages. Most out-of-the-box ML methods do not offer a method to ensure that second-stage controls are used for prediction in the ML-based first stage (and *in the correct functional form*). There are *ad hoc* solutions to this problem—writing custom functions that implement the ML algorithm *plus* a linear specification of the controls/fixed effects, or residualizing (i.e., Frisch-Waugh-Lovell). For an example, see the <code>fixest</code> package in R and its <code>feNmlm()</code> function, which is written to efficiently estimate maximum likelihood models with multiple fixed-effect (*i.e.*, large factor variables). This issue is particularly important for situations where conditioning on controls/fixed effects is integral to the instruments' exogeneity. Again, ML methods will, in this way, expose researchers to potential pitfalls.

neural networks in 2SLS can reproduce the original OLS bias, with little to no advantage over traditional approaches to recovering exogenous identifying variation through 2SLS.⁵

In Section 2 we formalize the theoretical settings and define the estimators. In Section 3 we introduce two data-generating processes. Because practitioners' use cases differ, we compare two general cases: one simple (and dense) case and a second, more complex and sparse case. In Section 4 we present the empirical results for the discussed estimators and DGPs. In Section 4 we also build a theoretical decomposition that explain how and why ML-based 2SLS procedures might increase second-stage bias.

While this paper focuses on a practical, *ad hoc* approach of inserting ML into the first stage of 2SLS, where ML methods might be thought of as complimenting traditional approaches to causal estimation, a separate strand of the literature offers tailor-made ML-based estimators in 2SLS-like structures. If one is willing to depart from a traditional 2SLS structure and accept different (typically stronger) identifying assumptions, these methods potentially capture more of the available "first-stage" variation without suffering from the issues we highlight in this paper. In Section 5 we discuss such a solution to the problems inherent to the *ad hoc* ML-based 2SLS approach. In Section 6 we offer concluding remarks.

Ultimately we conclude that while ML methods offer many promises for a range of applications, most out-of-the-box ML methods are not well suited for the first stage of two-stage least squares. Moreover, applying the *wrong* ML method in the first stage can actually generate more bias in parameter estimates than entirely ignoring endogeneity.

2 Models

2.1 The problem

Applied researchers commonly apply 2SLS to estimate the causal effect of some x on some y in a setting where the exogeneity of x cannot reasonably be assumed. In other words, where

$$y = \beta_0 + \beta_1 x + u , \tag{1}$$

there is concern over the potential for non-zero covariance between the variable of interest *x* and the disturbance *u* when estimating the parameter β_1 .

Let z denote a vector of *instrumental variables*. We express the first stage of a 2SLS estimates x as a function of these instruments:

$$x = f(\mathbf{z}) + \varepsilon . \tag{2}$$

In its traditional OLS-based implementation, $f(\mathbf{z})$ is linear in \mathbf{z} .

⁵ This ignores the practical as well. With our resources, the simulation for neural network frequently took several hours to complete which is considerably longer than the time it takes to run a traditional 2SLS.

Defining the predictions from (2) as $\hat{x} = f(\mathbf{z})$, the second stage of the 2SLS procedure then regresses the outcome variable *y* on \hat{x} ,

$$y = \gamma_0 + \gamma_1 \hat{x} + w , \qquad (3)$$

to achieve an estimate for β_1 in (1). We define $\hat{\gamma}_1$ as this estimate of β_1 . If the instruments are valid (*i.e.*, predictive of *x* and uncorrelated with *u*) and \hat{x} results from an OLS regression, then \hat{x} will also be exogenous.⁶ The second stage of OLS-implemented 2SLS then generates consistent estimates of β_1 , interpreted as the causal effect of *x* on *y*.

So why introduce ML? Applications of 2SLS identify the effect of x on y by extracting only a fraction of the "good" (exogenous) variation in x. The hope for ML-based 2SLS methods is that researchers can extract more of the good variation in x—a more flexible fit of the exogenous variation—while still omitting the bad variation. This desire has likely increased following Lee et al. (2022), who argue that many traditional evaluations of instrumental variables considerably overestimate their significance.

2.2 Estimators

In the analysis below we examine three classes of 2SLS-motivated estimators:

Class 1: "Traditional" two-stage regression methods: This set of estimators covers the standard two-stage regression estimators in an econometrician's toolbox: two-stage least squares, (unbiased) split-sample IV (Angrist and Krueger, 1995), the Fuller implementation of limitedinformation maximum likelihood (LIML) (Anderson and Rubin, 1949; Fuller, 1977), and jackknife IV (JIVE) (Angrist et al., 1999). These methods overlap in three important ways: they (*i*) employ a two-stage approach (*ii*) whose first stage creates a linear combination of the instruments (*iii*) with no formal variable selection.

Class 2: Machine-curated variable selections in standard 2SLS: This second class augments the standard OLS-based version of 2SLS with variable selection/synthesis. Specifically, these methods feature an additional procedure, *prior to the first stage*, that downselects or combines z into a more parsimonious set of variables. The elements of this more parsimonious expression of z then appear in the first stage. The rest of the 2SLS process proceeds as usual (*i.e.*, OLS). Importantly, while these models feature variable selection or synthesis, they also preserve linearity in both stages. Because these estimates result from linear combinations of z, the original exclusion restriction of z passes through to the selected/synthesized instruments.

Our first machine-curated method is the post-Lasso procedure of Belloni et al. (2012), which first estimates the linear relationship between *x* and **z** (a linearized version of 2) using penalized regression. This penalized regression minimizes the sum of squared error (SSE) *plus* a penalty proportional to the sum of the coefficients' magnitudes. That is, $\lambda \times ||\gamma||$, where γ is the vector of

⁶ We assume homogeneous treatment effects, which removes the requirement of monotonicity.

coefficients on the (standardized) instruments and λ is a the shrinkage parameter chosen by the researcher (typically via cross validation). Because each instrument's coefficient-based penalty changes discontinuously when moving away from $\gamma_i = 0$, Lasso can be used to select a set of *stronger* instruments (whose coefficients are non-zero). Post-Lasso selects the instruments whose coefficients are non-zero and then estimates standard, OLS-based 2SLS using those selected instruments.⁷

Principal-component analysis (PCA) offers an alternative route to simplifying z by selecting z's first k principal components (Pearson, 1901). Thus, as the second machine-curated method we consider, principal-component analysis (PCA) applied to 2SLS (as in Ng and Bai (2009) and Winkelried and Smith (2011)). This approach passes a set of principal components into the first stage of standard OLS-based 2SLS. While PCA may reduce the first stage's interpretability, this approach can drastically reduce the number of first-stage instruments while retaining considerable explanatory power.

Class 3: ML-based first stages in 2SLS: Our final class of estimators retains the general two-step framework of 2SLS but replaces the first stage with a variety of cross-validated ML algorithms. We evaluate a meaningful subset of machine-learning methods suitable for regression, including random forest (Ho, 1995; Breiman, 2001), boosted trees (Breiman, 1997; Mason et al., 1999; Friedman, 2001, 2002), neural networks (Turing, 2004; McCulloch and Pitts, 1943; Farley and Clark, 1954), and Lasso (Tibshirani, 1996; Santosa and Symes, 1986).^{8,9} Notably, most of these algorithms offer considerable flexibility (*e.g.*, nonlinearity in **z**) and variable selection (to varying degrees). This class offers considerable insights into the merits of off-the-shelf ML methods for machine-assisted 2SLS.

3 Data-generating processes

In order to examine the performance of ML in the predictive stage of 2SLS—in absolute terms and relative to "traditional" options—we employ two general data-generating processes (DGPs). For reasons described below we refer to the two DGPs as the *low-complexity* case and the *high-complexity* case. In practice, the researcher rarely knows the extent to which her case is *complex,* particularly in terms of extent of nonlinearity or the efficient number of instruments. While "complexity" is certainly subjective, our intention is to bookend the settings for which an applied researcher might apply ML-based 2SLS.

⁷ Angrist and Frandsen (2022) notes that this methodology may suffer from potentially unseen pre-test bias. Because our model comes from relatively strong instruments, as with the intuition of Zhao et al. (2021), we do not estimate de-biased Lasso models. We therefore allow post-Lasso to serve as a representation of both.

⁸ For our purposes, the contributions of Srivastava et al. (2014) (dropout), Ioffe and Szegedy (2015) (batch normalization), and Kingma and Ba (2017) (stochastic optimization) are particularly relevant.

⁹ For a nice review of ML methods in applied economics, including Lasso, tree-based methods, and neural networks, please see Storm et al. (2019). For broader and more in-depth coverage (from the authors of many of the methods), see James et al. (2013) and Hastie et al. (2009).

3.1 A low-complexity case

In this case, we aim to depict the performance of various estimators when the DGP is simple (few instruments and non-sparse) and closely matches the ideal scenario for OLS-based 2SLS: an endogenous regressor that is a linear combination of a relatively small set of strongly predictive, exogenous instruments. This case is applicable to researchers seeking to estimate the causal effect of a variable of interest x_1 on outcome y,

$$y = \beta_0 + \beta_1 x_1 + \varepsilon_y , \qquad (4)$$

but facing the challenge (e.g., omitted variables, simultaneity) that x_1 is endogenous and $E[\varepsilon_y|x_1] \neq 0$ prevents OLS from cleanly identifying β_1 in (4). Importantly, the causal effect β_1 is common across all individuals, which ensures differences across estimators are not due to the estimators recovering different local average treatment effects (LATEs).

In this low-complexity scenario, ML-based 2SLS methods are overkill: neither variable selection nor nonlinearity are necessary. In fact, our results demonstrate that ML methods can increase bias relative to 2SLS, even relative to endogenous OLS.

Formally, to model a scenario with a single endogenous regressor (x_1) and a small set of valid (and *individually* strong) instruments, we define the DGP as

$$\varepsilon_y = \beta_2 x_2 + \eta$$
,
 $x_2 = 1 + \varepsilon_c$, and
 $x_1 = g_x(\mathbf{z}) + \varepsilon_c$,

drawing special attention to the inclusion of ε_c as the disturbance common to both x_1 (the variable of interest) and x_2 (the *omitted* variable). This common error follows a standard normal distribution; η is distributed uniformly between -1 and 1.

We assume that a set of valid instruments \mathbf{z} exists such that $E[\varepsilon_y | \mathbf{z}] = 0$ and $E[x_1 | \mathbf{z}] \neq 0$ (we focus on the case where $|\mathbf{z}| = 7$). We also anticipate that the researcher has no beliefs or insights about the functional form of $g_x(\cdot)$, as is often the case in practice. In the true DGP for this case, $g_x(\mathbf{z}) = \sum_{i=1}^{7} z_i$. That is, $g_x(\cdot)$ is linear.

