

Inference on treatment effects after selection amongst high- dimensional controls

Alexandre Belloni
Victor Chernozhukov
Christian Hansen

The Institute for Fiscal Studies
Department of Economics, UCL

cemmap working paper CWP26/13

INFERENCE ON TREATMENT EFFECTS AFTER SELECTION AMONGST HIGH-DIMENSIONAL CONTROLS

A. BELLONI, V. CHERNOZHUKOV, AND C. HANSEN

ABSTRACT. We propose robust methods for inference on the effect of a treatment variable on a scalar outcome in the presence of very many controls. Our setting is a partially linear model with possibly non-Gaussian and heteroscedastic disturbances where the number of controls may be much larger than the sample size. To make informative inference feasible, we require the model to be approximately sparse; that is, we require that the effect of confounding factors can be controlled for up to a small approximation error by conditioning on a relatively small number of controls whose identities are unknown. The latter condition makes it possible to estimate the treatment effect by selecting approximately the right set of controls. We develop a novel estimation and uniformly valid inference method for the treatment effect in this setting, called the “post-double-selection” method. Our results apply to Lasso-type methods used for covariate selection as well as to any other model selection method that is able to find a sparse model with good approximation properties.

The main attractive feature of our method is that it allows for imperfect selection of the controls and provides confidence intervals that are valid uniformly across a large class of models. In contrast, standard post-model selection estimators fail to provide uniform inference even in simple cases with a small, fixed number of controls. Thus our method resolves the problem of uniform inference after model selection for a large, interesting class of models. We also present a simple generalization of our method to a fully heterogeneous model with a binary treatment variable. We illustrate the use of the developed methods with numerical simulations and an application that considers the effect of abortion on crime rates.

Key Words: treatment effects, partially linear model, high-dimensional-sparse regression, inference under imperfect model selection, uniformly valid inference after model selection, average treatment effects, average treatment effects for the treated

Date: First version: May 2010. This version is of July 19, 2013. This is a revision of a 2011 ArXiv/CEMMAP paper entitled “Estimation of Treatment Effects with High-Dimensional Controls”.

We thank Ted Anderson, Takeshi Amemiya, Stéphane Bonhomme, Mathias Cattaneo, Gary Chamberlain, Denis Chetverikov, Graham Elliott, Eric Tchetgen, Bruce Hansen, James Hamilton, Jin Hahn, Han Hong, Guido Imbens, Zhipeng Liao, Tom MaCurdy, Anna Mikusheva, Whitney Newey, Alexei Onatsky, Joseph Romano, Shinich Sakata, Andres Santos, Mathew Shum, Chris Sims, and participants of 10th Econometric World Congress in Shanghai 2010, Caltech, CIREQ-Montreal, Harvard-MIT, UCL, UCLA UC San-Diego, USC, UC-Davis, Princeton, Stanford, Informetrics Workshop, and Yale NSF conference on “High-Dimensional Statistics” for extremely helpful comments.

1. INTRODUCTION

Many empirical analyses focus on estimating the structural, causal, or treatment effect of some variable on an outcome of interest. For example, we might be interested in estimating the causal effect of some government policy on an economic outcome such as employment. Since economic policies and many other economic variables are not randomly assigned, economists rely on a variety of quasi-experimental approaches based on observational data when trying to estimate such effects. One important method is based on the assumption that the variable of interest can be taken as randomly assigned after controlling for a sufficient set of other factors. Economists, for example, might argue that changes in state-level public policies can be taken as randomly assigned relative to unobservable factors that could affect changes in state-level outcomes after controlling for aggregate macroeconomic activity, state-level economic activity, and state-level demographics; see, for example, Heckman, LaLonde, and Smith (1999) or Imbens (2004).

A problem empirical researchers face when relying on a conditional-on-observables identification strategy for estimating a structural effect is knowing which controls to include. Typically, economic intuition will suggest a set of variables that might be important but will not identify exactly which variables are important or the functional form with which variables should enter the model. This lack of clear guidance about what variables to use leaves researchers with the problem of selecting a set of controls from a potentially vast set of variables including raw regressors available in the data as well as interactions and other transformations of these regressors. A typical economic study will rely on an *ad hoc* sensitivity analysis in which a researcher reports results for several different sets of controls in an attempt to show that the parameter of interest that summarizes the causal effect of the policy variable is insensitive to changes in the set of control variables. See Donohue III and Levitt (2001), which we use as the basis for the empirical study in this paper, or examples in Angrist and Pischke (2008) among many other references.

We present an approach to estimating and performing inference on structural effects in an environment where the treatment variable may be taken as exogenous conditional on observables that complements existing strategies. We pose the problem in the framework of a partially linear model

$$y_i = d_i \alpha_0 + g(z_i) + \zeta_i \quad (1.1)$$

where d_i is the treatment/policy variable of interest, z_i is a set of control variables, and ζ_i is an unobservable that satisfies $E[\zeta_i \mid d_i, z_i] = 0$.¹ The goal of the econometric analysis is to conduct inference on the treatment effect α_0 . We examine the problem of selecting a set of variables from among p potential controls $x_i = P(z_i)$, which may consist of z_i and transformations of z_i , to adequately approximate $g(z_i)$ allowing for $p > n$. Of course, useful inference about α_0 is unavailable

¹We note that d_i does not need to be binary. This structure may also arise in the context of randomized treatment in the case where treatment assignment depends on underlying control variables, potentially in a complicated way. See, for example, Duflo, Glennerster, and Kremer (2008), especially Section 6.1, and Kremer and Glennerster (2011).

in this framework without imposing further structure on the data. We impose such structure by assuming that exogeneity of d_i may be taken as given once one controls linearly for a relatively small number $s < n$ of variables in x_i whose identities are *a priori* unknown. This assumption implies that linear combinations of these s unknown controls provide approximations to $g(z_i)$ and to $E[d_i|z_i] = m(z_i)$ which produce relatively small approximation errors for each object.² This assumption, which is termed approximate sparsity or simply sparsity, allows us to approach the problem of estimating α_0 as a variable selection problem. This framework allows for the realistic scenario in which the researcher is unsure about exactly which variables or transformations are important confounds and so must search among a broad set of controls.

The assumed sparsity includes as special cases the most common approaches to parametric and nonparametric regression analysis. Sparsity justifies the use of fewer variables than there are observations in the sample. When the initial number of variables is high, the assumption justifies the use of variable selection methods to reduce the number of variables to a manageable size. In many economic applications, formal and informal strategies are used to select such smaller sets of control variables. Many of these standard variable selection strategies are non-robust and may produce poor inference.³ In an effort to demonstrate robustness of their conclusions, researchers often employ *ad hoc* sensitivity analyses which examine the robustness of inferential conclusions to variations in the set of controls. Such sensitivity analyses are useful but lack rigorous justification. As a complement to these *ad hoc* approaches, we propose a formal, rigorous approach to inference allowing for selection of controls. Our proposal uses modern variable selection methods in a novel manner which results in valid inference following model selection.

The main contributions of this paper are providing an estimation and inference method within a partially linear model with potentially very high-dimensional controls and developing the supporting theory establishing its validity across a rich class of data-generating-processes (dgps). The method relies on the use of Lasso-type or other sparsity-inducing procedures for variable selection. Our approach differs from usual post-model-selection methods that rely on a single selection step. Rather, we use two different variable selection steps followed by a final estimation step as follows:

1. In the first step, we select a set of control variables that are useful for predicting the treatment d_i . This step helps to insure validity of post-model-selection-inference by finding control variables that are strongly related to the treatment and thus potentially important confounding factors.

²We carefully define what we mean by small approximation errors in Section 2.

³An example of these poor inferential properties is given in Figure 1 (left panel), presented in the next section, where a standard post-model selection estimator has a bimodal distribution which sharply deviates from the standard normal distribution. More examples are given in Section 6 where we document the poor inferential performance of a standard post-model selection method.

2. In the second step, we select additional variables by selecting control variables that predict y_i . This step helps to insure that we have captured important elements in the equation of interest, ideally helping keep the residual variance small, as well as providing an additional chance to find important confounds.
3. In the final step, we estimate the treatment effect α_0 of interest by the linear regression of y_i on the treatment d_i and the union of the set of variables selected in the two variable selection steps.

We provide theoretical results on the properties of the resulting treatment effect estimator and show that it provides inference that is uniformly valid over large classes of models and also achieves the semi-parametric efficiency bound under some conditions. Importantly, our theoretical results allow for imperfect variable selection in either of the two variable selection steps as well as allowing for non-Gaussianity and heteroscedasticity of the model's errors.⁴

We illustrate the theoretical results through an examination of the effect of abortion on crime rates following Donohue III and Levitt (2001). In this example, we find that the formal variable selection procedure produces a qualitatively different result than that obtained through the *ad hoc* set of sensitivity results presented in the original paper. By using formal variable selection, we select a small set of between six and nine variables depending on the outcome, compared to the set of eight variables considered by Donohue III and Levitt (2001). Once this set of variables is linearly controlled for, the estimated abortion effect is rendered imprecise. It is interesting that the key variable selected by the variable selection procedure is the initial condition for the abortion rate. The selection of this initial condition and the resulting imprecision of the estimated treatment effect suggest that one cannot determine precisely whether the effect attributed to abortion found when this initial condition is omitted from the model is due to changes in the abortion rate or some other persistent state-level factor that is related to relevant changes in the abortion rate and current changes in the crime rate.⁵ It is interesting that Foote and Goetz (2008) raise a similar concern based on intuitive grounds and additional data in a comment on Donohue III and Levitt (2001). Foote and Goetz (2008) find that a linear trend interacted with crime rates before abortion could have had an effect renders the estimated abortion effects imprecise.⁶ Overall, finding that a formal, rigorous approach to variable selection produces a qualitatively different result than a more

⁴In a companion paper that presents an overview of results for ℓ_1 - and post ℓ_1 -penalized estimators and inference methods, Belloni, Chernozhukov, and Hansen (2011), we provide similar results in the idealized Gaussian homoscedastic framework.

⁵Note that all models are estimated in first-differences to eliminate any state-specific factors that might be related to both the relevant level of the abortion rate and the level of the crime rate.

⁶Donohue III and Levitt (2008) provide yet more data and a more flexible specification in response to Foote and Goetz (2008). In a supplement available at <http://faculty.chicagobooth.edu/christian.hansen/research/>, we provide additional results based on Donohue III and Levitt (2008). The conclusions are similar to those obtained in this paper in that we find the estimated abortion effect becomes imprecise once one allows for a broad set of controls and selects among them. However, the specification of Donohue III and Levitt (2008) relies on a large number of district

ad hoc approach suggests that these methods might be used to complement economic intuition in selecting control variables for estimating treatment effects in settings where treatment is taken as exogenous conditional on observables.

Relationship to literature. We contribute to several existing literatures. First, we contribute to the literature on semi-parametric estimation of partially linear models; see Donald and Newey (1994), Härdle, Liang, and Gao (2000), Robinson (1988), and others.⁷ We differ from most of the existing literature which considers $p \ll n$ series terms by allowing $p \gg n$ series terms from which we select $\hat{s} \ll n$ terms to construct the regression fits. Considering an initial broad set of terms allows for more refined approximations of regression functions relative to the usual approach that uses only a few low-order terms. See, for example, Belloni, Chernozhukov, and Hansen (2011) for a wage function example and Section 5 for theoretical examples. However, our most important contribution is to allow for data-dependent selection of the appropriate series terms. The previous literature on inference in the partially linear model generally takes the series terms as given without allowing for their data-driven selection. However, selection of series terms is crucial for achieving consistency when $p \gg n$ and is needed for increasing efficiency even when $p = Cn$ with $C < 1$. That the standard estimator can be highly inefficient in the latter case follows from results in Cattaneo, Jansson, and Newey (2010).⁸ We focus on Lasso for performing this selection as a theoretically and computationally attractive device but note that any other method, such as selection using the traditional generalized cross-validation criteria, will work as long as the method guarantees sufficient sparsity in its solution. After model selection, one may apply conventional standard errors or the refined standard errors proposed by Cattaneo, Jansson, and Newey (2010).⁹

Second, we contribute to the literature on the estimation of treatment effects. We note that the policy variable d_i does not have to be binary in our framework. However, our method has a useful interpretation related to the propensity score when d_i is binary. In the first selection step, we select terms from x_i that predict the treatment d_i , i.e. terms that explain the propensity score. We also select terms from x_i that predict y_i , i.e. terms that explain the outcome regression function. Then we run a final regression of y_i on the treatment d_i and the union of selected terms. Thus, our procedure relies on the selection of variables relevant for both the propensity score and

cross time fixed effects and so does not immediately fit into our regularity conditions. We conjecture the methodology continues to work in this case but leave verification to future research.

⁷Estimation of the parameters of the linear part of a partially linear model is typically done by regressing $y_i - \hat{E}[y_i|z_i]$ on $d_i - \hat{E}[d_i|z_i]$ where $\hat{E}[y_i|z_i]$ and $\hat{E}[d_i|z_i]$ are preliminary nonparametric estimators of the conditional expectations of y_i and d_i given z_i under the assumption that $\dim(z_i)$ is small. Our approach fits neatly within this framework where we are offering selection based estimators of the conditional expectation functions.

⁸Cattaneo, Jansson, and Newey (2010) derive properties of series estimator under $p = Cn$, $C < 1$, asymptotics. It follows from their results that under homoscedasticity the series estimator achieves the semiparametric efficiency bound only if $C \rightarrow 0$.

⁹If the selected number of terms \hat{s} is a substantial fraction of n , we recommend using Cattaneo, Jansson, and Newey (2010) standard errors after applying our model selection procedure.

the outcome regression. The structure of our approach also implies that it offers another useful approach to estimating and then conditioning on a propensity score while performing regression adjustment for additional controls in the equation for the outcome given the propensity score. Relying on selecting variables that are important for both objects allows us to achieve two goals: we obtain uniformly valid confidence sets for α_0 despite imperfect model selection, and we achieve full efficiency for estimating α_0 in the homoscedastic case. The relation of our approach to the propensity score brings about interesting connections to the treatment effects literature. Hahn (1998), Heckman, Ichimura, and Todd (1998), and Abadie and Imbens (2011) have constructed efficient regression or matching-based estimates of average treatment effects. Hahn (1998) also shows that conditioning on the propensity score is unnecessary for efficient estimation of average treatment effects. Hirano, Imbens, and Ridder (2003) demonstrate that one can efficiently estimate average treatment effects using estimated propensity score weighting alone. Robins and Rotnitzky (1995) have shown that using propensity score modeling coupled with a parametric regression model leads to efficient estimates if either the propensity score model or the parametric regression model is correct. While our contribution is quite distinct from these approaches, it also highlights the important robustness role played by the propensity score model in the selection of the right control terms for the final regression.

Third, we contribute to the literature on estimation and inference with high-dimensional data and to the uniformity literature. There has been extensive work on estimation and perfect model selection in both low and high-dimensional contexts,¹⁰ but there has been little work on inference after imperfect model selection. Perfect model selection relies on extremely unrealistic assumptions, and even moderate model selection mistakes can have serious consequences for inference as has been shown in Pötscher (2009b), Leeb and Pötscher (2008), and others. In work on instrument selection for estimation of a linear instrumental variables model, Belloni, Chen, Chernozhukov, and Hansen (2012) have shown that moderate model selection mistakes do not prevent valid inference about low-dimensional structural parameters due to the inherent adaptivity of the problem: Omission of a relevant instrument does not affect consistency of an IV estimator as long as there is another relevant instrument. The partially linear regression model (1.1) does not have the same adaptivity structure, and model selection based on the outcome regression alone produces confidence intervals with poor coverage properties. However, our post-double selection procedure creates the necessary adaptivity by performing two separate model selection steps. Performing the two selection steps helps reduce omitted variable bias sufficiently that it is possible to perform uniform inference after model selection. The uniformity holds over large, interesting classes of high-dimensional sparse models.¹¹ In that regard, our contribution is in the spirit of and builds upon the classical

¹⁰For reviews focused on econometric applications, see, e.g., Hansen (2005) and Belloni, Chernozhukov, and Hansen (2010).