In particular, we draw the instruments **z** from a multivariate normal distribution centered at zero (*i.e.*, $E[\mathbf{z}] = \mathbf{0}$) with variance-covariance matrix $\Sigma_{\mathbf{z}}$ where $\widehat{\text{Cov}}(z_i, z_j) = 0.6^{|h-k|}$ (and thus $\text{Var}(z_i) = 1$ for each *i*). By implication, $x_1 \sim N(0, \text{ Grand Sum}(\Sigma_{\mathbf{z}}) + 1)$.

In full, then, the data represents the following system of equations:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \eta$$

$$x_2 = 1 + \varepsilon_s ,$$

$$x_1 = g_x(\mathbf{z}, \varepsilon_s) = \sum_{i=1}^7 z_i + \varepsilon_s .$$

Notably, the specification of the instruments in this DGP produces a very strong first-stage with a relatively large *concentration parameter* (Belloni et al., 2012). Put simply, the concentration parameter μ^2 describes the extent to which the weak-instrument problem may arise within a given DGP. A higher value of μ^2 implies that 2SLS, without variable selection, will converge to the true β_1 at relatively small sample sizes.¹⁰ Consequently, the low-complexity case allows us to test how machine-curated first stages perform when there is little to be gained from variable selection/synthesis.

3.2 A high-complexity case

As our high-complexity case, we follow Belloni et al. (2012)'s DPG with two extensions. This DGP allows the researcher to tailor instruments' strengths with *many* instruments. Following Belloni et al. (2012), the DGP in our high-complexity case results from

$$y = \beta_0 + \beta_1 x_1 + \varepsilon_y$$
$$x_1 = \pi \mathbf{z} + \varepsilon_v ,$$

where

$$(\varepsilon_{y}, \varepsilon_{v}) \sim N\left(0, \begin{bmatrix} \sigma_{y}^{2} & \sigma_{y}\sigma_{v} \\ \sigma_{v}\sigma_{y} & \sigma_{v}^{2} \end{bmatrix}\right),$$

$$\mathbf{z} = \begin{bmatrix} z_{1} & z_{2} & \cdots & z_{100} \end{bmatrix} \sim N(\mathbf{0}, \Sigma_{z}),$$

$$\Sigma_{z}[j, j] = \operatorname{Var}(z_{j}) = \sigma_{j}^{2} = 1, \quad \forall j \in \{1, \dots, 100\}, and$$

$$\Sigma_{z}[j, k] = \widehat{\operatorname{Cov}}(z_{j}, z_{k}) = \operatorname{Cor}(z_{j}, z_{k}) = 0.6^{|j-k|}, \quad \forall (j, k) \in \{1, \dots, 100\}$$

As before, we imagine the researcher's interest in this case focuses on identifying β_1 . However, unlike the earlier DGP, the high-complexity case produces sets of relevant and exogenous instruments that vary in their correlation and individual strength (*i.e.*, π_i).

In defining the "exponential" design of the *first-stage* coefficient vector π , we follow Belloni et al. (2012): π captures a "beta pattern" $\tilde{\pi} = (0.7^0, 0.7^1, 0.7^2, \dots, 0.7^{99})$ that is then multiplied by

¹⁰ In this "low-complexit" case, $\mu^2 \approx n \times 20.71$, which exceeds the values in Belloni et al. (2012). We discuss μ^2 's role further in the *high-complexity case* section below.

a constant *C*, *i.e.*, $\pi = C \times \tilde{\pi}$. The constant *C* implies a value for the concentration parameter, $\mu^2 = \frac{n\pi'\Sigma_c\pi}{\sigma_v^2}$.¹¹ In panels 1A–1C (Figure 1) we illustrate the three beta patterns that we adopt in the "high-complexity" DGP, generating three subcases of this DGP. As described above, the concentration parameter is useful for determining the behavior of IV estimators. Because we are less interested in the case of weak instruments, we use $\mu^2 = 180$, which creates a strong—though fairly sparse—set of instruments as outlined in Belloni et al. (2012).¹²

Belloni et al. (2012) arrange the coefficients π in descending order (*i.e.*, $\pi_1 > \pi_2 > \cdots > \pi_{100}$). However, the definition of Σ_z implies that "proximate" instruments are more correlated than "distant" instruments (*i.e.*, Cor(z_i , z_{i+1}) > Cor(z_i , z_{i+k}) for k > 1). Thus, the DGP of Belloni et al. (2012) ensures the strongest instruments correlate with each other. While this feature may be desirable in many contexts, we remain agnostic with regard to whether the strongest instruments are most correlated with each other or with other instruments. However, this agnosticism requires that we consider three sub-cases, each arising from alternative orderings of the coefficients in π and the covariance of the respective variables:

- Randomly shuffled: After generating the coefficients, we randomly re-order them to break the relationship between instruments' strengths and their covariance (Σ_z).
- **Descending from** z_1 : In this subcase, as in Belloni et al. (2012), $\pi_1 > \pi_2 > \cdots > \pi_{100}$.
- Descending from *z*₅₀: Here we modify Belloni et al. (2012) by defining *π*₅₀ as the largest coefficient: *π*₅₀ > *π*₅₁ > · · · > *π*₁₀₀ > *π*₁ > *π*₂ > · · · > *π*₄₉. Because "proximate" instruments are correlated in Σ_z, this subcase implies that the strongest instrument (*z*₅₀) is very correlated both with the second-strongest instrument (*z*₅₁) and with the weakest instrument (*z*₄₉).

Finally, we define $\sigma_v^2 = \pi' \Sigma_z \pi$ (which forces that $Var(x_1) = 1$) and $\sigma_y = 1$. In panels 1A-1C of Figure 1 we illustrate the cross-instrument correlations implied by Σ_z : in Panel 1D we show a correlation matrix among the 100 instruments, and in Panel 1E we highlight the correlation of z_1 and z_{50} to each of the other 100 instruments. Instruments are strongly correlated with their neighbors and weakly correlated with non-neighbors, which limits the information accessible from any single instrument.

4 Results

We now discuss the results of our simulations. In every simulation, we include an "oracle model" that extracts the entirety of the exogeneous component of x_1 (perfectly removing endogeneity) and a simple OLS model (where we entirely ignore endogeneity). While one might expect the oracle and plain OLS models to bookend the biases in 2SLS-related models, our simulations

¹¹ For a proof of this statement, see Belloni et al. (2012).

¹² It is important to select this value thoughtfully. Choosing a μ^2 that is too small will simulate a weak-instruments problem. Choosing a μ^2 that is too large will yield a scenario in which all instruments are "overpoweringly" valid, which reduces the effectiveness of selection or dimension-reduction techniques. See Hansen et al. (2008) for additional discussion of μ^2 .

demonstrate that they *do not*. That is, inserting machine learning into the first stage can lead to outcomes that are even worse than ignoring endogeneity.

In each case, we are interested in the performances of the estimators in terms of their biases and the precision of estimates. Recall that these estimators include three broad classes: (*i*) traditional methods (OLS-based 2SLS, split-sample IV, LIML, and jackknife IV), (*ii*) machine-curated 2SLS (variable-selection or -curation via post-Lasso and PCA), and (*iii*) 2SLS applications with ML-powered predictions in their first stages (*i.e.*, replacing first-stage OLS with either Lasso, boosted trees, random forests, or neural networks).

4.1 General results

In Figure 2 we depict the distributions of point estimates ($\hat{\beta}_1$) for a given method in the givezn DGPs. Panel 2A illustrates the low-complexity case; panels 2B–2D presents the results for our high-complexity cases. Table 1 summarizes each of these method-by-DGP combinations with the mean and standard error from each. The target parameter β_1 equals 1 throughout the simulations (indicated with a thin dashed line). Each distribution results from 1,000 iterations of the simulation.

To those with use-cases that resemble our "low-complexity case," the simulation results have a clear takeaway: PCA-based 2SLS and post-Lasso perform well and offer very safe choices.¹³ Important for the practitioner: All four nonlinear ML-in-the-first-stage methods (*i.e.*, Lasso, boosted trees, neural networks, and random forests) perform poorly in terms of bias and variance. In fact, random-forest-based 2SLS generates *more bias* in $\hat{\beta}_1$ than the OLS estimator that entirely ignores endogeneity—it is possible for an ML-based 2SLS estimator to *amplify* bias relative to plain OLS. We discuss the source of this bias amplification in the next section.

In the three high-complexity cases in Table 1 (columns *B–D*) and in panels 2B–2D of Figure 2, LIML and Jackknife IV generate very little bias in their estimates of β_1 , outperforming 2SLS. Across all three DGPs, 2SLS produces mean estimates roughly 2.3–5.8 percent larger than the true parameter, while the centers of LIML's and JIVE's distributions are within 0.4 percent of the true parameter. Injecting random forests into the first stage, on average, produces more biased estimates than naïve (endogenous) OLS, generating coefficient estimates that are 32–56 percent larger than the true estimates. Again, one can worsen endogeneity issues by using ML-based 2SLS estimators.

¹³ LIML also performs well, but with slightly larger variance. The Jackknife IV estimator yields *very* high variance in this low-complexity DGP, as do Neural Networks.

4.2 Decomposing the bias

To diagnose the sources of bias from different methods, we show one can decompose the wedge between β_1 and $\hat{\beta}_1^{2SLS}$ into three components,

Wedge =
$$\hat{\beta}_1^{2\text{SLS}} - \beta_1 = f\left(\beta_1 \widehat{\text{Cov}}(\hat{x}, e), \widehat{\text{Cov}}(\hat{x}, u), \frac{1}{\widehat{\text{Var}}(\hat{x})}\right),$$
 (5)

where *f* is non-decreasing with respect to each of its arguments, \hat{x} is the first-stage-based prediction of *x* from some set of valid instruments, *e* denotes the resulting first-stage residuals $(x - \hat{x})$, and *u* represents the population disturbance from regressing *y* on *x*, *i.e.*, $u = y - (\beta_0 + \beta_1 x)$. Below we derive and elaborate upon (5).

Each component of the wedge offers insights into how first-stage methods differentially produce biases—this delivers helpful intuition regarding the pitfalls that may arise in 2SLS applications that include ML-based first stages.

To see the component parts of the bias drawing $\hat{\beta}_1^{2\text{SLS}}$ away from β_1 , suppose again that the parameter of interest is β_1 , the causal effect of *x* on *y* in

$$y = \beta_0 + \beta_1 x + u . \tag{6}$$

Suppose also that x is endogenous, *i.e.*, $Cov(x, u) \neq 0$. The 2SLS estimate of β_1 comes from estimating

$$y = \beta_0 + \beta_1 \hat{x} + w , \qquad (7)$$

where again \hat{x} is the first-stage-based prediction of x from some set of valid instruments $\mathbf{z} = z_1, z_2, \dots, z_p$.