¹¹Note that this claim only applies to the main parameter α_0 and does not apply to the nuisance part g . Furthermore, our claim of uniformity only applies to models in which both g and m are approximately sparse, i.e.

contribution by Romano (2004) on the uniform validity of t-tests for the univariate mean. It also shares the spirit of recent contributions, among others, by Mikusheva (2007) on uniform inference in autoregressive models, by Andrews and Cheng (2011) on uniform inference in moment condition models that are potentially unidentified, and by Andrews, Cheng, and Guggenberger (2011) on a generic framework for uniformity analysis.

Finally, we contribute to the broader literature on high-dimensional estimation. For variable selection we use ℓ_1 -penalized methods, though our method and theory will allow for the use of other methods. ℓ_1 -penalized methods have been proposed for model selection problems in high-dimensional least squares problems, e.g. Lasso in Frank and Friedman (1993) and Tibshirani (1996), in part because they are computationally efficient. Many ℓ_1 -penalized methods and related methods have been shown to have good estimation properties even when perfect variable selection is not feasible; see, e.g., Candès and Tao (2007), Meinshausen and Yu (2009), Bickel, Ritov, and Tsybakov (2009), Huang, Horowitz, and Wei (2010), Belloni and Chernozhukov (2013) and the references therein. Such methods have also been shown to extend to nonparametric and non-Gaussian cases as in Bickel, Ritov, and Tsybakov (2009) and Belloni, Chen, Chernozhukov, and Hansen (2012). These methods produce models with a relatively small set of variables. The last property is important in that it leaves the researcher with a set of variables that may be examined further; in addition it corresponds to the usual approach in economics that relies on considering a small number of controls.

Paper Organization. In Section 2, we formally present the modeling environment including the key sparsity condition and develop our estimation and inference method. We establish the consistency and asymptotic normality of our estimator of α_0 uniformly over large classes of models in Section 3; then we present a generalization of the basic procedure to allow for model selection methods other than Lasso. In Section 4, we present a series of theoretical examples in which we provide primitive condition that imply the higher-level conditions of Section 3; and the present a series of numerical examples that verify our theoretical results numerically. In Section 5 we present other potential extensions and present a generalization to heterogeneous treatment effects when d_i is binary. We apply our method to the abortion and crime example of Donohue III and Levitt (2001) in Section 6. In appendices, we provide the proofs.

well-approximated by $s \ll n^{1/2}$ terms, so that they are both estimable at $o(n^{-1/4})$ rates. Approximate sparsity is more general than assumptions often used to justify series estimation; so in that regard, the uniformity regions - the sets of models over which inference is valid - are substantial. We demonstrate that these uniformity regions for our proposed approach are larger than the uniformity regions for standard post-selection methods theoretically and through Monte-Carlo experiments. We also note that any useful statistical inference method is uniformly valid over some regions and stops being uniformly valid once the regions are allowed to be large enough. Our aim is to find methods with uniformity regions that encompass interesting models, not to provide methods that would provide uniformly valid inference over all models. For example, our approach will not work in “dense” models, models where g or m are not well-approximated unless $s \gg n^{1/2}$ terms are used.

Notation. In what follows, we work with triangular array data $\{(\omega_{i,n}, i = 1, \dots, n), n = 1, 2, 3, \dots\}$ defined on probability space $(\Omega, \mathcal{A}, P_n)$, where $P = P_n$ can change with n . Each $\omega_{i,n} = (y'_{i,n}, z'_{i,n}, d'_{i,n})'$ is a vector with components defined below, and these vectors are i.n.i.d. – independent across i , but not necessarily identically distributed. Thus, all parameters that characterize the distribution of $\{\omega_{i,n}, i = 1, \dots, n\}$ are implicitly indexed by P_n and thus by n . We omit the dependence on these objects from the notation in what follows for simplicity. We use array asymptotics as doing so allows consideration of approximating sequences that better capture some finite-sample phenomena and to insure the robustness of conclusions with respect to perturbations of the data-generating process P along various sequences. This robustness, in turn, translates into uniform validity of confidence regions over certain regions of data-generating processes.

We use the following empirical process notation, $\mathbb{E}_n[f] := \mathbb{E}_n[f(\omega_i)] := \sum_{i=1}^n f(\omega_i)/n$, and $\mathbb{G}_n(f) := \sum_{i=1}^n (f(\omega_i) - \mathbb{E}[f(\omega_i)])/\sqrt{n}$. Since we want to deal with i.n.i.d. data, we also introduce the average expectation operator: $\bar{\mathbb{E}}[f] := \mathbb{E}\mathbb{E}_n[f] = \mathbb{E}\mathbb{E}_n[f(\omega_i)] = \sum_{i=1}^n \mathbb{E}[f(\omega_i)]/n$. The l_2 -norm is denoted by $\|\cdot\|$, and the l_0 -norm, $\|\cdot\|_0$, denotes the number of non-zero components of a vector. We use $\|\cdot\|_\infty$ to denote the maximal element of a vector. Given a vector $\delta \in \mathbb{R}^p$, and a set of indices $T \subset \{1, \dots, p\}$, we denote by $\delta_T \in \mathbb{R}^p$ the vector in which $\delta_{Tj} = \delta_j$ if $j \in T$ and $\delta_{Tj} = 0$ if $j \notin T$. We use the notation $(a)_+ = \max\{a, 0\}$, $a \vee b = \max\{a, b\}$, and $a \wedge b = \min\{a, b\}$. We also use the notation $a \lesssim b$ to denote $a \leq cb$ for some constant $c > 0$ that does not depend on n ; and $a \lesssim_P b$ to denote $a = O_P(b)$. For an event E , we say that E wp $\rightarrow 1$ when E occurs with probability approaching one as n grows. We also use \rightsquigarrow to denote convergence in distribution. Given a p -vector b , we denote $\text{support}(b) = \{j \in \{1, \dots, p\} : b_j \neq 0\}$.

2. INFERENCE ON TREATMENT AND STRUCTURAL EFFECTS CONDITIONAL ON OBSERVABLES

2.1. Framework. We consider the partially linear model

$$y_i = d_i \alpha_0 + g(z_i) + \zeta_i, \quad \mathbb{E}[\zeta_i | z_i, d_i] = 0, \quad (2.2)$$

$$d_i = m(z_i) + v_i, \quad \mathbb{E}[v_i | z_i] = 0, \quad (2.3)$$

where y_i is the outcome variable, d_i is the policy/treatment variable whose impact α_0 we would like to infer,¹² z_i represents confounding factors on which we need to condition, and ζ_i and v_i are disturbances. The parameter α_0 is the average treatment or structural effect under appropriate conditions given, for example, in Heckman, LaLonde, and Smith (1999) or Imbens (2004) and is the parameter of interest in many empirical studies.

The confounding factors z_i affect the policy variable via the function $m(z_i)$ and the outcome variable via the function $g(z_i)$. Both of these functions are unknown and potentially complicated.

¹²We consider the case where d_i is a scalar for simplicity. Extension to the case where d_i is a vector of fixed, finite dimension is accomplished by introducing an equation like (5.43) for each element of the vector.

We use linear combinations of control terms $x_i = P(z_i)$ to approximate $g(z_i)$ and $m(z_i)$, writing (5.42) and (5.43) as

$$y_i = d_i \alpha_0 + \underbrace{x_i' \beta_{g0} + r_{gi}}_{g(z_i)} + \zeta_i, \quad (2.4)$$

$$d_i = \underbrace{x_i' \beta_{m0} + r_{mi}}_{m(z_i)} + v_i, \quad (2.5)$$

where $x_i' \beta_{g0}$ and $x_i' \beta_{m0}$ are approximations to $g(z_i)$ and $m(z_i)$, and r_{gi} and r_{mi} are the corresponding approximation errors. In order to allow for a flexible specification and incorporation of pertinent confounding factors, the vector of controls, $x_i = P(z_i)$, we can have a dimension $p = p_n$ which can be large relative to the sample size. Specifically, our results only require $\log p = o(n^{1/3})$ along with other technical conditions. High-dimensional regressors $x_i = P(z_i)$ could arise for different reasons. For instance, the list of available controls could be large, i.e. $x_i = z_i$ as in e.g. Koenker (1988). It could also be that many technical controls are present; i.e. the list $x_i = P(z_i)$ could be composed of a large number of transformations of elementary regressors z_i such as B-splines, dummies, polynomials, and various interactions as in Newey (1997), Chen (2007), or Chen and Pouzo (2009; 2012).

Having very many controls creates a challenge for estimation and inference. A key condition that makes it possible to perform constructive estimation and inference in such cases is termed sparsity. Sparsity is the condition that there exist approximations $x_i' \beta_{g0}$ and $x_i' \beta_{m0}$ to $g(z_i)$ and $m(z_i)$ in (5.44)-(5.45) that require only a small number of non-zero coefficients to render the approximation errors r_{gi} and r_{mi} small relative to estimation error. More formally, sparsity relies on two conditions. First, there exist β_{g0} and β_{m0} such that at most $s = s_n \ll n$ elements of β_{m0} and β_{g0} are non-zero so that

$$\|\beta_{m0}\|_0 \leq s \text{ and } \|\beta_{g0}\|_0 \leq s.$$

Second, the sparsity condition requires the size of the resulting approximation errors to be small compared to the conjectured size of the estimation error:

$$\{\bar{E}[r_{gi}^2]\}^{1/2} \lesssim \sqrt{s/n} \text{ and } \{\bar{E}[r_{mi}^2]\}^{1/2} \lesssim \sqrt{s/n}.$$

Note that the size of the approximating model $s = s_n$ can grow with n just as in standard series estimation.

The high-dimensional-sparse-model framework outlined above extends the standard framework in the treatment effect literature which assumes both that the identities of the relevant controls are known and that the number of such controls s is much smaller than the sample size. Instead, we assume that there are many, p , potential controls of which at most s controls suffice to achieve a desirable approximation to the unknown functions $g(\cdot)$ and $m(\cdot)$ and allow the identity of these controls to be unknown. Relying on this assumed sparsity, we use selection methods to select approximately the right set of controls and then estimate the treatment effect α_0 .

2.2. The Method: Least Squares after Double Selection. We propose the following method for estimating and performing inference about α_0 . The most important feature of this method is that it does not rely on the highly unrealistic assumption of perfect model selection which is often invoked to justify inference after model selection. To the best of our knowledge, our result is the first of its kind in this setting. This result extends our previous results on inference under imperfect model selection in the instrumental variables model given in Belloni, Chen, Chernozhukov, and Hansen (2012). The problem is fundamentally more difficult in the present paper since model selection mistakes in which an instrument with a small coefficient is missed have relatively little impact in the IV model if there are instruments with large coefficients while similar model selection mistakes may substantively impact inference on α in the present problem. We overcome this difficulty by introducing additional model selection steps. The construction of our advocated procedure reflects our effort to offer a method that has attractive properties for inference across a wide range of dgps. The estimator is \sqrt{n} -consistent and asymptotically normal under mild conditions and provides confidence intervals that are robust to various perturbations of the dgp that preserve approximate sparsity.

To define the method, we first write the reduced form corresponding to (5.42)-(5.43) as:

$$y_i = x_i' \bar{\beta}_0 + \bar{r}_i + \bar{\zeta}_i, \quad (2.6)$$

$$d_i = x_i' \beta_{m0} + r_{mi} + v_i, \quad (2.7)$$

where $\bar{\beta}_0 := \alpha_0 \beta_{m0} + \beta_{g0}$, $\bar{r}_i := \alpha_0 r_{mi} + r_{gi}$, $\bar{\zeta}_i := \alpha_0 v_i + \zeta_i$.

We have two equations and hence can apply model selection methods to each equation to select control terms. In this paper, we focus on selection using the Lasso method described below but note other model selection procedures could also be used. Given the set of selected controls from (2.6) and (2.7), we can estimate α_0 by a least squares regression of y_i on d_i and the union of the selected controls. Inference on α_0 may then be performed using conventional methods for inference about parameters estimated by least squares. Intuitively, this procedure works well since we are more likely to recover key controls by considering selection of controls from both equations instead of just considering selection of controls from the single equation (5.44) or (2.6). In finite-sample experiments, single-selection methods essentially fail, providing poor inference relative to the double-selection method outlined above. This performance is also supported theoretically by the fact that the double-selection method requires weaker regularity conditions for its validity and for attaining the semi-parametric efficiency bound¹³ than the single selection method.

Now we formally define the post-double-selection estimator: Let $\widehat{I}_1 = \text{support}(\widehat{\beta}_1)$ denote the control terms selected by a feasible Lasso estimator $\widehat{\beta}_1$ computed using data $(\tilde{y}_i, \tilde{x}_i) = (d_i, x_i)$, $i = 1, \dots, n$. Let $\widehat{I}_2 = \text{support}(\widehat{\beta}_2)$ denote the control terms selected by a feasible Lasso estimator $\widehat{\beta}_2$

¹³The semi-parametric efficiency bound of Robinson (1988) is attained in the homoscedastic case whenever such a bound formally applies.

computed using data $(\tilde{y}_i, \tilde{x}_i) = (y_i, x_i)$, $i = 1, \dots, n$. The post-double-selection estimator $\check{\alpha}$ of α_0 is defined as the least squares estimator obtained by regressing y_i on d_i and the selected control terms x_{ij} with $j \in \hat{I} \supseteq \hat{I}_1 \cup \hat{I}_2$:

$$(\check{\alpha}, \check{\beta}) = \underset{\alpha \in \mathbb{R}, \beta \in \mathbb{R}^p}{\operatorname{argmin}} \{ \mathbb{E}_n[(y_i - d_i\alpha - x'_i\beta)^2] : \beta_j = 0, \forall j \notin \hat{I} \}. \quad (2.8)$$

The set \hat{I} may contain variables that were not selected in the variable selection steps with indices in \hat{I}_3 that the analyst thinks are important for ensuring robustness. We call \hat{I}_3 the amelioration set. Thus,

$$\hat{I} = \hat{I}_1 \cup \hat{I}_2 \cup \hat{I}_3; \quad (2.9)$$

let $\hat{s} = \|\hat{I}\|_0$ and $\hat{s}_j = \|\hat{I}_j\|_0$ for $j = 1, 2, 3$.

We define a feasible Lasso estimator below and note that other selection methods could be used as well. When feasible Lasso is used to construct \hat{I}_1 and \hat{I}_2 , we refer to the post-double-selection estimator as the *post-double-Lasso estimator*. When other model selection devices are used to construct \hat{I}_1 and \hat{I}_2 , we shall refer the estimator as the generic post-double-selection estimator.

The main theoretical result of the paper shows that the post-double-selection estimator $\check{\alpha}$ obeys

$$([\bar{E}v_i^2]^{-1}\bar{E}[v_i^2\zeta_i^2][\bar{E}v_i^2]^{-1})^{-1/2}\sqrt{n}(\check{\alpha} - \alpha_0) \rightsquigarrow N(0, 1) \quad (2.10)$$

under approximate sparsity conditions, *uniformly* within a rich set of data generating processes. We also show that the standard plug-in estimator for standard errors is consistent in these settings. All of these results imply uniform validity of confidence regions over large, interesting classes of models. Figure 2.2 (right panel) illustrates the result (2.10) by showing that the finite-sample distribution of our post-double-selection estimator is very close to the normal distribution. In contrast, Figure 2.2 (left panel) illustrates the problem with the traditional post-single-selection estimator based on (5.44), showing that its distribution is bimodal and sharply deviates from the normal distribution. Finally, it is worth noting that the estimator achieves the semi-parametric efficiency bound under homoscedasticity.

2.3. Selection of controls via feasible Lasso Methods. Here we describe feasible variable selection via Lasso. Note that each of the regression equations above is of the form

$$\tilde{y}_i = \underbrace{\tilde{x}_i'\beta_0 + r_i}_{f(\tilde{z}_i)} + \epsilon_i,$$

where $f(\tilde{z}_i)$ is the regression function, $\tilde{x}_i'\beta_0$ is the approximation based on the dictionary $\tilde{x}_i = P(\tilde{z}_i)$, r_i is the approximation error, and ϵ_i is the error. The Lasso estimator is defined as a solution to

$$\min_{\beta \in \mathbb{R}^p} \mathbb{E}_n[(\tilde{y}_i - \tilde{x}_i'\beta)^2] + \frac{\lambda}{n}\|\beta\|_1, \quad (2.11)$$

where $\|\beta\|_1 = \sum_{j=1}^p |\beta_j|$; see Frank and Friedman (1993) and Tibshirani (1996). The non-differentiability of the penalty function at zero induces the solution $\hat{\beta}$ to have components set exactly to zero, and

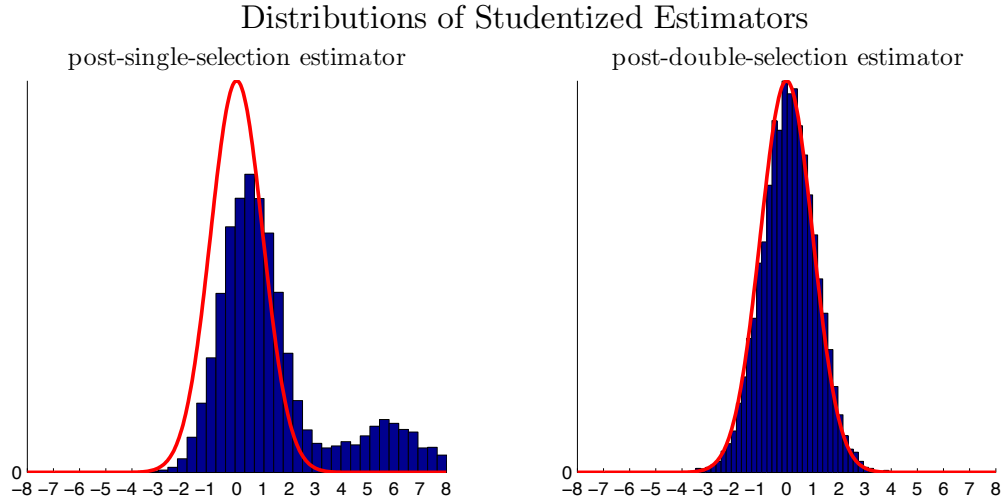


FIGURE 1. The finite-sample distributions (densities) of the standard post-single selection estimator (left panel) and of our proposed post-double selection estimator (right panel). The distributions are given for centered and studentized quantities. The results are based on 10000 replications of Design 1 described in Section 6, with R^2 's in equation (2.6) and (2.7) set to 0.5.

thus the Lasso solution may be used for model selection. The selected model $\hat{T} = \text{support}(\hat{\beta})$ is often used for further refitting by least squares, leading to the so called Post-Lasso or Gauss-Lasso estimator, see, e.g., Belloni and Chernozhukov (2013). The Lasso estimator is computationally attractive because it minimizes a convex function. In the homoskedastic Gaussian case, a basic choice for penalty level suggested by Bickel, Ritov, and Tsybakov (2009) is

$$\lambda = 2 \cdot c\sigma \sqrt{2n \log(2p/\gamma)}, \quad (2.12)$$

where $c > 1$, $1 - \gamma$ is a confidence level that needs to be set close to 1, and σ is the standard deviation of the noise. The formal motivation for this penalty is that it leads to near-optimal rates of convergence of the estimator under approximate sparsity; see Bickel, Ritov, and Tsybakov (2009) or Belloni, Chen, Chernozhukov, and Hansen (2012). The good behavior of the estimator of β_0 in turn implies good approximation properties of the selected model \hat{T} , as noted in Belloni and Chernozhukov (2013). Unfortunately, even in the homoskedastic case the penalty level specified above is not feasible since it depends on the unknown σ .

Belloni, Chen, Chernozhukov, and Hansen (2012) formulate a feasible Lasso estimator $\hat{\beta}$ geared for heteroscedastic, non-Gaussian cases, which solves

$$\min_{\beta \in \mathbb{R}^p} \mathbb{E}_n[(\tilde{y}_i - \tilde{x}_i' \beta)^2] + \frac{\lambda}{n} \|\hat{\Psi} \beta\|_1, \quad (2.13)$$

where $\widehat{\Psi} = \text{diag}(\widehat{l}_1, \dots, \widehat{l}_p)$ is a diagonal matrix of penalty loadings. The penalty level λ and loadings \widehat{l}_j 's are set as

$$\lambda = 2 \cdot c \sqrt{n} \Phi^{-1}(1 - \gamma/2p) \text{ and } \widehat{l}_j = l_j + o_P(1), \quad l_j = \sqrt{\mathbb{E}_n[\tilde{x}_{ij}^2 \epsilon_i^2]}, \text{ uniformly in } j = 1, \dots, p, \quad (2.14)$$

where $c > 1$ and $1 - \gamma$ is a confidence level.¹⁴ The l_j 's are ideal penalty loadings that are not observed, and we estimate l_j by \widehat{l}_j obtained via an iteration method given in Appendix A. The validity of using these estimates was established in Belloni, Chen, Chernozhukov, and Hansen (2012). We refer to the resulting feasible Lasso method as the *Iterated Lasso*. The estimator $\widehat{\beta}$ has statistical performance that is similar to that of the (infeasible) Lasso described above in Gaussian cases and delivers similar performance in non-Gaussian, heteroscedastic cases; see Belloni, Chen, Chernozhukov, and Hansen (2012). In this paper, we only use $\widehat{\beta}$ as a model selection device. Specifically, we only make use of

$$\widehat{T} = \text{support}(\widehat{\beta}),$$

the labels of the regressors with non-zero estimated coefficients. We show that the selected model \widehat{T} has good approximation properties for the regression function f under approximate sparsity in Section 3.

In what follows, we shall use the term *feasible Lasso* to refer to the Iterated Lasso estimator $\widehat{\beta}$ solving (2.13)-(2.14) with $c > 1$ and $1 - \gamma$ set such that

$$\gamma = o(1) \text{ and } \log(1/\gamma) \lesssim \log(p \vee n). \quad (2.15)$$

We note that many other selection devices could be used and provide a brief discussion and some formal results in Section 3.3.

2.4. Intuition for the Importance of Double Selection. To build intuition, we discuss the case where there is only one control; that is, $p = 1$. This scenario provides the simplest possible setting where variable selection might be interesting. In this case, Lasso-type methods act like conservative t -tests which allows the properties of selection methods to be explained easily.

With $p = 1$, the model is

$$y_i = \alpha_0 d_i + \beta_g x_i + \zeta_i, \quad (2.16)$$

$$d_i = \beta_m x_i + v_i. \quad (2.17)$$

For simplicity, all errors and controls are taken as normal,

$$\begin{pmatrix} \zeta_i \\ v_i \end{pmatrix} | x_i \sim N \left(0, \begin{pmatrix} \sigma_\zeta^2 & 0 \\ 0 & \sigma_v^2 \end{pmatrix} \right), \quad x_i \sim N(0, 1), \quad (2.18)$$

¹⁴Practical recommendations include the choice $c = 1.1$ and $\gamma = .05$.

where the variance of x_i is normalized to be 1. The underlying probability space is equipped with probability measure \mathbf{P} . Let \mathbf{P} denote the collection of all dgps \mathbf{P} where (2.16)-(2.18) hold with non-singular covariance matrices in (2.18). Suppose that we have an i.i.d. sample $(y_i, d_i, x_i)_{i=1}^n$ from the dgp $\mathbf{P}_n \in \mathbf{P}$. The subscript n signifies that the dgp and all true parameter values may change with n to better model finite-sample phenomena such as coefficients being “close to zero”. As in the rest of the paper, we keep the dependence of the true parameter values on n implicit. Under the stated assumption, x_i and d_i are jointly normal with variances $\sigma_x^2 = 1$ and $\sigma_d^2 = \beta_m^2 \sigma_x^2 + \sigma_v^2$ and correlation $\rho = \beta_m \sigma_x / \sigma_d$.

The standard post-single-selection method for inference proceeds by applying model selection methods – ranging from standard t -tests to Lasso-type selectors – to the first equation only, followed by applying OLS to the selected model. In the model selection stage, standard selection methods would necessarily omit x_i wp $\rightarrow 1$ if

$$|\beta_g| \leq \frac{\ell_n}{\sqrt{n}} c_n, \quad c_n := \frac{\sigma_\zeta}{\sigma_x \sqrt{1 - \rho^2}}, \quad \text{for some } \ell_n \rightarrow \infty, \quad (2.19)$$

where ℓ_n is a slowly varying sequence depending only on \mathbf{P} . On the other hand, these methods would necessarily include x_i wp $\rightarrow 1$, if

$$|\beta_g| \geq \frac{\ell'_n}{\sqrt{n}} c_n, \quad \text{for some } \ell'_n > \ell_n, \quad (2.20)$$

where ℓ'_n is another slowly varying sequence in n depending only on \mathbf{P} . With most standard model selection devices with sensible tuning choices, we have $\ell_n = C\sqrt{\log n}$ and $\ell'_n = C'\sqrt{\log n}$ with constants C and C' depending only on \mathbf{P} . In the case of Lasso methods, we prove this in Section 5. This is also true in the case of the conservative t -test, which omits x_i if the t -statistic $|t| = |\hat{\beta}_g|/\text{s.e.}(\hat{\beta}_g) \geq \Phi^{-1}(1 - 1/(2n))$, where $\hat{\beta}_g$ is the OLS estimator, and $\text{s.e.}(\hat{\beta}_g)$ is the corresponding standard error. In this case, we have $\Phi^{-1}(1 - 1/(2n)) = \sqrt{2\log n}(1 + o(1))$ so the test will have power approaching 1 for alternatives of the form (2.20) with $\ell'_n = 2\sqrt{\log n}$ and power approaching 0 for alternatives of the form (2.20) with $\ell_n = \sqrt{\log n}$.¹⁵

A standard selection procedure would work with the first equation. Under “good” sequences of models \mathbf{P}_n such that (2.20) holds, x_i is included wp $\rightarrow 1$, and the estimator becomes the standard OLS estimator with the standard large sample asymptotics under \mathbf{P}_n

$$\sigma_n^{-1} \sqrt{n}(\hat{\alpha} - \alpha_0) = \underbrace{\sigma_n^{-1} \mathbb{E}_n[v_i^2]^{-1} \sqrt{n} \mathbb{E}_n[v_i \zeta_i]}_{=:i} + o_P(1) \rightsquigarrow N(0, 1),$$

where $\sigma_n^2 = \sigma_\zeta^2(\sigma_v^2)^{-1}$. On the other hand, when $\beta_g = 0$ or $\beta_g = o(\ell_n/\sqrt{n})$ and ρ is bounded away from 1, we have that

$$\sigma_n^{*-1} \sqrt{n}(\hat{\alpha} - \alpha_0) = \underbrace{\sigma_n^{*-1} \mathbb{E}_n[d_i^2]^{-1} \sqrt{n} \mathbb{E}_n[d_i \zeta_i]}_{=:i^*} + o_P(1) \rightsquigarrow N(0, 1),$$

¹⁵This result assumes that the canonical estimator of the standard error is used.

where $\sigma_n^{*2} = \sigma_\zeta^2(\sigma_d^2)^{-1}$. The variance σ_n^{*2} is smaller than the variance σ_n^2 from estimation with x_i included if $\beta_m \neq 0$. The potential reduction in variance is often used as a “motivation” for the standard selection procedure. The estimator is super-efficient, achieving a variance smaller than the semi-parametric efficiency bound under homoscedasticity. That is, the estimator is “too good.”

The “too good” behavior of the procedure that looks solely at the first equation has its price. There are plausible sequences of dgps P_n where $\beta_g = \frac{\ell'_n}{\sqrt{n}}c_n$, the coefficient on x_i is not zero but is close to zero,¹⁶ in which the control x_i is dropped wp $\rightarrow 1$ and

$$|\sigma_n^{*-1}\sqrt{n}(\hat{\alpha} - \alpha_0)| \rightsquigarrow \infty. \quad (2.21)$$

That is, the standard post-selection estimator is not asymptotically normal and even fails to be uniformly consistent at the rate of \sqrt{n} . This poor behavior occurs because the omitted variable bias created by dropping x_i may be large even when the magnitude of the regression coefficient, $|\beta_m|$, in the confounding equation (2.17) is small but is not exactly zero. To see this, note

$$\sigma_n^{*-1}\sqrt{n}(\hat{\alpha} - \alpha_0) = \underbrace{\sigma_n^{*-1}\mathbb{E}_n[d_i^2]^{-1}\sqrt{n}\mathbb{E}_n[d_i\zeta_i]}_{=:i^*} + \underbrace{\sigma_n^{*-1}\mathbb{E}_n[d_i^2]^{-1}\sqrt{n}\mathbb{E}_n[d_ix_i]\beta_g}_{=:ii}.$$

The term i^* has standard behavior; namely $i^* \rightsquigarrow N(0, 1)$. The term ii generates the *omitted variable bias*, and it may be arbitrarily large since, wp $\rightarrow 1$,

$$|ii| \geq \frac{1}{2} \frac{|\rho|}{\sqrt{1-\rho^2}} \ell_n \nearrow \infty,$$

if $\ell_n|\rho| \nearrow \infty$.¹⁷ This yields the conclusion (2.21) by the triangle inequality.

In contrast to the standard approach, the post-double-selection method for inference proceeds by applying model selection methods, such as standard t -tests or Lasso-type selectors, to both equations and taking the selected controls as the union of controls selected from each equation. This selection is then followed by applying OLS to the selected controls. Thus, our approach drops x_i only if the omitted variable bias term ii is small. To see this, note that the double-selection-methods *include* x_i wp $\rightarrow 1$ if its coefficient in either (2.16) or (2.17) is not very small. Mathematically, x_i is included if

$$\text{either } |\beta_g| \geq \frac{\ell'_n}{\sqrt{n}} \left(\frac{\sigma_\zeta}{\sigma_x \sqrt{1-\rho^2}} \right) \text{ or } |\beta_m| \geq \frac{\ell'_n}{\sqrt{n}} \left(\frac{\sigma_v}{\sigma_x} \right) \quad (2.22)$$

where ℓ'_n is a slowly varying sequence in n . As already noted, $\ell_n \propto \ell'_n \propto \sqrt{\log n}$ would be standard for Lasso-type methods as well as for using simple t -tests to do model selection. Considering t -tests and using these rates, we would omit x_i if both $|t_g| = |\hat{\beta}_g|/\text{s.e.}(\hat{\beta}_g) \leq \Phi^{-1}(1 - 1/(2n))$ and $|t_m| = |\hat{\beta}_m|/\text{s.e.}(\hat{\beta}_m) \leq \Phi^{-1}(1 - 1/(2n))$ where $\hat{\beta}_g$ and $\hat{\beta}_m$ denote the OLS estimator from each

¹⁶Such sequences are very relevant in that they are designed to generate approximations that better capture the fact that one cannot distinguish an estimated coefficient from 0 arbitrarily well in any given finite sample.

¹⁷Recall that $\rho = \beta_m \sigma_x / \sigma_d$, so $\ell_n|\rho| \nearrow \infty$ as long as $\ell_n|\beta_m| \nearrow \infty$ assuming that σ_x/σ_d is bounded away from 0 and ∞ .

equation and s.e. denotes the corresponding estimated standard errors. Note that the critical value used in the t-tests above is conservative in the sense that the false rejection probability is tending to zero because $\Phi^{-1}(1 - 1/(2n)) = \sqrt{2 \log n}(1 + o(1))$. Again, note that Lasso-type methods operate similarly.