Because we estimate the second stage in (7) via OLS, the estimate for β_1 can be written

$$\hat{\beta}_1^{2\text{SLS}} = \beta_1 + \frac{\widehat{\text{Cov}}(\hat{x}, w)}{\widehat{\text{Var}}(\hat{x})} .$$
(8)

where $\widehat{\text{Cov}}(\cdot)$ and $\widehat{\text{Var}}(\cdot)$ refer to the sample-based covariance and variance.

Using (6) and (7), we can rewrite w as

$$w = y - (\beta_0 + \beta_1 \hat{x}) = \beta_0 + \beta_1 x + u - \beta_0 - \beta_1 \hat{x} = \beta_1 (x - \hat{x}) + u = \beta_1 e + u ,$$
(9)

where *e* is the first-stage residual, the difference between *x* and \hat{x} .

Substituting (9) for *w*, we can decompose the covariance in (8) into two components:

$$\widehat{\text{Cov}}(\hat{x}, w) = \widehat{\text{Cov}}(\hat{x}, \beta_1 e + u)$$
$$= \beta_1 \widehat{\text{Cov}}(\hat{x}, e) + \widehat{\text{Cov}}(\hat{x}, u) .$$
(10)

If the first-stage predictions (\hat{x}) come from OLS, then $\widehat{\text{Cov}}(\hat{x}, e)$ is mechanically zero. The second term, $\widehat{\text{Cov}}(\hat{x}, u)$, is typically small when \hat{x} comes from a linear combination of valid instruments.

Finally, substituting (10) into (8) yields a helpful expression for the 2SLS estimate for β_1 ,

$$\hat{\beta}^{2\text{SLS}} = \beta_1 + \frac{\beta_1 \widehat{\text{Cov}}(\hat{x}, e) + \widehat{\text{Cov}}(\hat{x}, u)}{\widehat{\text{Var}}(\hat{x})} . \tag{11}$$

Again, first-stage OLS guarantees that $\widehat{\text{Cov}}(\hat{x}, e)$ is zero and, with valid instruments, that $\widehat{\text{Cov}}(\hat{x}, u)$ is small. Whether $\widehat{\text{Var}}(\hat{x})$ is "small" is typically of little consequence with OLS (as $\beta_1 \widehat{\text{Cov}}(\hat{x}, e) + \widehat{\text{Cov}}(\hat{x}, u)$ is typically small). However, all three points can generate important issues when we mix ML methods into the first stage of 2SLS. With ML methods, nothing guarantees that $\widehat{\text{Cov}}(\hat{x}, e)$ is zero or that $\widehat{\text{Cov}}(\hat{x}, u)$ is small. Moreover, many ML methods are constructed to *reduce* the variance of predictions, which further amplifies bias. This variance-reduction aspect is particularly relevant for nonlinear methods.

The $\beta_1 \text{Cov}(\hat{x}, e)$ component

For the term $\beta_1 \text{Cov}(\hat{x}, e)$ to differ from zero and generate bias, $\beta_1 \neq 0$ and $\text{Cov}(\hat{x}, e) \neq 0$. We assume that the population-regression coefficient β_1 differs from zero. Consequently, the term $\beta_1 \text{Cov}(\hat{x}, e)$ only generates bias when $\text{Cov}(\hat{x}, e) \neq 0$; β_1 scales the bias and affects its direction.

By construction, OLS produces predictions that are orthogonal to their residuals, *i.e.*, $Cov(\hat{x}, e) = 0$. This first term is therefore irrelevant when the first stage uses OLS. However, when practitioners adopt other methods in the first stage (*e.g.*, non-linear methods) nothing guarantees first-stage predictions are uncorrelated with their residuals. Indeed, this component relates to many researchers' definitions of *the forbidden regression*. This part of the bias results from using estimators whose predictions correlate with their residuals (rather than resulting from a violation of the exclusion restriction). While nonlinear methods can generate $\widehat{Cov}(\hat{x}, e) = 0$, many do not (as illustrated by the column *a* of Table 2).

In addition, because $\widehat{\text{Cov}}(\hat{x}, e)$ typically drops out of OLS regression, OLS-based empirical intuition does not help here. One implication from this non-OLS intuition of $\widehat{\text{Cov}}(\hat{x}, e)$ is that the bias generated by this component is proportional to the size of the target parameter β_1 . Where treatment effects are larger, the bias transmitted through this component is also larger.

To understand why some methods produce larger values of $\widehat{\text{Cov}}(\hat{x}, e)$ than other methods, first

decompose this covariance into $\widehat{\text{Cov}}(\hat{x}, x)$ and $\widehat{\text{Var}}(\hat{x})$:

$$\widehat{\text{Cov}}(\hat{x}, e) = \widehat{\text{Cov}}(\hat{x}, x - \hat{x}) = \widehat{\text{Cov}}(\hat{x}, x) - \widehat{\text{Var}}(\hat{x}) .$$
(12)

While $\widehat{\text{Cov}}(\hat{x}, e)$ is not generally signable, it is bounded between $-\widehat{\text{Var}}(\hat{x})$ and $\widehat{\text{Cov}}(\hat{x}, x)$.¹⁴ Further, we can sign $\beta_1 \widehat{\text{Cov}}(\hat{x}, e)$ in five subcases:¹⁵

$$\operatorname{Sign}\left\{\beta_{1}\widehat{\operatorname{Cov}}(\hat{x}, e)\right\} = \operatorname{Sign}\left\{\beta_{1}\widehat{\operatorname{Corr}}(\hat{x}, e)\right\}$$
$$= \operatorname{Sign}\left\{\beta_{1}\sigma_{e}^{-1}\left(\widehat{\operatorname{Corr}}(\hat{x}, x)\sigma_{x} - \sigma_{\hat{x}}\right)\right\}$$
$$= \operatorname{Sign}\left\{\beta_{1}\right\} \cdot \operatorname{Sign}\left\{\widehat{\operatorname{Corr}}(\hat{x}, x)\sigma_{x} - \sigma_{\hat{x}}\right\}$$
$$\begin{pmatrix} (+) & \text{if }\beta_{1} > 0 \text{ and }\widehat{\operatorname{Corr}}(\hat{x}, x)\sigma_{x} > \sigma_{\hat{x}} \\ (-) & \text{if }\beta_{1} > 0 \text{ and }\widehat{\operatorname{Corr}}(\hat{x}, x)\sigma_{x} < \sigma_{\hat{x}} \\ (-) & \text{if }\beta_{1} < 0 \text{ and }\widehat{\operatorname{Corr}}(\hat{x}, x)\sigma_{x} < \sigma_{\hat{x}} \\ (+) & \text{if }\beta_{1} < 0 \text{ and }\widehat{\operatorname{Corr}}(\hat{x}, x)\sigma_{x} < \sigma_{\hat{x}} \\ (+) & \text{if }\beta_{1} < 0 \text{ and }\widehat{\operatorname{Corr}}(\hat{x}, x)\sigma_{x} < \sigma_{\hat{x}} \\ 0 & \text{if }\beta_{1} = 0 \text{ or }\widehat{\operatorname{Corr}}(\hat{x}, x)\sigma_{x} = \sigma_{\hat{x}}, \end{cases}$$
(13)

where σ_x refers to the standard deviation of *x* ($\sigma_{\hat{x}}$ and σ_e are defined similarly).

As (13) reveals, the sign of $\beta_1 \widehat{\text{Cov}}(\hat{x}, e)$ depends on two quantities: (*i*) the sign of β_1 , and (*ii*) the sign of $\widehat{\text{Corr}}(\hat{x}, x) \sigma_x - \sigma_{\hat{x}}$. It is difficult to generalize the sign of $\widehat{\text{Cov}}(\hat{x}, e)$ without further assumptions. While one may be tempted to assume $\sigma_x > \sigma_{\hat{x}}$, this assumption is not sufficient for signing $\widehat{\text{Cov}}(\hat{x}, e)$, as it still depends upon the magnitude of $\widehat{\text{Corr}}(\hat{x}, x)$.¹⁶ The knife-edge case where a = 0 appears unlikely except in cases where either $\beta_1 = 0$ or where $\widehat{\text{Cov}}(\hat{x}, e)$ is mechanically zero (*e.g.*, OLS).

Across the twelve models that we consider in Table 2, only the non-OLS models produce $\widehat{\text{Cov}}(\hat{x}, e) \neq 0$ —this is unsurprising. Lasso, neural nets, boosted trees, and random forests all produce positive covariance between \hat{x} and e. In other words, in all of our DGPs, the term $\widehat{\text{Cov}}(\hat{x}, e)$ biases $\hat{\beta}$ upward (positively) whenever it is non-zero.¹⁷ Random forest models generate the largest covariance between \hat{x} and e (and consequently the largest $\widehat{\text{Cov}}(\hat{x}, e)$) in each of the DGPs. Depending upon the DGP, Lasso, neural nets, and boosted trees generate the second-highest covariance. Because our *shallow* subcase of neural nets approximates OLS, its covariance between \hat{x} and e is approximately zero.

One way to ensure that $\widehat{\text{Cov}}(\hat{x}, e) = 0$ for a nonlinear model is to linearize its output—this is achievable, for example, by using the ML-based prediction $\hat{x}(\mathbf{z})$ as an *instrument* for *x*, rather than

¹⁴ We assume predictions, \hat{x} , will have non-negative covariance with the true values, x.

¹⁵ We assume e, x, and \hat{x} have variation and that the predictions \hat{x} positively correlate with the true values x.

¹⁶ Further, this assumption is equivalent to making an assumption on $\widehat{\text{Cov}}(\hat{x}, e)$, which means one is essentially assuming the result. That is, $\widehat{\text{Var}}(x) = \widehat{\text{Var}}(\hat{x}) + \widehat{\text{Var}}(e) + 2 \widehat{\text{Cov}}(\hat{x}, e)$. That said, in every iteration of our simulations, $\widehat{\text{Var}}(x) > \widehat{\text{Var}}(\hat{x})$.

¹⁷ This upward bias is partly due to the true parameter β_1 being positive.

plugging it into the second stage (Angrist and Krueger, 2001; Chen et al., 2020). While this approach forces $\widehat{\text{Cov}}(\hat{x}, e) = 0$, it requires strengthening assumptions. We discuss this possibility further in Section 5.

More broadly, the component of bias due to covariance between first-stage predictions (\hat{x}) and their residuals (e) (*i.e.*, the $\widehat{\text{Cov}}(\hat{x}, e)$ term) accounts for the vast majority of the bias for Lasso and substantial amounts of the bias in random forests, boosted trees, and neural nets (the exact portion of the bias differs across DGPs and iterations). While $\widehat{\text{Cov}}(\hat{x}, e)$ does not account for all of the bias, the non-zero covariance between first-stage predictions and residuals is an important (potentially large) component of the bias of ML-based 2SLS models.