Given the discussion in the preceding paragraph, it is immediate that the post-double selection estimator satisfies

$$\sigma_n^{-1} \sqrt{n}(\tilde{\alpha} - \alpha_0) = i + o_P(1) \rightsquigarrow N(0, 1) \quad (2.23)$$

under any sequence of $P_n \in \mathbf{P}$. We get this approximating distribution whether or not x_i is omitted. That this is the approximate distribution when x_i is included follows as in the single-selection case. To see that we get the same approximation when x_i is omitted, note that we drop x_i only if

$$\text{both } |\beta_g| \leq \frac{\ell_n}{\sqrt{n}} c_n \text{ and } |\beta'_m| \leq \frac{\ell_n}{\sqrt{n}} (\sigma_v/\sigma_x), \quad (2.24)$$

i.e. coefficients in front of x_i in both equations are small. In this case,

$$\sigma_n^{*-1} \sqrt{n}(\tilde{\alpha} - \alpha_0) = \underbrace{\sigma_n^{*-1} \mathbb{E}_n[d_i^2]^{-1} \sqrt{n} \mathbb{E}_n[d_i \zeta_i]}_{=i^*} + \underbrace{\sigma_n^{*-1} \mathbb{E}_n[d_i^2]^{-1} \sqrt{n} \mathbb{E}_n[d_i x_i] \beta_g}_{=ii}.$$

Once again, the term ii is due to *omitted variable bias*, and it obeys $\text{wp} \rightarrow 1$ under (2.24)

$$|ii| \leq 2\sigma_\zeta^{-1} \sigma_d \sigma_d^{-2} \sqrt{n} \sigma_x^2 |\beta_m \beta_g| \leq 2 \frac{\sigma_v/\sigma_d}{\sqrt{1-\rho^2}} \frac{\ell_n^2}{\sqrt{n}} = 2 \frac{\ell_n^2}{\sqrt{n}} \rightarrow 0,$$

since $(\sigma_v/\sigma_d)^2 = 1 - \rho^2$. Moreover, we can show $i^* - i = o_P(1)$ under such sequences, so the first order asymptotics of $\tilde{\alpha}$ is the same whether x_i is included or excluded.

To summarize, the post-single-selection estimator may not be root- n consistent in sensible models which translates into bad finite-sample properties. The potential poor finite-sample performance may be clearly seen in Monte-Carlo experiments. The estimator $\hat{\alpha}$ is thus non-regular: its first-order asymptotic properties depend on the model sequence P_n in a strong way. In contrast, the post-double selection estimator $\tilde{\alpha}$ guards against omitted variables bias which reduces the dependence of the first-order behavior on P_n . This good behavior under sequences P_n translates into uniform with respect to $P \in \mathbf{P}$ asymptotic normality.

We should note, of course, that the post-double-selection estimator is first-order equivalent to the long-regression in this model.¹⁸ This equivalence disappears under approximating sequences with number of controls proportional to the sample size, $p \propto n$, or greater than the sample size, $p \gg n$. It is these scenarios that motivate the use of selection as a means of regularization. In these more complicated settings the intuition from this simple $p = 1$ example carries through, and

¹⁸This equivalence may be a reason double-selection was previously overlooked. There are higher-order differences between the long regression estimator and our estimator. In particular, our estimator has variance that is smaller than the variance of the long regression estimator by a factor that is proportional to $(1 - \ell_n/\sqrt{n})$ if $\beta_g = 0$; and when $\beta_g = \ell_n/\sqrt{n}$, there is a small reduction in variance that is traded against a small increase in bias.

the post-single selection method has a highly non-regular behavior while the post-double selection method continues to be regular.

It is also informative to consider semi-parametric efficiency in this simple example. The post-single-selection estimator is super-efficient when $\beta_m \neq 0$ and $\beta_g = 0$. The super-efficiency in this case is apparent upon noting that the estimator is root- n consistent and normal with asymptotic variance $E[\zeta_i^2]E[d_i^2]^{-1}$. This asymptotic variance is generally smaller than the semi-parametric efficiency bound $E[\zeta_i^2]E[v_i^2]^{-1}$. The price of this efficiency gain is the fact that the post-single-selection estimator breaks down when β_g may be small but non-zero. The corresponding confidence intervals therefore also break down. In contrast, the post-double-selection estimator remains well-behaved in any case, and confidence intervals based on the double-selection-estimator are uniformly valid for this reason.

3. THEORY OF ESTIMATION AND INFERENCE

3.1. Regularity Conditions. In this section, we provide regularity conditions that are sufficient for validity of the main estimation and inference result. We begin by stating our main condition, which contains the previously defined approximate sparsity assumption as well as other more technical assumptions. Throughout the paper, we let c , C , and q be absolute constants, and let $\ell_n \nearrow \infty$, $\delta_n \searrow 0$, and $\Delta_n \searrow 0$ be sequences of absolute positive constants. By absolute constants, we mean constants that are given, and do not depend on the dgp $P = P_n$.

We assume that for each n the following condition holds on dgp $P = P_n$.

Condition ASTE (P): Approximate Sparse Treatment Effects. (i) $\{(y_i, d_i, z_i), i = 1, \dots, n\}$ are i.n.i.d. vectors on (Ω, \mathcal{F}, P) that obey the model (5.42)-(5.43), and the vector $x_i = P(z_i)$ is a dictionary of transformations of z_i , which may depend on n but not on P . (ii) The true parameter value α_0 , which may depend on P , is bounded, $|\alpha_0| \leq C$. (iii) Functions m and g admit an approximately sparse form. Namely there exists $s \geq 1$ and β_{m0} and β_{g0} , which depend on n and P , such that

$$m(z_i) = x_i' \beta_{m0} + r_{mi}, \quad \|\beta_{m0}\|_0 \leq s, \quad \{\bar{E}[r_{mi}^2]\}^{1/2} \leq C\sqrt{s/n}, \quad (3.25)$$

$$g(z_i) = x_i' \beta_{g0} + r_{gi}, \quad \|\beta_{g0}\|_0 \leq s, \quad \{\bar{E}[r_{gi}^2]\}^{1/2} \leq C\sqrt{s/n}. \quad (3.26)$$

(iv) The sparsity index obeys $s^2 \log^2(p \vee n)/n \leq \delta_n$ and the size of the amelioration set obeys $\hat{s}_3 \leq C(1 \vee \hat{s}_1 \vee \hat{s}_2)$. (v) For $\tilde{v}_i = v_i + r_{mi}$ and $\tilde{\zeta}_i = \zeta_i + r_{gi}$ we have $|\bar{E}[\tilde{v}_i^2 \tilde{\zeta}_i^2] - \bar{E}[v_i^2 \zeta_i^2]| \leq \delta_n$, and $\bar{E}[|\tilde{v}_i|^q + |\tilde{\zeta}_i|^q] \leq C$ for some $q > 4$. Moreover, $\max_{i \leq n} \|x_i\|_\infty^2 s n^{-1/2+2/q} \leq \delta_n$ w.p. $1 - \Delta_n$.

Comment 3.1. The approximate sparsity (iii) and rate condition (iv) are the main conditions for establishing the key inferential result. We present a number of primitive examples to show that these conditions contain standard models used in empirical research as well as more flexible models. Condition (iv) requires that the size \hat{s}_3 of the amelioration set \hat{I}_3 should not be substantially larger

than the size of the set of variables selected by the Lasso method. Simply put, if we decide to include controls in addition to those selected by Lasso, the total number of additions should not dominate the number of controls selected by Lasso. This and other conditions will ensure that the total number \hat{s} of controls obeys $\hat{s} \lesssim_P s$. We also require that $s^2 \log^2(p \vee n)/n \rightarrow 0$. Note that s is the bound on the number of regressors used by a sparse model to achieve an approximation error of order $\sqrt{s/n}$ and that the rate of convergence for the estimated coefficients would be $\sqrt{s/n}$ if we knew the identities of these s variables. Thus, the estimated function converges to the population function at a rate of $\sqrt{s/n}$ in the idealized setting where we know the identities of the relevant variables, and we would achieve an approximation rate of $o(n^{-1/4})$ under the condition that $s^2/n \rightarrow 0$ in this case. When the identities of the relevant variables are unknown, we use the stronger rate condition $s^2 \log^2(p \vee n)/n \rightarrow 0$ where the additional logarithmic term is the cost of not knowing the correct set of variables. This decrease in the rate of convergence can be substantial for large p , for example if $\log p \propto n^\gamma$ for some positive $\gamma < 1/2$. This condition can be relaxed using the sample-splitting method of Fan, Guo, and Hao (2011), which is done in a supplementary appendix. Condition (v) is simply a set of sufficient conditions for consistent estimation of the variance of the double selection estimator. If the regressors are uniformly bounded and the approximation errors are going to zero a.s., it is implied by other conditions stated below; and it can also be demonstrated under other sorts of more primitive conditions. \square

The next condition concerns the behavior of the Gram matrix $\mathbb{E}_n[x_i x_i']$. Whenever $p > n$, the empirical Gram matrix $\mathbb{E}_n[x_i x_i']$ does not have full rank and in principle is not well-behaved. However, we only need good behavior of smaller submatrices. Define the minimal and maximal m -sparse eigenvalue of a semi-definite matrix M as

$$\phi_{\min}(m)[M] := \min_{1 \leq \|\delta\|_0 \leq m} \frac{\delta' M \delta}{\|\delta\|^2} \quad \text{and} \quad \phi_{\max}(m)[M] := \max_{1 \leq \|\delta\|_0 \leq m} \frac{\delta' M \delta}{\|\delta\|^2}. \quad (3.27)$$

To assume that $\phi_{\min}(m)[\mathbb{E}_n[x_i x_i']] > 0$ requires that all empirical Gram submatrices formed by any m components of x_i are positive definite. We shall employ the following condition as a sufficient condition for our results.

Condition SE (P): Sparse Eigenvalues. *There is an absolute sequence $\ell_n \rightarrow \infty$ such that the maximal and minimal $\ell_n s$ -sparse eigenvalues are bounded from below and away from zero, namely with probability at least $1 - \Delta_n$,*

$$\kappa' \leq \phi_{\min}(\ell_n s)[\mathbb{E}_n[x_i x_i']] \leq \phi_{\max}(\ell_n s)[\mathbb{E}_n[x_i x_i']] \leq \kappa'',$$

where $0 < \kappa' < \kappa'' < \infty$ are absolute constants.

Comment 3.2. It is well-known that Condition SE is quite plausible for many designs of interest. For instance, Condition SE holds if

- (a) x_i , $i = 1, \dots, n$, are i.i.d. zero-mean sub-Gaussian random vectors that have population Gram matrix $E[x_i x_i']$ with minimal and maximal $s \log n$ -sparse eigenvalues bounded away from zero and from above by absolute constants where $s(\log n)(\log p)/n \leq \delta_n \rightarrow 0$;
- (b) x_i , $i = 1, \dots, n$, are i.i.d. bounded zero-mean random vectors with $\|x_i\|_\infty \leq K_n$ a.s. that have population Gram matrix $E[x_i x_i']$ with minimal and maximal $s \log n$ -sparse eigenvalues bounded from above and away from zero by absolute constants where $K_n^2 s(\log^3 n)\{\log(p \vee n)\}/n \leq \delta_n \rightarrow 0$.

The claim (a) holds by Theorem 3.2 in Rudelson and Zhou (2011)¹⁹ and claim (b) holds by Lemma 1 in Belloni and Chernozhukov (2013) or by Theorem 1.8 Rudelson and Zhou (2011). Recall that a standard assumption in econometric research is to assume that the population Gram matrix $E[x_i x_i']$ has eigenvalues bounded from above and away from zero, see e.g. Newey (1997). The conditions above allow for this and more general behavior, requiring only that the $s \log n$ sparse eigenvalues of the population Gram matrix $E[x_i x_i']$ are bounded from below and from above. \square

The next condition imposes moment conditions on the structural errors and regressors.

Condition SM (P): Structural Moments. *There are absolute constants $0 < c < C < \infty$ and $4 < q < \infty$ such that for $(\tilde{y}_i, \epsilon_i) = (y_i, \zeta_i)$ and $(\tilde{y}_i, \epsilon_i) = (d_i, v_i)$ the following conditions hold:*

- (i) $\bar{E}[|d_i|^q] \leq C$, $c \leq E[\zeta_i^2 | x_i, v_i] \leq C$ and $c \leq E[v_i^2 | x_i] \leq C$ a.s. $1 \leq i \leq n$,
- (ii) $\bar{E}[|\epsilon_i|^q] + \bar{E}[\tilde{y}_i^2] + \max_{1 \leq j \leq p} \{\bar{E}[x_{ij}^2 \tilde{y}_i^2] + \bar{E}[|x_{ij}^3 \epsilon_i^3|] + 1/\bar{E}[x_{ij}^2]\} \leq C$,
- (iii) $\log^3 p/n \leq \delta_n$,
- (iv) $\max_{1 \leq j \leq p} \{ |(\mathbb{E}_n - \bar{E})[x_{ij}^2 \epsilon_i^2]| + |(\mathbb{E}_n - \bar{E})[x_{ij}^2 \tilde{y}_i^2]| \} + \max_{1 \leq i \leq n} \|x_i\|_\infty^2 \frac{s \log(n \vee p)}{n} \leq \delta_n$ wp $1 - \Delta_n$.

These conditions, which are rather mild, ensure good model selection performance of feasible Lasso applied to equations (2.6) and (2.7). These conditions also allow us to invoke moderate deviation theorems for self-normalized sums from Jing, Shao, and Wang (2003) to bound some important error components.

3.2. The Main Result. The following is the main result of this paper. It shows that the post-double selection estimator is root- n consistent and asymptotically normal. Under homoscedasticity this estimator achieves the semi-parametric efficiency bound. The result also verifies that plug-in estimates of the standard errors are consistent.

Theorem 1 (Estimation and Inference on Treatment Effects). *Let $\{P_n\}$ be a sequence of data-generating processes. Assume conditions ASTE (P), SM (P), and SE (P) hold for $P = P_n$ for each n . Then, the post-double-Lasso estimator $\tilde{\alpha}$, constructed in the previous section, obeys as $n \rightarrow \infty$*

$$\sigma_n^{-1} \sqrt{n}(\tilde{\alpha} - \alpha_0) \rightsquigarrow N(0, 1),$$

¹⁹See also Zhou (2009) and Baraniuk, Davenport, DeVore, and Wakin (2008).

where $\sigma_n^2 = [\bar{E}v_i^2]^{-1}\bar{E}[v_i^2\zeta_i^2][\bar{E}v_i^2]^{-1}$. Moreover, the result continues to apply if σ_n^2 is replaced by $\hat{\sigma}_n^2 = [\mathbb{E}_n\hat{v}_i^2]^{-1}\mathbb{E}_n[\hat{v}_i^2\hat{\zeta}_i^2][\mathbb{E}_n\hat{v}_i^2]^{-1}$, for $\hat{\zeta}_i := [y_i - d_i\hat{\alpha} - x_i'\hat{\beta}]\{n/(n - \hat{s} - 1)\}^{1/2}$ and $\hat{v}_i := d_i - x_i'\hat{\beta}$, $i = 1, \dots, n$ where $\hat{\beta} \in \arg \min_{\beta} \{\mathbb{E}_n[(d_i - x_i'\beta)^2] : \beta_j = 0, \forall j \notin \hat{I}\}$ where \hat{I} is defined as in equation (2.9).

Comment 3.3. Note that under i.i.d. sampling under P and under the conditional homoscedasticity $E[\zeta_i^2|z_i] = E[\zeta_i^2]$, the asymptotic variance σ_n^2 reduces to $\{E[v_i^2]\}^{-1}E[\zeta_i^2]$, which is the semi-parametric efficiency bound for the partially linear model of Robinson (1988). \square

Another and perhaps more noteworthy consequence is the following corollary.