The $\widehat{\text{Cov}}(\hat{x}, u)$ component

Unlike $\widehat{\text{Cov}}(\hat{x}, e)$, the second component of the wedge between $\hat{\beta}_1^{2\text{SLS}}$ and β_1 can be non-zero for both OLS-based methods and non-OLS models. However, methods that use *non-linear* predictions of *x* in the first stage (*i.e.*, ML-assisted 2SLS) require special care to reduce $\widehat{\text{Cov}}(\hat{x}, u)$ and produce low-bias estimates of β_1 .

This second term, $\widehat{\text{Cov}}(\hat{x}, u)$, is effectively the exclusion restriction, and any 2SLS-inspired estimator can reduce bias in $\hat{\beta}_1$ by ensuring $\widehat{\text{Cov}}(\hat{x}, u)$ is approximately zero. Assuming the instruments \mathbf{z} are valid, an arbitrary prediction algorithm can maintain $\widehat{\text{Cov}}(\hat{x}, u) \approx 0$ through either of three conditions:

- 1. Restrict the algorithm's choice set: By restricting the learning algorithm to choosing from a set/class of functions where each individual function satisfies the exclusion restriction, one mechanically ensures the first-stage predictions \hat{x} do not co-vary with the unobserved disturbance *u*. For example, traditional OLS-based 2SLS chooses the first stage from the set of linear combinations of the instruments. As long as the instruments are uncorrelated with the disturbance, any element of this set (the linear combination of the instruments) will also be uncorrelated with the disturbance—satisfying the exclusion restriction.
- 2. Extend the exclusion restriction: One may extend the commonly invoked "uncorrelated" assumption to the stronger and underlying exclusion restriction of conditional mean independence. Rather than only assuming all *linear* combinations of the instruments are exogenous, one could assume that **all functions** of the instruments (linear and nonlinear) satisfy the exclusion restriction. Put differently, this condition requires $\widehat{Cov}(f(\mathbf{z}), u) \approx 0$ for any function *f*.
- 3. **Really** *trust* **the ML algorithm:** The final option is to simply rely upon the algorithm to find a function that satisfies the exclusion restriction, irrespective of choice set—something akin to closing one's eyes and hoping for the best. While this option makes a heroic assumption, as ML algorithms are not typically designed to search for and find valid exclusion restrictions, it is the default scenario. If a practitioner does not enforce condition **1** and does not assume

condition **2**, then she is left hoping that the ML methods successfully chose a function that includes a valid exclusion restriction.

Put simply, sufficiently flexible learning algorithms can recover endogenous variation in *x*—even when *using linearly valid instruments*. Many ML training methods explicitly incentivize and enable algorithms to do this.

Column *b* of Table 2 highlights the tendency of flexible first-stage models (*e.g.*, tree methods and neural nets) to recover endogeneity. As the flexibility of algorithms increases, $\widehat{\text{Cov}}(\hat{x}, u)$ tends to increase as well (across all DGPs). This covariance and its associated bias are particularly large for tree-based methods (especially random forests) and neural nets with multiple hidden layers. Notably, in panels 2B–2D of Figure 2, the densities of unrestricted and narrow neural networks are bimodal. As Appendix Figure A3 illustrates, the bimodality results from whether the neural network (*i*) "chooses" zero hidden layers (the less biased mode) or (*ii*) goes deeper (learning the endogenous error and generating more bias).^{18,19} This covariance between predictions \hat{x} and the unobserved disturbance *u* accounts for a substantial amount of the bias in nonlinear methods, which demonstrates that the previously discussed first component $\widehat{\text{Cov}}(\hat{x}, e)$ is not the only issue facing these models.

The $\frac{1}{\operatorname{Var}(\hat{x})}$ component

While the first two bias components enter additively, the third component scales their sum. Any method that reduces the variance of the first-stage predictions (*i.e.*, reduces $\widehat{Var}(\hat{x})$) mechanically inflates the bias produced by $\beta_1 \widehat{Cov}(\hat{x}, e) + \widehat{Cov}(\hat{x}, u)$.

In the case of properly specified, OLS-based 2SLS, the variance of the predictions hardly affects bias in β_1 , since $\widehat{\text{Cov}}(\hat{x}, e) = 0$ and $\widehat{\text{Cov}}(\hat{x}, u) \approx 0$. However, most ML algorithms reduce the variance of their predictions while trading between out-of-sample bias and variance. This tradeoff between bias and variance happens *outside of a 2SLS framework*. Consequently, when practitioners insert variance-reducing ML methods into 2SLS, the variance reduction actually amplifies bias in the second-stage estimates.

Taking these insights to the results in Panel B of Table 2, notice that variance reduction can cause methods to perform poorly. For example, Lasso-assisted 2SLS produces the lowest variance \hat{x} in two of the three high-complexity cases. (In cases 2 and 3, Lasso-based 2SLS has the highest

¹⁸ This result highlights the importance of allowing neural networks to choose no hidden layers.

¹⁹ Another related concern familiar to the ML literature is *overfit*. Overfit models tend to produce larger values of $\widehat{Cov}(\hat{x}, u)$ than models that have been cross-validated. Though cross-validation is best/standard practice for machinelearning methods in prediction problems, here it retains importance by preventing the algorithms from overfitting the target variable *x* in the first stage (even when out-of-sample performance is no longer the goal). We use five-fold cross-validation (CV) to tune the hyperparameters for Lasso-, tree-, and neural-net-based methods. Our neural-net cross-validation departs from standard five-fold CV. In Appendix Section A.3 we detail our cross-validation process for training neural net. One might further avoid overfit by applying holdout-style methods—only generating predictions for observation *i* when *i* is not in the training set. JIVE, split-sample IV, and Chen et al. (2020) all feature this additional safeguard. We do not employ these holdout-based methods because our goal in this paper is to simulate the results of a researcher using off-the-shelf ML tools in the first stage of 2SLS.

 $1/Var\hat{x}$ -based amplifier of the bias.) This high degree of bias amplification generates notable bias in Lasso relative to many other methods. (This is also evident in Figure 2). So while Lasso's $\widehat{Cov}(\hat{x}, e)$ and $\widehat{Cov}(\hat{x}, u)$ are less than or equal to those of many other methods, the amplification produced by variance-reduction in \hat{x} ultimately causes Lasso to have substantial bias. Notably, post-Lasso-based 2SLS produces less bias, partly due to the fact that it includes less variance reduction.

Worse yet, tree-based methods substantially reduce variance *and* produce relatively large $\widehat{Cov}(\hat{x}, e)$ and $\widehat{Cov}(\hat{x}, u)$ components, resulting in very large bias in their parameter estimates (even larger than naïve OLS).

5 Discussion

In this paper, we examine the implications of plugging off-the-shelf ML methods into a 2SLS framework. In many cases, injecting ML into the first stage of 2SLS generates substantial bias.

While there are many approaches to combining instrumental-variable intuition and machine learning, they relax the traditional 2SLS structure and require different (generally stronger) identifying assumptions and/or tailor-made ML algorithms.²⁰ Among the current options, the closest in spirit to our question of "What are the implications of inserting ML into 2SLS?" is the "machine learning split-sample" (MLSS) estimator proposed by Chen et al. (2020).²¹

With two fairly simple expansions of the traditional 2SLS framework, MLSS mitigates many biases generated by naïvely plugging ML methods into the first stage. However, as with other more ML-forward methods, the solution is not without the cost of substantially strengthening the exclusion restriction. Specifically, Chen et al. (2020) proposes augmenting 2SLS with two simple techniques: restrict ML-based predictions to be explicitly out of sample (using split-sample methods) and use the ML-generated predictions as a "synthetic" instrument that then enters linearly in the first stage.

The idea for out-of-sample (split-sample) ML predictions follows the lead of Jackknife IV and Split Sample IV. By introducing out-of-sample methods to the ML-prediction exercise, Chen et al. aims to prevent the ML algorithm from fitting the first-stage errors, and shutting down the bias generated by $\widehat{Cov}(\hat{x}, u)$. One potential drawback, however, is that this out-of-sample step likely increases variability (as seen in the JIVE results of Figure 2).

The second component of Chen et al. involves a *zeroth stage* (i.e., before the first stage), in which the practitioner trains an ML algorithm to predict *x* using the instruments z.²² The predictions from

²⁰ For example, MLSS (Chen et al., 2020), DeepIV (Hartford et al., 2017), DeepGMM (Bennett et al., 2019), KIV (Singh et al., 2019), Adversarial Estimation of Riesz Representers (Chernozhukov et al., 2020), Neural Estimation of SEM (Liao et al., 2020), and Non-Parametric IV (Kilbertus et al., 2020).

²¹ Angrist and Frandsen (2022) also applies split-sample methods to several ML algorithms (*i.e.*, post-Lasso, random forest), both in the first stage of 2SLS and while synthesizing instruments in a stage that precedes the first stage.

 $^{^{22}}$ Note that this zero $^{\rm th}$ stage is identical to first stages that naïvely insert ML methods into 2SLS.

this zeroth stage are then used as the instrument within a traditional 2SLS framework. The benefit is that the resulting linear first stage (linearizing the results of a potentially forbidden regression) guarantees that $\widehat{\text{Cov}}(\hat{x}, e) = 0$ and shuts down one avenue through which bias enters.

Importantly, this zeroth stage of the MLSS approach requires that no learnable function of instruments meaningfully predicts the structural disturbance u. The function spaces of ML algorithms can cover all possible functions of the instruments (*e.g.*, most neural networks are universal approximators), which requires strengthening the exclusion restriction from the assumption of "no correlation" to the actual underlying assumption of conditional mean-independence between the instruments and disturbance. For example, this strengthening includes ML-learned step functions, threshold indicators, kinks, many-way interactions—*any* learnable function f(z) that is informative for x. This potentially infinite-dimensional set of exclusion restrictions is presumably more difficult to justify than the typical identifying assumption assumed in 2SLS applications.

Overall, solutions exist for the bias components that we identify in this paper, but each solution comes at a cost. Some solutions are fairly cheap (*e.g.*, out-of-sample prediction methods require sufficient overlap and can increase noise). Other solutions require more from the practitioner, *e.g.*, strengthening a linear exclusion restriction to an infinite-dimensional function space. Researchers may be comfortable invoking these stronger identifying assumptions in some settings. However, more work is needed in exploring and bounding the relative benefits and costs of these procedures in practice.²³

6 Conclusion

Our results show that naïvely inserting machine-learning algorithms into the first stage of a 2SLS structure often produces bias. While some channels of this bias will be familiar to an economics audience, others may be less familiar, and arise from the intersection of "classical" econometric tools with ML methods. In terms of minimum bias, we find that the most-successful cases of integrating machine learning into 2SLS are largely restricted to instrument selection or modification (*e.g.*, post-Lasso) before entering a traditional 2SLS framework. However, inserting highly nonlinear algorithms (*e.g.*, random forests, boosted trees, and relatively deep neural nets) into the first stage of 2SLS drives estimates of causal parameters to be *more* biased than much simpler and more-easily interpreted alternatives. Importantly, such methods can even yield more bias than a naïve, endogenous OLS regression.

We mechanically decompose the bias within ML-augmented 2SLS estimates into three terms:

• $\beta \operatorname{Cov}(\hat{x}, e)$: When first-stage predictions are not orthogonal to their residuals (relevant for nonlinear methods), $\operatorname{Cov}(\hat{x}, e) \neq 0$. This wedge increases with the magnitude of the target parameter, β .

²³ For example, future work could continue to explore the bias from instruments that nonlinearly violate an expanded exclusion restriction (potentially discoverable by nonlinear ML algorithms) but satisfy the traditional linear exclusion restriction. Future work could also consider how ML-based 2SLS approaches estimate potentially different LATEs.

- $\widehat{\text{Cov}}(\hat{x}, u)$: Flexible ML methods can overfit *x* in the first stage, unintentionally recovering endogenous variation in *x*. This *overfit* first stage produces fitted values (\hat{x}) that are no longer exogenous from the structural disturbance *u*, even when the first-stage instruments are exogenous.
- $1/Var(\hat{x})$: ML-based prediction methods typically reduce the variance of their predictions. As such, ML-augmented 2SLS often reduces the variance of \hat{x} , which amplifies any non-zero wedge between the estimate and the estimand coming from $\beta Cov(\hat{x}, e) + Cov(\hat{x}, u)$.

Overall, we find that many of the methods pioneered in the 1990s and 2000s (*e.g.*, split-sample IV, LIML, and JIVE) reduce bias resulting from weak or over-identified IV for a linear first stage.²⁴ At the same time, we find that less-traditional methods that *mindfully* incorporate machine learning techniques (*e.g.*, PCA-synthesized instruments and post-Lasso) generate improvements over 2SLS across a variety of DGPs.²⁵ Importantly, we show that carelessly injecting ML into the first stage can also produce substantial bias, possibly worsening bias over naïve OLS.

The fundamental problem comes from recognizing the first stage of 2SLS solely as a prediction problem without recognizing that it is also part of a larger estimation *system*.

ML methods can produce good predictions. However, they were neither developed nor optimized for two-stage procedures that generate minimally biased causal estimates. For example, variance reduction in prediction typically improves out-of-sample prediction accuracy but can actually inflate parameter bias in the second stage of 2SLS. Solutions have been proposed to these and other issues from ML-based 2SLS. However, if one approaches these problems without the caution that we suggest (*e.g.*, blindly relying upon ML to learn the variation in the first stage) bias often results. Unfortunately, the solutions to these ML-in-2SLS problems can often require the practitioner to strengthen the exclusion restriction where complexity is highest (*e.g.*, in high-dimensional data). Because ML methods are typically adopted in complex data environments that are not well understood, employing these methods without safeguards can entice researchers into making heroic assumptions.

Given the performance of existing off-the-shelf estimators, naïvely injecting ML into 2SLS appears to produce little gain relative to its costs/biases. More broadly, applying ML methods to 2SLS requires the practitioner to face issues present in both the prediction *and* causal-inference settings—and in addition, less-familiar issues arising from their interaction.

²⁴ Angrist and Frandsen (2022) comes to a similar conclusion. Focusing more on heterogeneous treatment effects and ML-assisted covariate selection, Angrist and Frandsen does not decompose the bias associated with ML methods applied to the instruments in the first-stage prediction problem (Section 4.2 and summarized above). Others have also found that JIVE and LIML are the best choice in the weak-many-instruments case—Hansen and Kozbur (2014) explores a regularized JIVE, and Carrasco and Tchuente (2015) develops a regularized LIML approach.

²⁵ See Ackerberg and Devereux (2009) for a discussion of tools for addressing over-identified instrumental variables estimation and Andrews et al. (2019) for a discussion of weak instruments in linear IV regression.

Figures

Figure 1: The "high-complexity" DGP: 100 instruments of varying strength

Order of instruments' coefficients (π)



The top panels—1A, 1B, and 1C—illustrate the instruments' coefficients (π) in the DGP for the first stage of each of the three subcases (*i.e.*, how the instruments **z** relate to *x* in each subcase). While the three subcases differ in their first-stage coefficients, they share the same variance-covariance structure Σ_z , which 1D depicts ($Var(z_i) = 1, \forall i$). Figure 1E individually plots the correlation between two instruments (z_1 and z_{50}) and all of the other instruments. The instruments z_1 and z_{50} are specifically of interest because they are the strongest instruments in subcases 2 and 3, respectively.

19

Figure 2: Main results— $\hat{\beta}$ distributions across competing two-stage methods



(B) High-complexity case 1: 100 instruments; coefficients 'shuffled'

Each individual distribution represents the density of estimates from 1,000 iterations of simulation. The true value of the target parameter β is 1 (the dashed vertical line). The DGP underlying subfigure 2A uses 7 strong and exogenous instruments. For subfigures 2B, 2C, and 2D, the models have access to 100 exogenous instruments of varying strengths. The three 'high-complexity' cases differ in their DGPs' first-stage coefficients but share a common variance structure among the 100 instruments (Σ_{τ})—following and extending Belloni et al. (2012). The Oracle model refers to a model using only the exogenous portion of x (the endogenous regressor). In Naive OLS, we estimate the parameter without regard for the endogeneity of x. The remaining methods each estimate the target parameter with a different variety of approaches toward 2SLS or related methods. Table 1 contains the mean and standard deviation for each distribution.



Figure 3: Exclusion-restriction violations via higher-order interactions among instruments ('low-complexity case' of 7 strong instruments)

(B) 'Selection' methods: Lasso, post-Lasso, and PCA

(A) 'Standard' linear estimators: OLS, LIML, SSIV, and JIVE

The DGPs underlying these figures add a *k*-term interaction between the instruments to the structural error. *E.g.*, when the *x*-axis equals 3, we add the interaction $x_1 \times x_2 \times x_3$ to the error; k = 1 is a linear exclusion-restriction violation. Interactions with k > 1 do not violate the exclusion restriction of linear methods but do violate the exclusion restriction for methods that require a higher-order/expanded exclusion restriction (*i.e.*, conditional independence). Each subfigure illustrates the performance of the given estimators from 1,000 iterations. The solid lines and dots mark the mean of the estimator-DGP combination; the shaded bands give the boundaries between the 2.5th and 97.5th percentiles. The true value is 1. *Normalized* methods ('norm.') use the stated method to synthesize a single instrument from the seven instruments *prior to the first stage* (as in Chen et al. (2020)). See Figure A2 for the densities of 2-, 3-, 4-, and 5-way interactions.

7 Tables

	Low-complexity case High-complexity cases					
	7 strong instruments (A)	100 m (B)	nents (D)			
Naive OLS	1.038	1.335	1.334	1.223		
	(0.007)	(0.020)	(0.019)	(0.046)		
First stage: OLS	1.000	1.058	1.056	1.023		
	(0.007)	(0.040)	(0.039)	(0.042)		
LIML (Fuller)	1.000	1.000	0.998	1.000		
	(0.008)	(0.044)	(0.043)	(0.043)		
Split-sample IV	1.001	1.062	1.060	1.025		
	(0.007)	(0.041)	(0.040)	(0.043)		
Jackknife IV (JIVE)	1.000	0.998	0.996	1.000		
	(0.017)	(0.044)	(0.044)	(0.044)		
First stage: PCA	1.000	1.032	1.026	1.016		
	(0.007)	(0.042)	(0.041)	(0.045)		
First stage: Post-Lasso selection	1.000	1.026	1.023	1.013		
	(0.007)	(0.044)	(0.042)	(0.042)		
First stage: Lasso	1.007	1.100	1.098	1.042		
	(0.007)	(0.045)	(0.045)	(0.045)		
First stage: Neural net	1.008	1.215	1.209	1.110		
	(0.029)	(0.180)	(0.176)	(0.105)		
First stage: Neural net, shallow	1.002	1.069	1.065	1.030		
	(0.018)	(0.066)	(0.067)	(0.049)		
First stage: Neural net, narrow	1.008	1.213	1.210	1.100		
	(0.027)	(0.183)	(0.185)	(0.103)		
First stage: Boosted trees	1.008	1.254	1.255	1.121		
	(0.007)	(0.041)	(0.039)	(0.047)		
First stage: Random forest, CV	1.071	1.562	1.563	1.316		
	(0.008)	(0.033)	(0.034)	(0.058)		

Table 1: Simulation results: Mean and standard deviation for methods and DGPs

This table provides the means and standard deviations of the distributions illustrated in Figure 2. Each **column** contains a separate DGP: (a) contains the *low-complexity* DGP with 7 (equally) strong instruments; (b)–(d) contain the three *high-complexity* cases with 100 instruments of mixed strengths. **Rows** differ by estimator. For each DGP-estimator combination, we summarize the estimates for the parameter of interest (β) across 1,000 iterations using a mean and standard deviation (the standard deviation is in parentheses).