Corollary 1 (Uniformly Valid Confidence Intervals). (i) Let \mathbf{P}_n be the collection of all data-generating processes P for which conditions $ASTE(P)$, $SM(P)$, and $SE(P)$ hold for given n . Let $c(1 - \xi) = \Phi^{-1}(1 - \xi/2)$. Then

$$\lim_{n \rightarrow \infty} \sup_{P \in \mathbf{P}_n} |P(\alpha_0 \in [\hat{\alpha} \pm c(1 - \xi)\hat{\sigma}_n/\sqrt{n}]) - (1 - \xi)| = 0.$$

(ii) Let $\mathbf{P} = \cap_{n \geq 1} \mathbf{P}_n$ be the collection of data-generating processes for which the conditions above hold for all $n \geq 1$. Then

$$\lim_{n \rightarrow \infty} \sup_{P \in \mathbf{P}_n} |P(\alpha_0 \in [\hat{\alpha} \pm c(1 - \xi)\hat{\sigma}_n/\sqrt{n}]) - (1 - \xi)| = 0.$$

By exploiting both equations (5.44) and (5.45) for model selection, the post-double-selection method creates the necessary adaptivity that makes it robust to imperfect model selection. Robustness of the post-double selection method is reflected in the fact that Theorem 1 permits the data-generating process to change with n . Thus, the conclusions of the theorem are valid for a wide variety of sequences of data-generating processes which in turn define the regions \mathbf{P} of uniform validity of the resulting confidence sets. These regions appear to be substantial, as we demonstrate via a sequence of theoretical and numerical examples in Section 5 and 6. In contrast, the standard post-selection method based on (5.44) produces confidence intervals that do not have close to correct coverage in many cases.²⁰

Comment 3.4. Our approach to uniformity analysis is most similar to that of Romano (2004), Theorem 4. It proceeds under triangular array asymptotics, with the sequence of dgps obeying certain constraints; then these results imply uniformity over sets of dgps that obey the constraints for all sample sizes. This approach is also similar to the classical central limit theorems for sample means under triangular arrays, and does not require the dgps to be parametrically (or otherwise tightly) specified, which then translates into uniformity of confidence regions. This approach is somewhat different in spirit to the generic uniformity analysis suggested by Andrews, Cheng, and Guggenberger (2011). \square

²⁰Negative results derived in Pötscher (2009a) for standard post-selection methods do not apply to the proposed estimator because of the use of both equations for model selection.

Comment 3.5. Uniformity holds over a large class of approximately sparse models, which cover conventional models used in series estimation of partially linear models as shown in Section 5. Of course, for every interesting class of models and any non-trivial inference method, one could find an even bigger class of models where the uniformity does not apply. In particular, our models do not cover models with many small coefficients. In the series case, a model with many small coefficients corresponds to a deviation from smoothness towards highly non-smooth functions, for example functions generated as realized paths of a white noise process. The fact that our results do not cover such models motivates further research work on inference procedures that would provide valid inference when one considers deviations from the given class of models that are deemed important. In the simulations in Section 6, we consider incorporating the ridge fit along the other controls to be selected over using lasso to build extra robustness against “many small coefficients” deviations away from approximately sparse models. \square

3.3. Inference after Double Selection by a Generic Selection Method. The conditions provided so far offer a set sufficient conditions that are tied to the use of Lasso as the model selector. The purpose of this section is to prove that the main results apply to any other model selection method that is able to select a sparse model with good approximation properties. As in the case of Lasso, we allow for imperfect model selection. Next we state a high-level condition that summarizes a sufficient condition on the performance of a model selection method that allows the post-double selection estimator to attain good inferential properties.

Condition HLMS (P): High-Dimensional Linear Model Selection. *A model selector provides possibly data-dependent sets $\hat{I}_1 \cup \hat{I}_2 \subseteq \hat{I} \subset \{1, \dots, p\}$ of covariate names such that, with probability $1 - \Delta_n$, $|\hat{I}| \leq Cs$ and*

$$\min_{\beta: \beta_j=0, j \notin \hat{I}_1} \sqrt{\mathbb{E}_n[(m(z_i) - x_i' \beta)^2]} \leq \delta_n n^{-1/4} \text{ and } \min_{\beta: \beta_j=0, j \notin \hat{I}_2} \sqrt{\mathbb{E}_n[(g(z_i) - x_i' \beta)^2]} \leq \delta_n n^{-1/4}.$$

Condition HLMS requires that with high probability the selected models are sparse and generate good approximations for the functions g and m . Examples of methods producing such models include the Dantzig selector (Candès and Tao, 2007), feasible Dantzig selector (Gautier and Tsybakov, 2011), Bridge estimator (Huang, Horowitz, and Ma, 2008), SCAD penalized least squares (Fan and Li, 2001), square-root-Lasso (Belloni, Chernozhukov, and Wang, 2011), and thresholded Lasso (Belloni and Chernozhukov, 2013), to name a few. We emphasize that, similarly to the previous arguments, these conditions allow for imperfect model selection. Nonetheless we note that Condition HLMS implicitly assumes that tuning parameters of the model selection procedure are set properly to achieve these conditions.

The following result establishes the inferential properties of a generic post-double-selection estimator.

Theorem 2 (Estimation and Inference on Treatment Effects under High-Level Model Selection). *Let $\{P_n\}$ be a sequence of data-generating processes and the model selection device be such that conditions $ASTE(P)$, $SM(P)$, $SE(P)$, and $HLMS(P)$ hold for $P = P_n$ for each n . Then the generic post-double-selection estimator $\hat{\alpha}$ based on \hat{I} , as defined in (2.8), obeys*

$$([\bar{E}v_i^2]^{-1}\bar{E}[v_i^2\zeta_i^2][\bar{E}v_i^2]^{-1})^{-1/2}\sqrt{n}(\hat{\alpha} - \alpha_0) \rightsquigarrow N(0, 1).$$

Moreover, the result continues to apply if $\bar{E}[v_i^2]$ and $\bar{E}[v_i^2\zeta_i^2]$ are replaced by $\mathbb{E}_n[\hat{v}_i^2]$ and $\mathbb{E}_n[\hat{v}_i^2\hat{\zeta}_i^2]$ for $\hat{\zeta}_i := [y_i - d_i\hat{\alpha} - x_i'\hat{\beta}]\{n/(n - \hat{s} - 1)\}^{1/2}$ and $\hat{v}_i := d_i - x_i'\hat{\beta}$, $i = 1, \dots, n$ where $\hat{\beta} \in \arg \min_{\beta} \{\mathbb{E}_n[(d_i - x_i'\beta)^2] : \beta_j = 0, \forall j \notin \hat{I}\}$.

Theorem 2 can also be used to establish uniformly valid confidence intervals as shown in the following corollary.

Corollary 2 (Uniformly Valid Confidence Intervals). (i) *Let \mathbf{P}_n be the collection of all data-generating processes P for which conditions $ASTE(P)$, $SM(P)$, $SE(P)$, and $HLMS(P)$ hold for given n . Let $c(1 - \xi) = \Phi^{-1}(1 - \xi/2)$. Then*

$$\lim_{n \rightarrow \infty} \sup_{P \in \mathbf{P}_n} |P(\alpha_0 \in [\hat{\alpha} \pm c(1 - \xi)\hat{\sigma}_n/\sqrt{n}]) - (1 - \xi)| = 0.$$

(ii) *Let $\mathbf{P} = \cap_{n \geq n_0} \mathbf{P}_n$ be the collection of data-generating processes for which the conditions above hold for all $n \geq n_0$ for some n_0 . Then*

$$\lim_{n \rightarrow \infty} \sup_{P \in \mathbf{P}_n} |P(\alpha_0 \in [\hat{\alpha} \pm c(1 - \xi)\hat{\sigma}_n/\sqrt{n}]) - (1 - \xi)| = 0.$$

4. THEORETICAL AND MONTE-CARLO EXAMPLES

4.1. Theoretical Examples. The purpose of this section is to give a sequence of examples – progressing from simple to somewhat involved – that highlight the range of the applicability and robustness of the proposed method. In these examples, we specify primitive conditions which cover a broad range of applications including nonparametric models and high-dimensional parametric models. We emphasize that our main regularity conditions cover even more general models which combine various features of these examples such as models with both nonparametric and high-dimensional parametric components.

In all examples, the model is

$$\begin{aligned} y_i &= d_i\alpha_0 + g(z_i) + \zeta_i, & \mathbb{E}[\zeta_i \mid z_i, v_i] &= 0, \\ d_i &= m(z_i) + v_i, & \mathbb{E}[v_i \mid z_i] &= 0, \end{aligned} \tag{4.28}$$

however, the structure for g and m will vary across examples, and so will the assumptions on the error terms ζ_i and v_i .

4.1.1. *Parametric model with fixed p .* We start out with a simple example, in which the dimension p of the regressors is fixed. In practical terms this example approximates cases with p small compared to n . This simple example is important since standard post-single-selection methods fail even in this simple case. Specifically, they produce confidence intervals that are *not* valid uniformly in the underlying data-generating process; see Leeb and Pötscher (2008). In contrast, the post-double-selection method produces confidence intervals that are valid uniformly in the underlying data-generating process.

Example 1. (Parametric Model with Fixed p .) Consider $(\Omega, \mathcal{A}, \mathbf{P})$ as the probability space, on which we have (y_i, z_i, d_i) as i.i.d. vectors for $i = 1, \dots, n$ obeying the model (4.28) with

$$\begin{aligned} g(z_i) &= \sum_{j=1}^p \beta_{g0j} z_{ij}, \\ m(z_i) &= \sum_{j=1}^p \beta_{m0j} z_{ij}. \end{aligned} \quad (4.29)$$

For estimation we use $x_i = (z_{ij}, j = 1, \dots, p)'$. We assume that there are absolute constants $0 < b < B < \infty$, $q_x \geq q > 4$, with $4/q_x + 4/q < 1$, such that

$$\begin{aligned} \mathbb{E}[\|x_i\|^{q_x}] &\leq B, \quad |\alpha_0| + \|\beta_{g0}\| + \|\beta_{m0}\| \leq B, \quad b \leq \lambda_{\min}(\mathbb{E}[x_i x_i']), \\ b &\leq \mathbb{E}[\zeta_i^2 \mid x_i, v_i], \quad \mathbb{E}[|\zeta_i^q| \mid x_i, v_i] \leq B, \quad b \leq \mathbb{E}[v_i^2 \mid x_i], \quad \mathbb{E}[|v_i^q| \mid x_i] \leq B. \end{aligned} \quad (4.30)$$

Corollary 3 (Parametric Example with Fixed p). *Let \mathbf{P} be the collection of all regression models P that obey the conditions set forth in Example 1 for all n for the given constants (p, b, B, q_x, q) . Then, any $P \in \mathbf{P}$ obeys Conditions ASTE (P) with $s = p$, SE (P), and SM (P) for all $n \geq n_0$, with the constants n_0 and $(\kappa', \kappa'', c, C)$ and sequences Δ_n and δ_n in those conditions depending only on (p, b, B, q_x, q) . Therefore, the conclusions of Theorem 1 hold for any sequence $P_n \in \mathbf{P}$, and the conclusions of Corollary 1 on the uniform validity of confidence intervals apply uniformly in $P \in \mathbf{P}$.*

4.1.2. *Nonparametric Examples.* The next examples are more substantial and include infinite-dimensional models which we approximate with linear functional forms with potentially very many regressors, $p \gg n$. The key to estimation in these models is a smoothness condition which requires regression coefficients to decay at some rates. In series estimation, this condition is often directly connected to smoothness of the regression function.

Let a and A be positive constants. We shall say that a sequence of coefficients

$$\theta = \{\theta_j, j = 1, 2, \dots\}$$

is a -smooth with constant A if

$$|\theta_j| \leq A j^{-a}, \quad j = 1, 2, \dots,$$

which will be denoted as $\theta \in S_A^a$. We shall say that a sequence of coefficients $\theta = \{\theta_j, j = 1, 2, \dots\}$ is a -smooth with constant A after p -rearrangement if

$$|\theta_{(j)}| \leq A j^{-a}, \quad j = 1, 2, \dots, p, \quad |\theta_j| \leq A j^{-a}, \quad j = p + 1, p + 2, \dots,$$

which will be denoted as $\theta \in S_A^a(p)$, where $\{|\theta_{(j)}|, j = 1, \dots, p\}$ denotes the decreasing rearrangement of the numbers $\{|\theta_j|, j = 1, \dots, p\}$. Since $S_A^a \subset S_A^a(p)$, the second kind of smoothness is strictly more general than the first kind.

Here we use the term “smoothness” motivated by Fourier series analysis where smoothness of functions often translates into smoothness of the Fourier coefficients in the sense that is stated above; see, e.g., Kerkycharian and Picard (1992). For example, if a function $h : [0, 1]^d \mapsto \mathbb{R}$ possesses $r > 0$ continuous derivatives uniformly bounded by a constant M and the terms P_j are compactly supported Daubechies wavelets, then h can be represented as $h(z) = \sum_{j=1}^{\infty} P_j(z)\theta_{hj}$, with $|\theta_{hj}| \leq A j^{-r/d-1/2}$ for some constant A ; see Kerkycharian and Picard (1992). We also note that the second kind of smoothness is considerably more general than the first since it allows relatively large coefficients to appear anywhere in the series of the first p coefficients. In contrast, the first kind of smoothness only allows relatively large coefficients among the early terms in the series. Lasso-type methods are specifically designed to deal with the generalized smoothness of the second kind and perform equally well under both kinds of smoothness. In the context of series applications, smoothness of the second kind allows one to approximate functions that exhibit oscillatory phenomena or spikes, which are associated with “high-order” series terms. An example of this is the wage function example given in Belloni, Chernozhukov, and Hansen (2011).

Before we proceed to other examples we discuss a way to generate sparse approximations in infinite-dimensional examples. Consider, for example, a function h that can be represented a.s. as $h(z_i) = \sum_{j=1}^{\infty} \theta_{hj} P_j(z_i)$ with coefficients $\theta_h \in S_A^a(p)$. In this case we can construct sparse approximations by simply thresholding to zero all coefficients smaller than $1/\sqrt{n}$ and with indices $j \geq p$. This generates a sparsity index $s \leq A^{\frac{1}{a}} n^{\frac{1}{2a}}$. The non-zero coefficients could be further reoptimized by using the least squares projection. More formally, given a sparsity index $s > 0$, a target function $h(z_i)$, and terms $x_i = (P_j(z_i) : j = 1, \dots, p)' \in \mathbb{R}^p$, we let

$$\beta_{h0} := \arg \min_{\|\beta\|_0 \leq s} E[(h(z_i) - x_i' \beta)^2], \quad (4.31)$$

and define $x_i' \beta_{h0}$ as the best s -sparse approximation to $h(z_i)$.

Example 2. (Gaussian Model with Very Large p .) Consider (Ω, \mathcal{A}, P) as the probability space on which we have (y_i, z_i, d_i) as i.i.d. vectors for $i = 1, \dots, n$ obeying the model (4.28) with

$$\begin{aligned} g(z_i) &= \sum_{j=1}^{\infty} \theta_{gj} z_{ij}, \\ m(z_i) &= \sum_{j=1}^{\infty} \theta_{mj} z_{ij}. \end{aligned} \quad (4.32)$$

Assume that the infinite dimensional vector $w_i = (\zeta_i, v_i, z_i)'$ with j^{th} element denoted $w_i(j)$ is jointly Gaussian with covariance operator $[\text{Cov}(w_i(j), w_i(k))]_{j,k \geq 1}$ that has minimal and maximal eigenvalues bounded below by an absolute constant $\underline{\kappa} > 0$ and above by an absolute constant $\bar{\kappa} < \infty$.

The main assumption that guarantees approximate sparsity is the smoothness condition on the coefficients. Let $a > 1$ and $0 < A < \infty$ be absolute constants. We require that the coefficients of the expansions in (4.32) are a -smooth with constant A after p -rearrangement, namely

$$\theta_m = (\theta_{mj}, j = 1, 2, \dots) \in S_A^a(p), \quad \theta_g = (\theta_{gj}, j = 1, 2, \dots) \in S_A^a(p).$$

For estimation purposes we shall use $x_i = (z_{ij}, j = 1, \dots, p)'$, and assume that $|\alpha_0| \leq B$ and $p = p_n$ obeys

$$n^{\frac{1-a}{a} + \chi} \log^2(p \vee n) \leq \bar{\delta}_n, \quad A^{1/a} n^{\frac{1}{2a}} \leq p \bar{\delta}_n, \quad \text{and} \quad \log^3 p/n \leq \bar{\delta}_n,$$

for some sequence of positive constants $\bar{\delta}_n \searrow 0$ and absolute constants B and $\chi > 0$.