 Table 2: Simulation results: Decomposing bias components (means from simulation)

Bias components									
	(a+b)c	a	b	с					
Estimator	Bias	$\widehat{\text{Cov}}(\hat{x}, e)$	$\widehat{\operatorname{Cov}}(\hat{x}, u)$	$1/\widehat{\operatorname{Var}}(\hat{x})$	$\widehat{\operatorname{Var}}(\hat{x})$	$\widehat{\operatorname{Var}}(x)$	$\widehat{\text{Cov}}(x, \hat{x})$	$\operatorname{Corr}(x,\hat{x})$	$\widehat{\operatorname{Cov}}(x, u)$
(A) Low-complexity case									
Naive OLS	0.04	0	1.01	0.04	26.41	26.41	26.41	1	1.01
First stage: OLS	0	0	0.01	0.04	25.41	26.41	25.41	0.98	1.01
Split-sample IV	0	0	0.02	0.04	25.42	26.41	25.42	0.98	1.01
Jackknife IV (JIVE)	0	0	0.01	0.05	20.76	21.74	20.76	0.98	1.01
First stage: PCA	0	0	0.01	0.04	25.41	26.41	25.41	0.98	1.01
First stage: Post-Lasso selection	0	0	0.01	0.04	25.41	26.41	25.41	0.98	1.01
First stage: Lasso selection	0.01	0.17	0.01	0.04	25.08	26.41	25.25	0.98	1.01
First stage: Neural net	0.01	0.1	0.05	0.05	20.52	21.67	20.61	0.98	0.99
First stage: Neural net, narrow	0.01	0.09	0.04	0.05	20.52	21.67	20.61	0.98	0.99
First stage: Neural net, shallow	0	0	0.03	0.05	20.71	21.67	20.71	0.98	0.99
First stage: Boosted trees	0.01	0.01	0.2	0.04	25.51	26.41	25.52	0.98	1.01
First stage: Random forest, CV	0.07	1.08	0.62	0.04	24.01	26.41	25.09	1	1.01
(B) High-complexity case 1									
Naive OLS	0.22	0	0.15	1.56	0.64	0.64	0.64	1	0.15
First stage: OLS	0.02	0	0.01	1.77	0.56	0.64	0.56	0.94	0.15
Split-sample IV	0.03	0	0.01	1.79	0.56	0.64	0.56	0.93	0.15
Jackknife IV (JIVE)	0	0	0	1.77	0.57	0.65	0.57	0.94	0.15
First stage: PCA	0.02	0	0.01	2.02	0.5	0.64	0.5	0.88	0.15
First stage: Post-Lasso selection	0.01	0	0.01	1.79	0.56	0.64	0.56	0.93	0.15
First stage: Lasso selection	0.04	0.02	0	1.92	0.52	0.64	0.54	0.93	0.15
First stage: Neural net, shallow	0.03	0	0.02	1.77	0.57	0.64	0.57	0.94	0.15
First stage: Neural net, narrow	0.1	0.01	0.05	1.75	0.57	0.64	0.58	0.96	0.15
First stage: Neural net	0.11	0.01	0.06	1.74	0.58	0.64	0.58	0.96	0.15
First stage: Boosted trees	0.12	0.03	0.04	1.91	0.52	0.65	0.55	0.95	0.15
First stage: Random forest, CV	0.32	0.06	0.09	2.04	0.49	0.65	0.56	0.99	0.15
(C) High-complexity case 2									
Naive OLS	0.33	0	0.35	0.95	1.06	1.06	1.06	1	0.35
First stage: OLS	0.06	0	0.03	1.66	0.6	1.06	0.6	0.76	0.35
Split-sample IV	0.06	0	0.03	1.76	0.57	1.06	0.57	0.73	0.35
Jackknife IV (JIVE)	0	0	0	1.65	0.61	1.06	0.61	0.76	0.35
First stage: PCA	0.03	0	0.02	1.75	0.57	1.06	0.57	0.74	0.35
First stage: Post-Lasso selection	0.02	0	0.01	1.75	0.57	1.06	0.57	0.74	0.35
First stage: Lasso selection	0.1	0.04	0.01	2.11	0.48	1.06	0.52	0.73	0.35
First stage: Neural net	0.21	0.03	0.13	1.49	0.7	1.06	0.73	0.84	0.35
First stage: Neural net, narrow	0.21	0.04	0.12	1.52	0.68	1.06	0.71	0.84	0.35
First stage: Neural net, shallow	0.06	0	0.04	1.63	0.62	1.06	0.62	0.76	0.35
First stage: Boosted trees	0.25	0.07	0.07	1.83	0.55	1.06	0.62	0.81	0.36
First stage: Random forest, CV	0.56	0.16	0.22	1.51	0.66	1.06	0.82	0.98	0.35
(D) High-complexity case 3									
Naive OLS	0.34	0	0.35	0.95	1.06	1.06	1.06	1	0.35
First stage: OLS	0.06	0	0.04	1.66	0.61	1.06	0.61	0.76	0.35
Split-sample IV	0.06	0	0.04	1.76	0.57	1.06	0.57	0.73	0.35
Jackknife IV (JIVE)	0	0	0	1.64	0.61	1.06	0.61	0.76	0.36
First stage: PCA	0.03	0	0.02	1.76	0.57	1.06	0.57	0.73	0.35
First stage: Post-Lasso selection	0.03	0	0.02	1.74	0.58	1.06	0.58	0.74	0.35
First stage: Lasso selection	0.1	0.04	0.01	2.1	0.48	1.06	0.52	0.73	0.35
First stage: Neural net	0.22	0.03	0.13	1.49	0.69	1.06	0.73	0.84	0.36
First stage: Neural net, narrow	0.21	0.04	0.12	1.53	0.67	1.06	0.71	0.84	0.36
First stage: Neural net, shallow	0.07	0	0.05	1.63	0.62	1.06	0.62	0.76	0.36
First stage: Boosted trees	0.25	0.07	0.07	1.83	0.55	1.06	0.62	0.81	0.36
First stage: Random forest, CV	0.56	0.16	0.22	1.51	0.66	1.06	0.82	0.98	0.35

A cell's value provides the given statistic's mean across 1,000 iterations. The panels (A–D) denote the four separate DGPs; rows refer to estimators. We omit LIML as it is not a two-stage method and thus does not produce \hat{x} .

References

- Ackerberg, D. A. and Devereux, P. J. (2009). Improved JIVE Estimators for Overidentified Linear Models with and without Heteroskedasticity. *Review of Economics and Statistics*, 91(2):351–362.
- Anderson, T. W. and Rubin, H. (1949). Estimation of the Parameters of a Single Equation in a Complete System of Stochastic Equations. *The Annals of Mathematical Statistics*, 20(1):46–63.
- Andrews, I., Stock, J., and Sun, L. (2019). Weak Instruments in IV Regression: Theory and Practice. Annual Review of Economics, 11:727–753.
- Angrist, J. and Pischke, J.-S. (2009). *Mostly Harmless Econometrics: An Empiricist's Companion*. Princeton University Press, Princeton, New Jersey.
- Angrist, J. D. and Frandsen, B. (2022). Machine Labor. Journal of Labor Economics, 40(S1):S97–S140.
- Angrist, J. D., Imbens, G. W., and Krueger, A. B. (1999). Jackknife instrumental variables estimation. *Journal of Applied Econometrics*, 14(1):57–67.
- Angrist, J. D. and Krueger, A. B. (1995). Split-Sample Instrumental Variables Estimates of the Return to Schooling. *Journal of Business & Economic Statistics*, 13(2):225–235.
- Angrist, J. D. and Krueger, A. B. (2001). Instrumental Variables and the Search for Identification: From Supply and Demand to Natural Experiments. *Journal of Economic Perspectives*, 15(4):69–85.
- Belloni, A., Chen, D., Chernozhukov, V., and Hansen, C. (2012). Sparse Models and Methods for Optimal Instruments With an Application to Eminent Domain. *Econometrica*, 80(6):2369–2429.
- Belloni, A., Chernozhukov, V., and Hansen, C. (2011). Lasso Methods for Gaussian Instrumental Variables Models. *arXiv preprint arXiv:1012.1297*.
- Belloni, A., Chernozhukov, V., and Hansen, C. (2013). Inference on Treatment Effects after Selection among High-Dimensional Controls. *The Review of Economic Studies*, 81(2):608–650.
- Bennett, A., Kallus, N., and Schnabel, T. (2019). Deep generalized method of moments for instrumental variable analysis. *Advances in Neural Information Processing Systems*, 32.
- Bevis, L. E. and Villa, K. (2020). Intergenerational Transmission of Maternal Health: Evidence from Cebu, the Philippines. *Journal of Human Resources*, pages 0819–10372R2.
- Biewen, M. and Kugler, P. (2021). Two-stage least squares random forests with an application to Angrist and Evans (1998). *Economics Letters*, 204:109893.
- Breiman, L. (1997). Arcing the edge. Technical report, Technical Report 486, Statistics Department, University of California at Berkeley.
- Breiman, L. (2001). Random Forests. Machine Learning, 45(1):5–32.
- Carrasco, M. and Tchuente, G. (2015). Regularized LIML for many instruments. *Journal of Econometrics*, 186(2):427–442.
- Chen, D. L. and Yeh, S. (2020). Government Expropriation Increases Economic Growth and Racial Inequality: Evidence from Eminent Domain. TSE Working Paper No. 16-693.

- Chen, J., Chen, D. L., and Lewis, G. (2020). Mostly Harmless Machine Learning: Learning Optimal Instruments in Linear IV Models. *arXiv preprint arXiv:2011.06158*.
- Chen, W., Chen, X., Hsieh, C.-T., and Song, Z. (2019). A Forensic Examination of China's National Accounts. *Brookings Papers on Economic Activity*, 2019(1):77–141.
- Chernozhukov, V., Chetverikov, D., Demirer, M., Duflo, E., Hansen, C., Newey, W., and Robins, J. (2018). Double/debiased machine learning for treatment and structural parameters. *The Econometrics Journal*, 21(1):C1–C68.
- Chernozhukov, V., Hansen, C., and Spindler, M. (2015). Valid Post-Selection and Post-Regularization Inference: An Elementary, General Approach. *Annual Review of Economics*, 7(1):649–688.
- Chernozhukov, V., Newey, W., Singh, R., and Syrgkanis, V. (2020). Adversarial estimation of riesz representers. *arXiv preprint arXiv:2101.00009*.
- Derenoncourt, E. (2022). Can You Move to Opportunity? Evidence from the Great Migration. *American Economic Review*, 112(2):369–408.
- Farley, B. and Clark, W. (1954). Simulation of self-organizing systems by digital computer. *Transactions of the IRE Professional Group on Information Theory*, 4(4):76–84.
- Friedman, J. H. (2001). Greedy function approximation: A gradient boosting machine. *Annals of Statistics*, pages 1189–1232.
- Friedman, J. H. (2002). Stochastic gradient boosting. *Computational Statistics & Data Analysis*, 38(4):367–378.
- Fuller, W. A. (1977). Some Properties of a Modification of the Limited Information Estimator. *Econometrica*, 45(4):939–953.
- Hansen, C., Hausman, J., and Newey, W. (2008). Estimation With Many Instrumental Variables. Journal of Business & Economic Statistics, 26(4):398–422.
- Hansen, C. and Kozbur, D. (2014). Instrumental variables estimation with many weak instruments using regularized JIVE. *Journal of Econometrics*, 182(2):290–308.
- Hartford, J., Lewis, G., Leyton-Brown, K., and Taddy, M. (2017). Deep IV: A Flexible Approach for Counterfactual Prediction. In Precup, D. and Teh, Y. W., editors, *Proceedings of the 34th International Conference on Machine Learning, PMLR*, volume 70, pages 1414–1423.
- Hastie, T., Tibshirani, R., and Friedman, J. (2009). The Elements of Statistical Learning. Springer New York.
- Heckman, J. J., Urzua, S., and Vytlacil, E. (2006). Understanding instrumental variables in models with essential heterogeneity. *The Review of Economics and Statistics*, 88(3):389–432.
- Heckman, J. J. and Vytlacil, E. (2005). Structural equations, treatment effects, and econometric policy evaluation. *Econometrica*, 73(3):669–738.
- Ho, T. K. (1995). Random Decision Forests. In Proceedings of the Third International Conference on Document Analysis and Recognition (Volume 1), pages 278–282. IEEE Computer Society.