Corollary 4 (Gaussian Nonparametric Model). *Let \mathbf{P}_n be the collection of all dgp \mathbf{P} that obey the conditions set forth in Example 2 for a given n and for the given constants $(\underline{\kappa}, \bar{\kappa}, a, A, B, \chi)$ and sequences $p = p_n$ and $\bar{\delta}_n$. Then, as established in Appendix G, any $\mathbf{P} \in \mathbf{P}_n$ obeys Conditions ASTE (\mathbf{P}) with $s = A^{1/a} n^{\frac{1}{2a}}$, $SE(\mathbf{P})$, and $SM(\mathbf{P})$ for all $n \geq n_0$, with constants n_0 and $(\kappa', \kappa'', c, C)$ and sequences Δ_n and δ_n in those conditions depending only on $(\underline{\kappa}, \bar{\kappa}, a, A, B, \chi)$, p , and $\bar{\delta}_n$. Therefore, the conclusions of Theorem 1 hold for any sequence $\mathbf{P}_n \in \mathbf{P}_n$, and the conclusions of Corollary 1 on the uniform validity of confidence intervals apply uniformly for any $\mathbf{P} \in \mathbf{P}_n$. In particular, these conclusions apply uniformly in $\mathbf{P} \in \mathbf{P} = \bigcap_{n \geq n_0} \mathbf{P}_n$.*

Example 3. (Series Model with Very Large p .) Consider $(\Omega, \mathcal{A}, \mathbf{P})$ as the probability space, on which we have (y_i, z_i, d_i) as i.i.d. vectors for $i = 1, \dots, n$ obeying the model:

$$\begin{aligned} g(z_i) &= \sum_{j=1}^{\infty} \theta_{gj} P_j(z_i), \\ m(z_i) &= \sum_{j=1}^{\infty} \theta_{mj} P_j(z_i), \end{aligned} \tag{4.33}$$

where z_i has support $[0, 1]^d$ with density bounded from below by constant $\underline{f} > 0$ and above by constant \bar{f} , and $\{P_j, j = 1, 2, \dots\}$ is an orthonormal basis on $L^2[0, 1]^d$ with bounded elements, i.e. $\max_{z \in [0, 1]^d} |P_j(z)| \leq B$ for all $j = 1, 2, \dots$. Here all constants are taken to be absolute. Examples of such orthonormal bases include canonical trigonometric bases, e.g. $\{1, \sqrt{2} \cos(2\pi jz), \sqrt{2} \sin(2\pi jz) : j \geq 1\}$ where $z \in [0, 1]$.

Let $a > 1$ and $0 < A < \infty$ be absolute constants. We require that the coefficients of the expansions in (4.33) are a -smooth with constant A after p -rearrangement, namely

$$\theta_m = (\theta_{mj}, j = 1, 2, \dots) \in S_A^a(p), \quad \theta_g = (\theta_{gj}, j = 1, 2, \dots) \in S_A^a(p).$$

For estimation purposes we shall use $x_i = (P_j(z_i), j = 1, \dots, p)'$, and assume that $p = p_n$ obeys

$$n^{(1-a)/a} \log^2(p \vee n) \leq \bar{\delta}_n, \quad A^{1/a} n^{\frac{1}{2a}} \leq p \bar{\delta}_n \quad \text{and} \quad \log^3 p/n \leq \bar{\delta}_n,$$

for some sequence of absolute constants $\bar{\delta}_n \searrow 0$. We assume that there are some absolute constants $b > 0$, $B < \infty$, $q > 4$, with $(1-a)/a + 4/q < 0$, such that

$$|\alpha_0| \leq B, \quad b \leq E[\zeta_i^2 | x_i, v_i], \quad E[|\zeta_i^q| | x_i, v_i] \leq B, \quad b \leq E[v_i^2 | x_i], \quad E[|v_i^q| | x_i] \leq B. \tag{4.34}$$

Corollary 5 (Nonparametric Model with Sieve-type Regressors). *Let \mathbf{P}_n be the collection of all regression models P that obey the conditions set forth above for a given n . Then any $P \in \mathbf{P}_n$ obeys Conditions ASTE (P) with $s = A^{1/a} n^{\frac{1}{2a}}$, $SE(P)$, and $SM(P)$ for all $n \geq n_0$, with absolute constants in those conditions depending only on $(\underline{f}, \bar{f}, a, A, b, B, q)$ and $\bar{\delta}_n$. Therefore, the conclusions of Theorem 1 hold for any sequence $P_n \in \mathbf{P}_n$, and the conclusions of Corollary 1 on the uniform validity of confidence intervals apply uniformly for any $P \in \mathbf{P}_n$. In particular, as a special case, the same conclusion applies uniformly in $P \in \mathbf{P} = \cap_{n \geq n_0} \mathbf{P}_n$.*

4.2. Monte-Carlo Examples. In this section, we examine the finite-sample properties of the post-double-selection method in a partially linear model through a series of simulation exercises and compare its performance to that of a standard post-single-selection method.

All of the simulation results are based on the structural model

$$y_i = d_i' \alpha_0 + x_i' \theta_g + \sigma_y(d_i, x_i) \zeta_i, \quad \zeta_i \sim N(0, 1) \quad (4.35)$$

where $p = \dim(x_i) = 200$, the covariates $x_i \sim N(0, \Sigma)$ with $\Sigma_{kj} = (0.5)^{|j-k|}$, $\alpha_0 = .5$, and the sample size n is set to 100. In each design, we generate

$$d_i = x_i' \theta_m + \sigma_d(x_i) v_i, \quad v_i \sim N(0, 1) \quad (4.36)$$

with $E[\zeta_i v_i] = 0$. Inference results for all designs are based on conventional t-tests with standard errors calculated using the heteroscedasticity consistent jackknife variance estimator discussed in MacKinnon and White (1985).

We report results from three different dgp's. In the first two dgp's, we set $\theta_{g,j} = c_y \beta_{0,j}$ and $\theta_{m,j} = c_d \beta_{0,j}$ with $\beta_{0,j} = (1/j)^2$ for $j = 1, \dots, 200$. The first dgp, which we label "Design 1," uses homoscedastic innovations with $\sigma_y = \sigma_d = 1$. The second dgp, "Design 2," is heteroscedastic with $\sigma_{d,i} = \sqrt{\frac{(1+x_i' \beta_0)^2}{E_n(1+x_i' \beta_0)^2}}$ and $\sigma_{y,i} = \sqrt{\frac{(1+\alpha_0 d_i + x_i' \beta_0)^2}{E_n(1+\alpha_0 d_i + x_i' \beta_0)^2}}$. The constants c_y and c_d are chosen to generate desired population values for the reduced form R^2 's, i.e. the R^2 's for equations (2.6) and (2.7). For each equation, we choose c_y and c_d to generate $R^2 = 0, .2, .4, .6$, and $.8$. In the heteroscedastic design, we choose c_y and c_d based on R^2 as if (4.35) and (4.36) held with v_i and ζ_i homoscedastic and label the results by R^2 as in Design 1. In the third design ("Design 3"), we use a combination of deterministic and random coefficients. For the deterministic coefficients, we set $\theta_{g,j} = c_y (1/j)^2$ for $j \leq 5$ and $\theta_{m,j} = c_d (1/j)^2$ for $j \leq 5$. We then generate the remaining coefficients as iid draws from $(\theta_{g,j}, \theta_{m,j})' \sim N(0_{2 \times 1}, (1/p) I_2)$. For each equation, we choose c_y and c_d to generate $R^2 = 0, .2, .4, .6$, and $.8$ in the case that all of the random coefficients are exactly equal to 0 and label the results by R^2 as in Design 1. We draw new x 's, ζ 's, and v 's at every simulation replication, and we also generate new θ 's at every simulation replication in Design 3.

We consider Designs 1 and 2 to be baseline designs. These designs do not have exact sparse representations but have coefficients that decay quickly so that approximately sparse representations

are available. Design 3 is meant to introduce a modest deviation from the approximately sparse model towards a model with many small, uncorrelated coefficients. Using this we shall document that our proposed procedure still performs reasonably well, although it could be improved by incorporation of a ridge fit as one of regressors over which selection occurs.²¹

We report results for five different procedures. Two of the procedures are infeasible benchmarks: Oracle and Double-Selection Oracle estimators, which use knowledge of the true coefficient structures θ_g and θ_m and are thus unavailable in practice. The Oracle estimates α by running ordinary least squares of $y_i - x'_i\theta_g$ on d_i , and the Double-Selection Oracle estimates α by running ordinary least squares of $y - x'_i\theta_g$ on $d_i - x'_i\theta_m$. The other procedures we consider are feasible. One procedure is the standard post-single selection estimator – the Post-Lasso – which applies Lasso to equation (4.35) without penalizing α , the coefficient on d , to select additional control variables from among x . Estimates of α_0 are then obtained by OLS regression of y on d and the set of additional controls selected in the Lasso step and inference using the Post-Lasso estimator proceeds using conventional heteroscedasticity robust OLS inference from this regression. Post-Double-Selection or Post-Double-Lasso is the feasible procedure advocated in this paper. We run Lasso of y on x to select a set of predictors for y and run Lasso of d on x to select a set of predictors for d . α_0 is then estimated by running OLS regression of y on d and the union of the sets of regressors selected in the two Lasso runs, and inference is simply the usual heteroscedasticity robust OLS inference from this regression.²² Post-Double-Selection + Ridge is an *ad hoc* variant of Post-Double-Selection in which we add the ridge fit from equation (4.36) as an additional potential regressor that may be selected by Lasso. The ridge fit for d_i is $x'_i(X'X + \lambda_d I_p)^{-1}X'D$ where λ_d is obtained by 10-fold cross-validation. This procedure is motivated by a desire to add further robustness in the case that many small coefficients are suspected and zeroing out these small coefficients may be undesirable. Further exploration of procedures that perform well, both theoretically and in simulations, in the presence of many small coefficients is an interesting avenue for additional research.

We start by summarizing results in Table 1 for $(R_y^2, R_d^2) = (0, .2), (0, .8), (.8, .2)$, and $(.8, .8)$ where R_y^2 is the population R^2 from regressing y on x (Structure R^2) and R_d^2 is the population R^2 from regressing d on x (First Stage R^2). We report root-mean-square-error (RMSE) for estimating α_0 and size of 5% level tests (Rej. Rate). As should be the case, the Oracle and Double-Selection Oracle,

²¹In a supplementary appendix, we present results for 26 additional designs. The results presented in this section are sufficient to illustrate the general patterns from the larger set of results. In particular, the post-double-Lasso performed very well across all simulation designs where approximate sparsity provides a reasonable description of the dgp. Unsurprisingly, the performance deteriorates as one deviates from the smooth/approximately sparse case. However, the post-double-Lasso outperformed all other feasible procedures considered in all designs.

²²All Lasso estimates require the choice of penalty parameter and loadings. We use the iterative procedure of Belloni, Chen, Chernozhukov, and Hansen (2012) to estimate the penalty loadings using a maximum of five iterations. We set the penalty parameter according to equation (18) in Belloni and Chernozhukov (2011b) with $c = 1.1$, $\alpha = .05$ and $\sigma = 1$ since the variance of the score is accounted for in the penalty loadings.

which are reported to provide the performance of an infeasible benchmark, perform well relative to the feasible procedures across the three designs. We do see that the feasible Post-Double-Selection procedures perform similarly to the Double-Selection Oracle without relying on *ex ante* knowledge of the coefficients that go in to the control functions, θ_g and θ_m . On the other hand, the Post-Lasso procedure generally does not perform as well as Post-Double-Selection and is very sensitive to the value of R_d^2 . While Post-Lasso performs adequately when R_d^2 is small, its performance deteriorates quickly as R_d^2 increases. This lack of robustness of traditional variable selection methods such as Lasso which were designed with forecasting, not inference about treatment effects, in mind is the chief motivation for our advocating the Post-Double-Selection procedure when trying to infer structural or treatment parameters.

We provide further details about the performance of the feasible estimators in Figures 1, 2, and 3 which plot size of 5% level tests, bias, and standard deviation for the Post-Lasso, Double-Selection (DS), and Double-Selection Oracle (DS Oracle) estimators of the treatment effect across the full set of R^2 values considered. Figure 1, 2, and 3 respectively report the results from Design 1, 2, and 3. The figures are plotted with the same scale to aid comparability, and rejection frequencies for Post-Lasso were censored at .5 for readability. Perhaps the most striking feature of the figures is the poor performance of the Post-Lasso estimator. The Post-Lasso estimator performs poorly in terms of size of tests across many different R^2 combinations and can have an order of magnitude more bias than the corresponding Post-Double-Selection estimator. The behavior of Post-Lasso is quite non-uniform across R^2 combinations, and Post-Lasso does not reliably control size distortions or bias except in the case where the controls are uncorrelated with the treatment (where First-Stage R^2 equals 0) and thus ignorable. In contrast, the Post-Double-Selection estimator performs relatively well across the full range of R^2 combinations considered. The Post-Double-Selection estimator's performance is also quite similar to that of the infeasible Double-Selection Oracle across the majority of R^2 values considered. Comparing across Figures 1 and 2, we see that size distortions for both the Post-Double-Selection estimator and the Double-Selection Oracle are somewhat larger in the presence of heteroscedasticity but that the basic patterns are more-or-less the same across the two figures. Looking at Figure 3, we also see that the addition of small independent random coefficients results in somewhat larger size distortions for the Post-Double-Selection estimator than in the other homoscedastic design, Design 1, though the procedure still performs relatively well.

In the final figure, Figure 4, we compare the performance of the Post-Double-Selection procedure to the *ad hoc* Post-Double-Selection procedure which selects among the original set of variables augmented with the ridge fit obtained from equation (4.36). We see that the addition of this variable does add robustness relative to Post-Double-Selection using only the raw controls in the sense of producing tests that tend to have size closer to the nominal level. This additional robustness is a good feature, though it comes at the cost of increased RMSE which is especially prominent for small values of the first-stage R^2 .

The simulation results are favorable to the Post-Double-Selection estimator. In the simulations, we see that the Post-Double-Selection procedure provides an estimator of a treatment effect in the presence of a large number of potential confounding variables that performs similarly to the infeasible estimator that knows the values of the coefficients on all of the confounding variables. Overall, the simulation evidence supports our theoretical results and suggests that the proposed Post-Double-Selection procedure can be a useful tool to researchers doing structural estimation in the presence of many potential confounding variables. It also shows, as a contrast, that the standard Post-Single-Selection procedure provides poor inference and therefore is not a reliable tool to these researchers.

5. EXTENSIONS: OTHER PROBLEMS AND HETEROGENEOUS TREATMENT EFFECTS

5.1. Other Problems. In order to discuss extensions in a very simple manner, we assume i.i.d sampling as well as assume away approximation errors; i.e. we let $g(z_i) = x_i' \beta_{g0}$ and $m(z_i) = x_i' \beta_{m0}$ where parameters β_{g0} and β_{m0} are high-dimensional and let $x_i = P(z_i)$. In the development of results for the partially linear model, we implicitly considered a moment condition for the target parameter α_0 given by

$$E[\varphi(y_i - d_i \alpha_0 - x_i' \beta) v_i] = 0, \quad (5.37)$$

where $\varphi(u) = u$ and v_i are measurable functions of z_i . We selected the instrument v_i such that the equation is first-order insensitive to the parameter β at $\beta = \beta_{g0}$:

$$\frac{\partial}{\partial \beta} E[\varphi(y_i - d_i \alpha_0 - x_i' \beta) v_i] \Big|_{\beta = \beta_{g0}} = 0. \quad (5.38)$$

Note that $\varphi(u) = u$ and $v_i = d_i - m(z_i)$ implement this condition. If (5.38) holds, the estimator of α_0 gets “immunized” against nonregular estimation of β_0 , for example, via a post-selection procedure or other regularized estimators. Such immunization ideas are in fact behind the classical Frisch-Waugh-Robinson partialling out technique in the linear setting and the Neyman (1979)’s $C(\alpha)$ test in the nonlinear setting. One way to view our contribution is as a recognition of the importance of this immunization in the context of post-selection inference leading to thinking about a post-selection approach to inference on the target parameter aimed at satisfying this condition. Our approach uses modern selection methods to estimate g and the function m defining the instrument v . Generalizations to nonlinear models, where φ is non-linear and can correspond to a likelihood score or quantile check function are given in Belloni, Chernozhukov, and Kato (2013) and Belloni, Chernozhukov, and Wei (2013); in these generalizations, achieving (5.38) is also critical.

Within the context of this paper, a potentially important extension is to consider a general treatment effect model, where d_i is interacted with transformations of z_i . As long as the interest lies in a particular regression coefficient, the current framework covers this implicitly since x_i could contain interactions of d_i with transformations of controls z_i . In the case where a fixed number of

such regression coefficients is of interest, we can estimate each of the coefficients by re-labeling the corresponding regressor as d_i and other regressors as x_i and then applying our procedure under each re-labeling. This component-wise procedure is valid as long as the assumed regularity conditions hold for each of the small number of resulting regression models. One can also accomodate vector d_i in this fashion.

A related topic is the estimation of average treatment effects when treatment effects are fully heterogeneous. When the treatment variable $d_i \in \{0, 1\}$ is binary (or discrete more generally), our approach is readily amenable to this problem. In this case, a parameter of interest is the average treatment effect,

$$\alpha_0 = E[g(1, z_i) - g(0, z_i)],$$

where $g(d_i, z_i) = E[y_i | d_i, z_i]$. We can write $g(d_i, z_i) = d_i x_i' \beta_{g0,1} + (1 - d_i) x_i' \beta_{g0,0}$. If the propensity score $P(d_i = 1 | z_i)$ is $m(z_i) = \Lambda(x_i' \beta_{m0})$ where $\Lambda(u)$ is a link such as logit or linear, we can use the moment equations of Hahn (1998),

$$E[\varphi(\alpha_0, y_i, d_i, g(0, z_i), g(1, z_i), m(z_i))] = 0, \quad (5.39)$$

where $\varphi(\alpha, y, d, g_0, g_1, m) = \alpha - \frac{d(y-g_1)}{m} + \frac{(1-d)(y-g_0)}{1-m} - (g_1 - g_0)$ to estimate α_0 . It is straightforward to check that

$$\frac{\partial}{\partial \beta_j} E[\varphi_i(\alpha_0, y_i, d_i, x_i' \beta_1, x_i' \beta_2, x_i' \beta_3)] = 0, \quad (5.40)$$

for each $j \in \{1, 2, 3\}$ when coefficients are set equal to their true values, i.e. when $(\beta_1, \beta_2, \beta_3) = (\beta_{g0,0}, \beta_{g0,1}, \beta_{m0})$.²³ Thus, we reduce dependence on the estimated values of β_{g0} and β_{m0} by using Hahn (1998)'s equation just as in the partially linear case. This property suggests that one can use the selection approach to regularization in order to estimate the parameter of interest α_0 . Note that one could adopt a double-selection method as in the partially linear model where selection is over terms explaining the propensity score and terms explaining the regression function. Using the results of this paper and setting $\Lambda(u) = u$, it is not difficult to show that the post-double-selection estimator $\check{\alpha}$ that solves $E_n[\varphi(\check{\alpha}, y_i, d_i, \hat{g}(0, z_i), \hat{g}(1, z_i), \hat{m}(z_i))] = 0$ satisfies

$$\sigma_n^{-1} \sqrt{n}(\check{\alpha} - \alpha_0) \rightsquigarrow N(0, 1), \quad \sigma_n^2 = E[\varphi^2(\alpha_0, y_i, d_i, g(0, z_i), g(1, z_i), m(z_i))], \quad (5.41)$$

where σ_n^2 is the semiparametric efficiency bound of Hahn (1998) assuming that both g and m satisfy the sparsity conditions assumed previously and that additional assumptions needed to guarantee consistency of post-Lasso estimators \hat{g} and \hat{m} in the uniform norm are satisfied.

²³In fact, some higher order derivatives also vanish. This higher order property can be exploited in conjunction with sample-splitting to relax requirements on s . The split-sample approach would estimate the parameters of the conditional expectation functions using one part of the sample, then use these estimates with the data from the other part of the sample to form the moment condition for that part of the sample, and then use this moment condition to estimate α_0 . This procedure would then be repeated switching the roles of the two parts of the sample. Combining the two resulting estimators would result in a fully efficient estimator. We do not discuss the details of such an approach here for brevity.

5.2. Theoretical Results on ATE with Heterogeneity. Consider an i.i.d. sample $(y_i, d_i, z_i)_{i=1}^n$ on the probability space $(\Omega, \mathfrak{F}, P)$. Suppose the treatment variable is binary, $d_i \in \{0, 1\}$, and that the outcome and propensity equations are

$$y_i = g(d_i, z_i) + \zeta_i, \quad E[\zeta_i \mid z_i, d_i] = 0, \quad (5.42)$$

$$d_i = m(z_i) + v_i, \quad E[v_i \mid z_i] = 0. \quad (5.43)$$

The first target parameter is the average treatment effect, $\alpha_0 = E[g(1, z_i) - g(0, z_i)]$, which is implicitly indexed by P like other parameters. In this model, d_i is not additively separable. The purpose of this section is to show that our analysis easily extends to this case.

The confounding factors z_i affect the policy variable via the propensity score $m(z_i)$ and the outcome variable via the function $g(d_i, z_i)$. Both of these functions are unknown and potentially complicated. As in the main text, we use linear combinations of control terms $x_i = P(z_i)$ to approximate $g(z_i)$ and $m(z_i)$, writing (5.42) and (5.43) as

$$y_i = \underbrace{\tilde{x}_i' \beta_{g0} + r_{gi}}_{g(d_i, z_i)} + \zeta_i, \quad (5.44)$$

$$d_i = \underbrace{\Lambda(\tilde{x}_i' \beta_{m0}) + r_{mi}}_{m(z_i)} + v_i, \quad (5.45)$$

where r_{gi} and r_{mi} are the approximation errors, and

$$\tilde{x}_i := (d_i x_i', (1 - d_i) x_i')', \quad \beta_{g0} := (\beta_{g0,1}', \beta_{g0,0}')', \quad x_i := P(z_i), \quad (5.46)$$

where $x_i' \beta_{g0,1}$, $x_i' \beta_{g0,0}$, and $x_i' \beta_{m0}$ are approximations to $g(1, z_i)$, $g(0, z_i)$, and $m(z_i)$, and $\Lambda(u) = u$ for the case of linear link and $\Lambda(u) = e^u / (1 + e^u)$ for the case of the logistic link. In order to allow for a flexible specification and incorporation of pertinent confounding factors, the dimension of the vector of controls, $x_i = P(z_i)$, can be large relative to the sample size.

Using the efficient moment condition (5.39) derived by Hahn (1998), we can define the post-double-selection estimator $\check{\alpha}$ as the solution to

$$\mathbb{E}_n [\varphi(\check{\alpha}, y_i, d_i, \hat{g}(0, z_i), \hat{g}(1, z_i), \hat{m}(z_i))] = 0, \quad (5.47)$$

where $\hat{g}(d_i, z_i)$ and $\hat{m}(z_i)$ are post-Lasso estimators of functions g and m based upon equations (5.44)-(5.45). In case of the logistic link Λ , Lasso for logistic regression is as defined in van de Geer (2008) and Bach (2010), and the associated post-Lasso estimators are as defined in Belloni, Chernozhukov, and Wei (2013).

In what follows, we use $\|w_i\|_{P,q}$ to denote the $L^q(P)$ norm of a random variable w_i with law determined by P , and we $\|w_i\|_{\mathbb{P}_n,q}$ to denote the empirical $L^q(\mathbb{P}_n)$ norm of a random variable with law determined by the empirical measure $\mathbb{P}_n = n^{-1} \sum_{i=1}^n \delta_{w_i}$, i.e., $\|w_i\|_{\mathbb{P}_n,q} = (n^{-1} \sum_{i=1}^n \|w_i\|^q)^{1/q}$.

Consider fixed sequences of positive numbers $\delta_n \searrow 0$ and $\Delta_n \searrow 0$ and constants $C > 0, c > 0, 1/2 > c' > 0$ which will not vary with P .

Condition HTE (P) . Heterogeneous Treatment Effects. *Consider i.i.d. sample $(y_i, d_i, z_i)_{i=1}^n$ on the probability space $(\Omega, \mathfrak{F}, P)$, where we shall call P the data-generating process, such that equations (5.44)-(5.45) holds, with $d_i \in \{0, 1\}$. (i) Approximation errors satisfy $\|r_{gi}\|_{P,2} \leq \delta_n n^{-1/4}$, $\|r_{gi}\|_{P,\infty} \leq \delta_n$, and $\|r_{mi}\|_{P,2} \leq \delta_n n^{-1/4}$, $\|r_{mi}\|_{P,\infty} \leq \delta_n$. (ii) With P -probability no less than $1 - \Delta_n$, estimation errors satisfy $\|\tilde{x}'_i(\hat{\beta}_g - \beta_{g0})\|_{\mathbb{P}_{n,2}} \leq \delta_n n^{-1/4}$, $\|x'_i(\hat{\beta}_m - \beta_{m0})\|_{\mathbb{P}_{n,2}} \leq \delta_n n^{-1/4}$, $K_n \|\hat{\beta}_m - \beta_m\|_1 \leq \delta_n$, $K_n \|\hat{\beta}_m - \beta_{m0}\|_1 \leq \delta_n$, estimators and approximations are sparse, namely $\|\hat{\beta}_g\|_0 \leq Cs$, $\|\hat{\beta}_m\|_0 \leq Cs$, $\|\beta_{g0}\|_0 \leq Cs$, and $\|\beta_{m0}\|_0 \leq Cs$, and the empirical and populations norms are equivalent on sparse subsets, namely $\sup_{\|\delta\|_0 \leq 2Cs} \|\tilde{x}'_i \delta\|_{\mathbb{P}_{n,2}} / \|\tilde{x}'_i \delta\|_{P,2} - 1 \leq \delta_n$. (iii) The following boundedness conditions hold: $\|x_{ij}\|_{P,\infty} \leq K_n$ for each j , $\|g\|_{P,\infty} \leq C$, $\|y_i\|_{P,\infty} \leq C$, $P(c' \leq m(z_i) \leq 1 - c') = 1$, and $\|\zeta_i^2\|_{P,2} \geq c$. (iv) The sparsity index obeys the following growth condition, $(s \log(p \vee n))^2/n \leq \delta_n$.*

These conditions are simple high-level conditions which encode both the approximate sparsity of the models as well as impose some reasonable behavior on the sparse estimators of m and g . These conditions are implied by other more primitive conditions in the literature; see van de Geer (2008) and Belloni, Chen, Chernozhukov, and Hansen (2012). Sufficient conditions for the equivalence between population and empirical norms over sparse subsets follow from Lemmas 10 and 11. The boundedness conditions are made to simplify arguments, and they could be removed at the cost of more complicated proofs and more stringent side conditions.

Theorem 3 (Uniform Post-Double Selection Inference on ATE). *Consider the set \mathbf{P}_n of data generating processes P such that equations (5.42)-(5.43) and Condition HTE (P) holds. (1) Then under any sequence $P \in \mathbf{P}_n$,*

$$\sigma_n^{-1} \sqrt{n}(\tilde{\alpha} - \alpha_0) \rightsquigarrow N(0, 1), \quad \sigma_n^2 = E[\varphi^2(\alpha_0, y_i, d_i, g(0, z_i), g(1, z_i), m(z_i))]. \quad (5.48)$$

- (2) *The result continues to hold with σ_n^2 replaced by $\hat{\sigma}_n^2 := E_n[\varphi^2(\alpha_0, y_i, d_i, \hat{g}(0, z_i), \hat{g}(1, z_i), \hat{m}(z_i))]$.*
(3) *Moreover, the confidence regions based upon post-double selection estimator $\tilde{\alpha}$ have uniform asymptotic validity:*

$$\lim_{n \rightarrow \infty} \sup_{P \in \mathbf{P}_n} |P(\alpha_0 \in [\tilde{\alpha} \pm \Phi^{-1}(1 - \xi/2) \hat{\sigma}_n / \sqrt{n}]) - (1 - \xi)| = 0.$$

The next target parameter is the average treatment effect on the treated:

$$\gamma_0 = E[g(1, z_i) - g(0, z_i) | d_i = 1].$$

The efficient moment condition, derived by Hahn (1998), for parameter γ_0 is as follows:

$$E[\tilde{\varphi}(\gamma_0, y_i, d_i, g(0, z_i), g(1, z_i), m(z_i), \mu)] = 0, \quad (5.49)$$

where $\mu = \mathbb{E}[m(z_i)] = \mathbb{P}(d_i = 1)$, and

$$\tilde{\varphi}(\gamma, y, d, g_0, g_1, m, \mu) = \frac{d(y - g_1)}{\mu} - \frac{m(1 - d)(y - g_0)}{(1 - m)\mu} + \frac{d(g_1 - g_0)}{\mu} - \gamma \frac{d}{\mu}.$$

In this case the post-double-selection estimator $\tilde{\gamma}$ that solves

$$\mathbb{E}_n [\tilde{\varphi}(\tilde{\gamma}, y_i, d_i, \hat{g}(0, z_i), \hat{g}(1, z_i), \hat{m}(z_i), \hat{\mu})] = 0, \quad (5.50)$$

where $\hat{g}(d_i, z_i)$ and $\hat{m}(z_i)$ are post-Lasso estimators or other sparse estimators obeying the regularity conditions posed in HTE of functions g and m based upon equations (5.44)-(5.45), and $\hat{\mu} = \mathbb{E}_n[d_i]$. As above, post-Lasso estimators in the case of the logistic link are defined in Belloni, Chernozhukov, and Wei (2013).

Theorem 4 (Uniform Post-Double Selection Inference on ATT). *Consider the set \mathbf{P}_n of data generating processes P such that equations (5.42)-(5.43) and Condition HTE (P) holds. (1) Then under any sequence $P \in \mathbf{P}_n$,*

$$\sigma_n^{-1} \sqrt{n}(\tilde{\gamma} - \gamma_0) \rightsquigarrow N(0, 1), \quad \sigma_n^2 = \mathbb{E}[\tilde{\varphi}^2(\gamma_0, y_i, d_i, g(0, z_i), g(1, z_i), m(z_i), \mu)]. \quad (5.51)$$

- (2) *The result continues to hold with σ_n^2 replaced by $\hat{\sigma}_n^2 := \mathbb{E}_n[\tilde{\varphi}^2(\gamma_0, y_i, d_i, \hat{g}(0, z_i), \hat{g}(1, z_i), \hat{m}(z_i), \mu)]$.*
 (3) *Moreover, the confidence regions based upon post-double selection estimator $\tilde{\alpha}$ have uniform asymptotic validity:*

$$\lim_{n \rightarrow \infty} \sup_{P \in \mathbf{P}_n} |\mathbb{P}(\gamma_0 \in [\tilde{\gamma} \pm \Phi^{-1}(1 - \xi/2) \hat{\sigma}_n / \sqrt{n}]) - (1 - \xi)| = 0.$$

These results contribute to recent results on estimation of the ATE and its variants in Cattaneo (2010), who considers this problem in a series framework where the number of series terms obeys $p^2/n \rightarrow 0$, and in Rothe and Firpo (2013), who provide results based on kernel estimators. These approaches are very useful but do not target “data-rich-environments” and do not study uniformity. Our framework allows for consideration of a large number of series terms, potentially much larger than the sample size, but requires that a relatively small number of these terms are needed through the sparsity condition $s^2 K_n (\log(p \vee n))^2 / n \rightarrow 0$. Our framework also covers semi-parametric models with a large number of raw regressors x_{ij} as long as $|x_{ij}| \leq K_n$ where K_n does not grow too quickly. Finally, we establish validity of our inferential results uniformly in P . Hence, the results provided above offer a useful contribution to this literature.