- Ioffe, S. and Szegedy, C. (2015). Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning*, pages 448–456. PMLR.
- James, G., Witten, D., Hastie, T., and Tibshirani, R. (2013). *An Introduction to Statistical Learning*. Springer New York.
- Kilbertus, N., Kusner, M. J., and Silva, R. (2020). A class of algorithms for general instrumental variable models. Advances in Neural Information Processing Systems, 33:20108–20119.
- Kingma, D. P. and Ba, J. (2017). Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Lee, D. S., McCrary, J., Moreira, M. J., and Porter, J. (2022). Valid t-ratio Inference for IV. American Economic Review, 112(10):3260–3290.
- Liao, L., Chen, Y.-L., Yang, Z., Dai, B., Kolar, M., and Wang, Z. (2020). Provably efficient neural estimation of structural equation models: An adversarial approach. *Advances in Neural Information Processing Systems*, 33:8947–8958.
- Liu, R., Shang, Z., and Cheng, G. (2020). On deep instrumental variables estimate. *arXiv preprint arXiv:2004.14954*.
- Mason, L., Baxter, J., Bartlett, P., and Frean, M. (1999). Boosting algorithms as gradient descent. Advances in Neural Information Processing Systems, 12:512–518.
- McCulloch, W. S. and Pitts, W. (1943). A logical calculus of the ideas immanent in nervous activity. *The bulletin of mathematical biophysics*, 5(4):115–133.
- Mueller-Smith, M. (2015). The criminal and labor market impacts of incarceration. Unpublished.
- Mullainathan, S. and Spiess, J. (2017). Machine Learning: An Applied Econometric Approach. *Journal of Economic Perspectives*, 31(2):87–106.
- Ng, S. and Bai, J. (2009). Selecting Instrumental Variables in a Data Rich Environment. *Journal of Time Series Econometrics*, 1(1).
- Pearson, K. (1901). LIII. on lines and planes of closest fit to systems of points in space. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science,* 2(11):559–572.
- Ribeiro, M. T., Singh, S., and Guestrin, C. (2016). "Why Should I Trust You?". In Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. ACM.
- Santosa, F. and Symes, W. W. (1986). Linear Inversion of Band-Limited Reflection Seismograms. *SIAM Journal on Scientific and Statistical Computing*, 7(4):1307–1330.
- Singh, A., Hosanagar, K., and Gandhi, A. (2020). Machine Learning Instrument Variables for Causal Inference. In *Proceedings of the 21st ACM Conference on Economics and Computation*. ACM.
- Singh, R., Sahani, M., and Gretton, A. (2019). Kernel instrumental variable regression. *Advances in Neural Information Processing Systems*, 32.
- Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. (2014). Dropout: A Simple Way to Prevent Neural Networks from Overfitting. *The Journal of Machine Learning Research*, 15(1):1929–1958.

- Storm, H., Baylis, K., and Heckelei, T. (2019). Machine learning in agricultural and applied economics. *European Review of Agricultural Economics*, 47(3):849–892.
- Tibshirani, R. (1996). Regression Shrinkage and Selection Via the Lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288.
- Turing, A. (2004). Intelligent Machinery (1948), pages 395-432. Oxford University Press.
- Wang, S., Wang, Q., and Zhao, J. (2020). Deep neural networks for choice analysis: Extracting complete economic information for interpretation. *Transportation Research Part C: Emerging Technologies*, 118.
- Winkelried, D. and Smith, R. (2011). Principal Components Instrumental Variable Estimation. Cambridge Working Papers in Economics: *CWPE1119*.
- Wooldridge, J. (2010). Econometric Analysis of Cross Section and Panel Data. MIT Press, Cambridge, Mass.
- Xu, L., Chen, Y., Srinivasan, S., de Freitas, N., Doucet, A., and Gretton, A. (2020). Learning deep features in instrumental variable regression. *arXiv preprint arXiv:2010.07154*.
- Zhao, S., Witten, D., and Shojaie, A. (2021). In defense of the indefensible: A very naive approach to high-dimensional inference. *Statistical science: a review journal of the Institute of Mathematical Statistics*, 36(4):562.

Appendix A

A.1 Appendix: Figures

Figure A1: Predictions vs. estimation: Comparing cross-validated prediction performance with bias in random-forest-based 2SLS



(A) In- and out-of-sample MSE for predictions of *x*













(A) Two-term 'exclusion-restriction violation': $x_1 \times x_2$

(B) Three-term 'exclusion-restriction violation': $x_1 \times x_2 \times x_3$

This figure illustrates the densities of the estimates portrayed in Figure 3. Post-lasso's high-variance in this experiment come from using the so-called plug-in lambda value that prevents overselection of poor-performing instruments. To calculate this plug-in penalty, two parameters are user-chosen: c and γ and multiple theoretical values are used. We used the suggested values as detailed in Belloni et al. (2012) and Belloni et al. (2013). This lambda plug-in value is calculated to be equal to $(1.1 * 2) \sqrt{N} * \Phi(\frac{1-(\gamma)}{2nz_c})$ and $\gamma \equiv \frac{1}{ln(N \lor nz)}$. In our case, the number of instruments nz = 7, the number of observations N = 1000, and $\Phi \sim Quantile Normal$. This results in $\lambda \approx 36$. By using the suggested higher levels of λ , post-lasso can result in a higher chance of under-fit, with the extreme example being a second-stage regression model of the form f(X|Z) = 0. Given our variables are independent and relatively strong, regularization is highly unlikely to improve fit, and large lambda may result in under-specification for the first stage.

A.2 Appendix: Math

A.2.1 Wedge *a* covariance and correlation

Recall that $e = x - \hat{x}$, *i.e.*, *e* is the first-stage prediction's residual.

$$\operatorname{Corr}(\hat{x}, e) = \frac{\operatorname{Cov}(\hat{x}, e)}{\sigma_{\hat{x}}\sigma_{e}}$$
$$= \frac{\operatorname{Cov}(\hat{x}, x)}{\sigma_{\hat{x}}\sigma_{e}} - \frac{\operatorname{Var}(\hat{x})}{\sigma_{\hat{x}}\sigma_{e}}$$
$$= \frac{\operatorname{Cov}(\hat{x}, x)}{\sigma_{\hat{x}}\sigma_{e}} \frac{\sigma_{x}}{\sigma_{x}} - \frac{\sigma_{\hat{x}}}{\sigma_{e}}$$
$$= \frac{\operatorname{Corr}(\hat{x}, x)\sigma_{x}}{\sigma_{e}} - \frac{\sigma_{\hat{x}}}{\sigma_{e}}$$
$$= \sigma_{e}^{-1} \Big(\operatorname{Corr}(\hat{x}, x)\sigma_{x} - \sigma_{\hat{x}} \Big)$$
(A1)

A.3 Appendix: MLP/Neural cross-validation procedure

Unlike many of the other methods explored in this paper, MLP (Multi-Layer Perceptrons) are difficult to cross-validate in a consistent way. This is for three main reasons.

The first reason is due to one of neural methods' advantages for prediction problems—that they are highly adaptable to many different problem spaces, varying both in more traditional hyper-parameters such as learning-rate and neural network width, but also in much-more nuanced choices such as optimization method or input structure. "Neural Networks," despite the term's usage in many settings, is actually less of a single model and more a label placed on an entire class of iteratively-optimized models. The work's aim has been to use "off-the-shelf" machine learning methods to understand what empirical concerns exist when placing these models naively in an otherwise-recognizable econometric instrumental variables setting, but for a neural network, the off-the-shelf model is highly dependent on the problem at hand. Unfortunately, this advantage of MLPs and other Neural Networks makes a full grid-search of the hyper-parameter space intractable. This requires us to restrict the grid-space somewhat to create a tractable solution, while allowing the model a good shot at choosing the "correct" specification. This restriction potentially handicaps the neural model's flexibility, and may produce higher average out-of-sample loss than the full set of model specifications could potentially produce.

Second, neural networks are computationally expensive to train—each combination of hyperparameters must be trained separately over many iterations, and for most optimization procedures it is useful to utilize different re-orderings of the data-set to reach a satisfactory loss-minimizing point. Even for less-complex data such as ours with relatively few observations, cross-validating even 100 hyperparameter combinations over 1000 synthetic datasets leads to prohibitively long training periods given our computational resources.

Last, neural networks are sensitive to their initialization point, which is not the case for any of our other methods included in this analysis. Unlike most other methods, neural-class models can find values for one of a number of potential loss-minimizing local optima for its parameters, and while those local optima can perform similarly well, they do not necessarily produce identical predictions and can fail on different subsets of data in different ways.

These differences make analyzing how cross-validating a neural network dictates hyper-parameters given a reasonable loss function more interesting because the choices a five-fold cross-validation approach might make are indicative of how the most flexible model chooses a specification given different search spaces. For our cross-validation, we use a five-fold cross-validation procedure to match our other methods, and use mean-squared error as our loss function. We fix a few hyperparameters in place—we use no regularization on

Figure A3: Explaining unrestricted/narrow neural networks' bimodal distributions of $\hat{\beta}$

(A) Comparing bias in $\hat{\beta}$: Approximately linear (no hidden (B) Comparing bias in $\hat{\beta}$ and out-of-sample loss: layers) vs. 'deeper' neural networks

Approximately linear vs. 'deeper' models



The y axis of Panel **a** depicts the second-stage estimate $\hat{\beta}$, and the x axis represents the depths of the neural networks cross validated in each of the 1,000 iterations. "Depth 1" implies no hidden layers-directly linking the input and output (approximating linear regression). Horizontal line segments in a connect the two possible depths that the model chose between. The solid dot marks the chosen depth (by cross validation). In a and b we color the two subsets separately to illustrate the source of the bimodal distributions in Figure 2. Panel c gives the probability that a model chose the more shallow (less deep) model for each combination of choices (probabilities are equal to the fractions in the cells). For instance, the top-left corner tells us that when facing a choice between a depth-1 model and a depth-6 model, CV chose the depth-1 model 81 out of 81 times. Each of these figures uses the results from 1,000 iterations of neural-net based 2SLS high-complexity case 2 (Panel C of Figure 2C; 100 mixed-strength instruments with strength decreasing from z_1) with no restrictions on the hyperparameter space (unrestricted in tables 1–2). The results are very similar for the other high-complexity DGPs.

the weights, and use an "Adam" optimizer with it's out-of-box/off-the-shelf learning rate of .001. We trained all models over 40 epochs, and used a batch-size of 10 observations. One unusual step we take is to introduce a leaky rectified linear activation function (ReLU) activation function to connect hidden layers. This was done to prevent the model from suffering from "dying weights" which is when parameters accidentally force a large number of activations to inappropriately "ignore" activations due to a(x) = 0 for all or many values of x.²⁶

We began by creating three separate hyper-parameter search spaces distinguished by maximum-allowable width and depth. The 'shallow/wide' neural networks are allowed to choose from hidden-layer representations that are 16, 32, 64, 256, and 512 nodes in hidden layer width respectively. This model is then restricted to contain at most a single hidden layer, but is also allowed to choose from a model that maps inputs to outputs directly, using a linear activation function. This functional form is sensible given x is linear in z for our DGPs. Excluding the 0-hidden layer case would prevent the cross-validation procedure from finding the easiest approximation for a linear functional form. The cross-validation procedure allows differences in regulation by choosing between a dropout rate of .1 or .2.

The "narrow/deep" neural network is instead allowed to choose from a representation with two, three, four, or five hidden layers each with a number of nodes (width) equal to 16, 32, or 64. This model, too, is allowed to choose from the simple linear mapping of inputs to output, for the same reasons as above.

The last search space, referred to as the "unrestricted" neural network is allowed to choose from any combination of the hyperparameters offered to the narrow or shallow networks—from zero through five hidden layers and using the full complement of widths.

To get a sense of how these search-spaces choose models on average, these three search spaces were used to cross-validate over our first 25 datasets, from which we generated a full list of 125 folds. These cross-dataset folds were then used to find, for each search space, the average out of sample MSE across all 125 folds.

From the set of available models available to every search space and for each iteration, we chose two models at random weighted by their average out of sample MSE. These probabilities were chosen using the weighted upper-tail normal CDF, normalized such that all weights for a given search space sum to one. Formally, where *i* is a given set of hyperparameters, *j* is a search-space and μ_j is the mean out of sample MSE for a given search space:

$$p_{i,j}^{chosen} = \int_{z_{i,j}^{mse}}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{z}{2}} dz , \qquad (A2)$$

$$z_{i,j}^{mse} := \frac{mse_{i,j} - \mu_j}{se(mse_j)} .$$
(A3)

For each iteration and using the probabilities above, two models are chosen at random for each search space, and then cross-validated again using a 5-fold procedure. The "winning" model is chosen by lowest out-of-sample MSE, and is used to predict \hat{x} for the first stage.

For a visual explanation and overview of the results from the selection method, see Figure A3.

A.3.1 Low-bias methods

Because machine learning algorithms are designed to minimize loss (maximizing fit), the fact that $Cov(\hat{x}, e) \neq 0$ is partially by design. To see this fact, consider any prediction method that minimizes mean-squared error (MSE)—conditional on the training data $\{x, z\}$:

²⁶ ReLU was tried initially, however using ReLU on our data seemingly led to a complete shutdown of the predictive power when we attempted it on our weaker-instrument setting. Even using dropout, the MLP's performance was poor in predicting out of sample. We cannot be certain that the dying weights problem was the cause of this issue, but adding a slight negative slope for activations less than 0 seems to have mitigated the behavior.

$$MSE(\hat{x}, x|x, \mathbf{z}) = \underbrace{\left(x - E[\hat{x} \mid x, \mathbf{z}]\right)^{2}}_{(\text{Bias of } x \text{ for } \hat{x})^{2}} + \underbrace{E\left[(\hat{x} - E[\hat{x} \mid x, \mathbf{z}])^{2} \mid x, \mathbf{z}\right]}_{\text{Cond. Var}(\hat{x})} + \underbrace{E\left[\varepsilon^{2} \mid x, \mathbf{z}\right]}_{\text{Cond. Var}(\varepsilon)}$$
(A4)

where ε is the irreducible error—the unknowable disturbance from the DGP of *x*, *i.e.*, $x = f(\mathbf{z}) + \varepsilon$.²⁷

Equation (A4) highlights that in an MSE-minimization problem, \hat{x} is the only component of MSE that a learning algorithm can change (x, z, and ε are all data dependent). This fact leads to the widely discussed variance-bias tradeoff—an arbitrary estimator will generally face a negotiate between low-variance predictions and low-bias predictions. Many traditional econometric estimators result from prioritizing zero bias and then selecting the minimum-variance estimator from this class of unbiased estimators. As (A4) points out, these estimators could reduce their out-of-sample MSE by *accepting* some bias and reducing variance. This tradeoff is at the heart of the prediction improvement many out-of-the-box algorithms offer relative to plain OLS.

A.4 Appendix: Neural approaches to measuring causal response

Neural networks and their offspring offer just such a route to explore data that falls outside of the traditional bounds, whether that is to use transformers for text data, or convolutional neural networks (CNNs) for imaging. However, many of the same problems apply to these tools as were outlined in the paper. In order to make full use of them in a two stage approach, the same stringent restrictions are required to generate meaningful and unbiased coefficients in a two-stage framework. Indeed, when running a cross-validated feed-forward network to produce a meaningful first stage estimation with our high-complexity data most simulations with any hidden layers simply reproduced an approximation of β close to that of naive OLS.²⁸ There is a burgeoning field of research in machine learning that strives to understand IV problems under less-parametric (though generally still somewhat parametric) causal structures and these methods are seeing success in both simulated and real-world data. The downside to using these powerful methods is that they require a new framework in which to understand them, and make interpretation of treatment effects more challenging.

The first of the recent batch of machine learning instrumental variables papers is referred to as "Deep IV" Hartford et al. (2017). The authors throw away the linear functional form for x = f(z, u) in the "first stage," but assume linearly additive confounding variables and learn the causal structure with a two-part one-pass neural network model. The authors do this by recasting the econometric approach to instrumental variables into two interlinked problem spaces - estimating the conditional distribution g(x|z) and then using the approximation of x given z to predict y. This creates difficulties because such methods produce good counterfactual predictions, but have a harder time matching the clean interpretable causal effect of X on y when compared to traditional econometric approaches. Further, because of the flexibility in functional form, models of this category tend to have more trouble outside of $supp(z_{train})$ or $supp(g(x|z_{train}))$ that are observed in a training sample - and it's difficult to apply a post-analysis structure to such a model to gather understanding on counterfactuals where z_{test} is considerably different than z_{train} . Further, many other methods have been created and can be used to estimate instrument-identified causal effects using a similar semi-parametric two-stage function that can identify complex functional forms in either first or second stages Bennett et al. (2019) and Xu et al. (2020) and improve on edge-of-support marginal effects.²⁹ Both of these papers and Deep IV are able to produce causal inferences using images as instruments—something that a regression would not be able to meaningfully do without some form of pre-model dimensionality reduction.

Neural approaches to causal inference are also not limited to use semi-parametric structural forms for het-

²⁷ These expectations are conditional on the given dataset; ε and \hat{x} are conditionally independent by definition. The expectation term is conditional on data observed, so for simplicity, the term $E_D(\hat{x}(z; D))$ will simply be referred to as \hat{x} .

²⁸ See Appendix Section A.3 for full details of the MLP methods.

²⁹ The methodology contained in Xu et al. (2020) is particularly useful, because it does not predict X directly, but rather applies Neural Networks to the task of learning polynomial forms to pass through first and second stages in a 2SLS (with an L2 penalty) and may mostly avoid components a and b as described earlier.

erogenous treatment effects. Kilbertus et al. (2020) created a neural network to identify the total set of conditional causal effects given a fully non-parametric instrumental variable analysis. (With the very reasonable assumption placed on the function of unobserved noise that it does not feature infinite discontinuities, for example.)

In spite of the massive technical improvements these models have made, trying to extract a beta-equivalent from the existing models is difficult, though interpretable machine learning methods do exist. Unfortunately, extracting meaningful information using prediction-explanation methods such as Ribeiro et al. (2016) about how a result is produced, or to infer what kind of economic information can be gathered from the weights within a model, neural networks are hard to interpret Wang et al. (2020). This makes comparing such models to traditional 2SLS or econometric approaches for ATE approximation difficult—and produces complications in choosing benchmarks as to how to evaluate them.

A.5 Appendix: Complications of ML and monotonicity

One often overlooked assumption in instrumental variables irrelevant under the assumption of constant treatment effects, but, without guaranteed constant treatment effects, is referred to as 'monotonicity' of heterogenous treatment effects. If a researcher is willing to simply find the 'Local Average Treatment Effects' or LATE Angrist and Pischke (2009), 2SLS can under certain circumstances recover that estimate. Simply put, this means that while treatment effects can vary across our population, response of endogenous variable x_i to instrument z_i must move in the same direction for all individuals *i*. For canonical cases of binary instruments and treatments, this boils down to an assumption of 'no defiers'.

Our results as written do not conclude one way or another about the implications of the monotonicity assumption with regards to the machine learning in the first stage, as our primary datasets feature homogenous treatment effects are constant and equal to β . However, relaxing from constant treatment effects to continuous, monotonicity-preserving and heterogenous treatment effects, under some circumstances using nonlinear methods can lead to inefficiency in estimates of the LATE.

To illustrate a very simple example, imagine a case estimating the coefficient β_1 where $x_i = 1 + \sum_{\nu=1}^{7} z_i * \gamma_i + \varepsilon_i + u_i$, where $\gamma_i \sim \gamma(.5, 4)$, thus $E(\gamma_i) = 2$, $min(\gamma_i) = 0$, $Y = 1 + X_i\beta_1 + u_i + \epsilon_i$ and all error terms $\epsilon_i, \varepsilon_i, u_i N(0, 1)$. This extends the strengthened exclusion restriction from $E(z|\varepsilon) = 0$ to $E(z|\varepsilon) = E(z|\eta X) = 0$ where $\eta = \gamma_i - E(\gamma_i)$.³⁰ In this case, the instruments are not themselves weak, but have weak effects on varying members of the sample. If a machine learning algorithm conditions its predictions on the instruments directly, and the coefficients γ_i are sufficiently varied, x_{ssnl} will only be a better estimate in expectation.

In this case; under the stronger exclusion restriction described above and using a MLSLS strategy as described in Chen et al. (2020) with a non-linear meta-model, will produce an unbiased result, but also may result in inflated standard errors relative to a 2SLS or SSIV approach. This is because, under this kind of monotonicity, linear IV becomes a weighted average of marginal treatment effects.³¹ If the econometrician believes this type of variation exists, the assumptions required for 2SLS' validity are strong, and MLSLS' are stronger. ML makes no guarantees on recovering this weighted sum in the same manner, so there exists important future work to examine how exactly this may impact structural estimates of β .

³⁰ See Heckman et al. (2006) for the original treatment of η in the IV case.

³¹ which the reader should turn to Heckman and Vytlacil (2005) and Heckman et al. (2006) to see a formal treatment of.