6. EMPIRICAL EXAMPLE: ESTIMATING THE EFFECT OF ABORTION ON CRIME

In the preceding sections, we have provided results demonstrating how variable selection methods, focusing on the case of Lasso-based methods, can be used to estimate treatment effects in models in which we believe the variable of interest is exogenous conditional on observables. We further illustrate the use of these methods in this section by reexamining Donohue III and Levitt’s (2001)

study of the impact of abortion on crime rates. In the following, we briefly review Donohue III and Levitt (2001) and then present estimates obtained using the methods developed in this paper.

Donohue III and Levitt (2001) discuss two key arguments for a causal channel relating abortion to crime. The first is simply that more abortion among a cohort results in an otherwise smaller cohort and so crime 15 to 25 years later, when this cohort is in the period when its members are most at risk for committing crimes, will be otherwise lower given the smaller cohort size. The second argument is that abortion gives women more control over the timing of their fertility allowing them to more easily assure that childbirth occurs at a time when a more favorable environment is available during a child's life. For example, access to abortion may make it easier to ensure that a child is born at a time when the family environment is stable, the mother is more well-educated, or household income is stable. This second channel would mean that more access to abortion could lead to lower crime rates even if fertility rates remained constant.

The basic problem in estimating the causal impact of abortion on crime is that state-level abortion rates are not randomly assigned, and it seems likely that there will be factors that are associated to both abortion rates and crime rates. It is clear that any association between the current abortion rate and the current crime rate is likely to be spurious. However, even if one looks at say the relationship between the abortion rate 18 years in the past and the crime rate among current 18 year olds, the lack of random assignment makes establishing a causal link difficult without adequate controls. An obvious confounding factor is the existence of persistent state-to-state differences in policies, attitudes, and demographics that are likely related to overall state level abortion and crime rates. It is also important to control flexibly for aggregate trends. For example, it could be the case that national crime rates were falling over some period while national abortion rates were rising but that these trends were driven by completely different factors. Without controlling for these trends, one would mistakenly associate the reduction in crime to the increase in abortion. In addition to these overall differences across states and times, there are other time varying characteristics such as state-level income, policing, or drug-use to name a few that could be associated with current crime and past abortion.

To address these confounds, Donohue III and Levitt (2001) estimate a model for state-level crime rates running from 1985 to 1997 in which they condition on a number of these factors. Their basic specification is

$$y_{cit} = \alpha_c a_{cit} + w'_{it} \beta_c + \delta_{ci} + \gamma_{ct} + \varepsilon_{cit} \quad (6.52)$$

where i indexes states, t indexes times, $c \in \{\text{violent, property, murder}\}$ indexes type of crime, δ_{ci} are state-specific effects that control for any time-invariant state-specific characteristics, γ_{ct} are time-specific effects that control flexibly for any aggregate trends, w_{it} are a set of control variables to control for time-varying confounding state-level factors, a_{cit} is a measure of the abortion

rate relevant for type of crime c ,²⁴ and y_{cit} is the crime-rate for crime type c . Donohue III and Levitt (2001) use the log of lagged prisoners per capita, the log of lagged police per capita, the unemployment rate, per-capita income, the poverty rate, AFDC generosity at time $t - 15$, a dummy for concealed weapons law, and beer consumption per capita for w_{it} , the set of time-varying state-specific controls. Tables IV and V in Donohue III and Levitt (2001) present baseline estimation results based on (6.52) as well as results from different models which vary the sample and set of controls to show that the baseline estimates are robust to small deviations from (6.52). We refer the reader to the original paper for additional details, data definitions, and institutional background.

For our analysis, we take the argument that the abortion rates defined above may be taken as exogenous relative to crime rates once observables have been conditioned on from Donohue III and Levitt (2001) as given. Given the seemingly obvious importance of controlling for state and time effects, we account for these in all models we estimate. We choose to eliminate the state effects via differencing rather than including a full set of state dummies but include a full set of time dummies in every model.²⁵ Thus, we will estimate models of the form

$$y_{cit} - y_{cit-1} = \alpha_c(a_{cit} - a_{cit-1}) + z'_{cit}\kappa_c + g_{ct} + \eta_{cit}. \quad (6.53)$$

where g_{ct} are time effects. We use the same state-level data as Donohue III and Levitt (2001) but delete Alaska, Hawaii, and Washington, D.C. which gives a sample with 48 cross-sectional observations and 12 time series observations for a total of 576 observations. With these deletions, our baseline estimates using the same controls as in (6.52) are quite similar to those reported in Donohue III and Levitt (2001). Baseline estimates from Table IV of Donohue III and Levitt (2001) and our baseline estimates based on the differenced version of (6.52) are given in the first and second row of Table 2 respectively.

Our main point of departure from Donohue III and Levitt (2001) is that we allow for a much richer set z_{cit} than allowed for in w_{it} in model (6.52). Our z_{cit} includes higher-order terms and interactions of the control variables defined above. In addition, we put initial conditions and initial differences of w_{it} and a_{cit} and within-state averages of w_{it} into our vector of controls z_{cit} . This

²⁴This variable is constructed as weighted average of abortion rates where weights are determined by the fraction of the type of crime committed by various age groups. For example, if 60% of violent crime were committed by 18 year olds and 40% were committed by 19 year olds in state i , the abortion rate for violent crime at time t in state i would be constructed as .6 times the abortion rate in state i at time $t - 18$ plus .4 times the abortion rate in state i at time $t - 19$. See Donohue III and Levitt (2001) for further detail and exact construction methods.

²⁵Part of the motivation for considering first-differences is that our theoretical results are for independent data. For both violent crime and property crime, this assumption seems like a better approximation in differences than in levels. The first three estimated autocorrelations of the first-difference residuals from the baseline specification using only the controls from Donohue III and Levitt (2001) based on violent crime, property crime, and murder are respectively (.0155, .0574, -.0487), (-.0736, .0651, .0540), and (-.3954, -.0813, .0066). Discussion of results obtained estimating the model in levels and using fixed effects are available in a supplementary appendix. Extending the formal results to accommodate dependence would be a useful extension for future work.

addition allows for the possibility that there may be some feature of a state that is associated both with its growth rate in abortion and its growth rate in crime. For example, having an initially high-levels of abortion could be associated with having high-growth rates in abortion and low growth rates in crime. Failure to control for this factor could then lead to misattributing the effect of this initial factor, perhaps driven by policy or state-level demographics, to the effect of abortion. Finally, we allow for more general trends by allowing for an aggregate quadratic trend in z_{cit} as well as interactions of this quadratic trend with control variables. This gives us a set of 284 control variables to select among in addition to the 12 time effects that we include in every model.²⁶

Note that interpreting estimates of the effect of abortion from model (6.52) as causal relies on the belief that there are no higher-order terms of the control variables, no interaction terms, and no additional excluded variables that are associated both to crime rates and the associated abortion rate. Thus, controlling for a large set of variables as described above is desirable from the standpoint of making this belief more plausible. At the same time, naively controlling lessens our ability to identify the effect of interest and thus tends to make estimates far less precise. The effect of estimating the abortion effect conditioning on the full set of 284 potential controls described above is given in the third row of Table 2. As expected, all coefficients are estimated very imprecisely. Of course, very few researchers would consider using 284 controls with only 576 observations due to exactly this issue.

We are faced with a tradeoff between controlling for very few variables which may leave us wondering whether we have included sufficient controls for the exogeneity of the treatment and controlling for so many variables that we are essentially mechanically unable to learn about the effect of the treatment. The variable selection methods developed in this paper offer one resolution to this tension. The assumed sparse structure maintains that there is a small enough set of variables that one could potentially learn about the treatment but adds substantial flexibility to the usual case where a researcher considers only a few control variables by allowing this set to be found by the data from among a large set of controls. Thus, the approach should complement the usual careful specification analysis by providing a researcher an efficient, data-driven way to search for a small set of influential confounds from among a sensibly chosen broad set of potential confounding variables.

In the abortion example, we use the post-double-selection estimator defined in Section 2.2 for each of our dependent variables.²⁷ For violent crime, six variables are selected in the abortion

²⁶The exact identities of the 284 potential controls is available upon request. It consists of linear and quadratic terms of each continuous variable in w_{it} , interactions of every variable in w_{it} , initial levels and initial differences of w_{it} and a_{cit} , the within-state averages of w_{it} , and interactions of these variables with a quadratic trend.

²⁷Implementation requires selection of a penalty parameter and loadings. We estimate the loadings using the iterative procedure proposed in Belloni, Chen, Chernozhukov, and Hansen (2012) with 100 as the maximum number of iterations. For each model, the iterative procedure converges after 21 or fewer iterations. We set the penalty parameter according to (2.12) with $c = 1.1$ and $\gamma = .05$.

equation,²⁸ and no variables are selected in the crime equation. For property crime, seven variables are selected in the abortion equation,²⁹ and two are selected in the crime equation.³⁰ For murder, six variables are selected in the abortion equation,³¹ and none were selected in the crime equation.

Estimates of the causal effect of abortion on crime obtained by searching for confounding factors among our set of 284 potential controls are given in the fourth row of Table 2. Each of these estimates is obtained from the least squares regression of the crime rate on the abortion rate and the six, nine, and six controls selected by the double-post-Lasso procedure for violent crime, property crime, and murder respectively. All of these estimates for the effect of abortion on crime rates are quite imprecise, producing 95% confidence intervals that encompass large positive and negative values. Note that the double-post-Lasso produces models that are not of vastly different size than the “intuitive” model (6.52). As a final check, we also report results that include all of the original variables from (6.52) in the amelioration set in the fifth row of the table. These results show that the conclusions made from using only the variable selection procedure do not qualitatively change when the variables used in the original Donohue III and Levitt (2001) are added to the equation. For a quick benchmark relative to the simulation examples, we note that the R^2 obtained by regressing the crime rate on the selected variables are .0251, .1179, and .0039 for violent crime, property crime, and the murder rate respectively and that the R^2 ’s from regressing the abortion rate on the selected variables are .8420, .6116, and .7781 for violent crime, property crime, and the murder rate respectively. These values correspond to regions of the R^2 space considered in the simulation where the double selection procedure substantially outperformed simple Lasso procedures.

It is interesting that one would draw qualitatively different conclusions from the estimates obtained using formal variable selection than from the estimates obtained using a small set of intuitively selected controls. Looking at the set of selected control variables, we see that initial conditions and interactions with trends are selected across all dependent variables. We also see that we cannot precisely determine the effect of the abortion rate on crime rates once one accounts for these initial conditions. Of course, this does not mean that the effects of the abortion rate provided in the first two rows of Table 2 are not representative of the true causal effects. It does,

²⁸The selected variables are lagged prisoners per capita, the initial change in beer consumption interacted with a linear trend, the initial change in income squared interacted with a linear trend, the within-state mean of income, the within-state mean of income interacted with a linear trend, and the initial level of the abortion rate.

²⁹The selected variables are lagged prisoners per capita, lagged income, the initial level of income, the initial change of income squared interacted with a linear trend, the within-state average of income, the within-state average of income interacted with a linear trend, and the initial level of abortion interacted with a linear trend.

³⁰The two variables are the initial level income squared interacted with a linear trend and the within-state average of AFDC generosity.

³¹The selected variables are lagged unemployment, the initial change in unemployment squared, the initial level of prisoners per capita, the within-state average of the number of prisoners per capita interacted with a linear trend, the within-state average of income interacted with a linear trend, and the initial level of the abortion rate interacted with a linear trend.

however, imply that this conclusion is strongly predicated on the belief that there are not other unobserved state-level factors that are correlated to both initial values of the controls and abortion rates, abortion rate changes, and crime rate changes. Interestingly, a similar conclusion is given in Foote and Goetz (2008) based on an intuitive argument.

We believe that the example in this section illustrates how one may use modern variable selection techniques to complement causal analysis in economics. In the abortion example, we are able to search among a large set of controls and transformations of variables when trying to estimate the effect of abortion on crime. Considering a large set of controls makes the underlying assumption of exogeneity of the abortion rate conditional on observables more plausible, while the methods we develop allow us to produce an end-model which is of manageable dimension. Interestingly, we see that one would draw quite different conclusions from the estimates obtained using formal variable selection. Looking at the variables selected, we can also see that this change in interpretation is being driven by the variable selection method’s selecting different variables, specifically initial values of the abortion rate and controls, than are usually considered. Thus, it appears that the usual interpretation hinges on the prior belief that initial values should be excluded from the structural equation for the differences.

7. CONCLUSION

In this paper, we consider estimation of treatment effects or structural parameters in a partially linear model, where the treatment is believed to be exogenous conditional on observables. We do not impose the conventional assumption that the identities of the relevant conditioning variables and the functional form with which they enter the model are known. Rather, we assume that the researcher believes there is a relatively small number of important factors whose identities are unknown within a much larger known set of potential variables and transformations. This sparsity assumption is used explicitly or implicitly in much of empirical research in economics. It allows the researcher to estimate the desired treatment effect and infer a set of important variables upon which one needs to condition by using modern variable selection techniques without *ex ante* knowledge of which are the important conditioning variables. Since naive application of variable selection methods in this context may result in very poor properties for inferring the treatment effect of interest, we propose a “double-selection” estimator of the treatment effect, provide a formal demonstration of its properties for estimating the treatment effect, and provide its approximate distribution under technical regularity conditions and the assumed sparsity in the model.

In addition to the theoretical development, we illustrate the potential usefulness of our proposal through a number of simulation studies and an empirical example. In Monte Carlo simulations, our procedure outperforms a simple variable selection strategy for estimating the treatment effect across the designs considered and does relatively well compared to an infeasible estimator that uses

the identities of the relevant conditioning variables. We then apply our estimator to attempt to estimate the causal impact of abortion on crime following Donohue III and Levitt (2001). We find that our procedure selects a small number of conditioning variables. After conditioning on these selected variables, one would draw qualitatively different inference about the effect of abortion on crime than would be drawn if one assumed that the correct set of conditioning variables was known and the same as those variables used in Donohue III and Levitt (2001). Taken together, the empirical and simulation examples demonstrate that the proposed method may provide a useful complement to other sorts of specification analysis done in applied research.

APPENDIX A. ITERATED LASSO ESTIMATION

Feasible implementation of Lasso under heteroscedasticity requires a choice of penalty parameter λ and estimation of penalty loadings (2.14). λ depends only on the known p and n and the researcher specified c and γ . In all examples, we use $c = 1.1$ and $\gamma = .05$. In this appendix, we state algorithms for estimating the penalty loadings.

Let I_0 be an initial set of regressors with a bounded number of elements, including for example the intercept. Let $\bar{\beta}(I_0)$ be the least squares estimator of the coefficients on the covariates associated with I_0 , and define $\hat{l}_{jI_0} := \sqrt{\mathbb{E}_n[x_{ij}^2(y_i - x_i'\bar{\beta}(I_0))^2]}$.

An algorithm for estimating the penalty loadings using Post-Lasso is as follows:

Algorithm 1 (Estimation of Lasso loadings using Post-Lasso iterations). *Set $\hat{l}_{j,0} := \hat{l}_{jI_0}$, $j = 1, \dots, p$. Set $k = 0$, and specify a small constant $\nu \geq 0$ as a tolerance level and a constant $K > 1$ as an upper bound on the number of iterations. (1) Compute the Post-Lasso estimator $\tilde{\beta}$ based on the loadings $\hat{l}_{j,k}$. (2) For $\hat{s} = \|\tilde{\beta}\|_0 = |\hat{T}|$ set $\hat{l}_{j,k+1} := \sqrt{\mathbb{E}_n[x_{ij}^2(y_i - x_i'\tilde{\beta})^2]} \sqrt{n/(n - \hat{s})}$. (3) If $\max_{1 \leq j \leq p} |\hat{l}_{j,k} - \hat{l}_{j,k+1}| \leq \nu$ or $k > K$, set the loadings to $\hat{l}_{j,k+1}$, $j = 1, \dots, p$ and stop; otherwise, set $k \leftarrow k + 1$ and go to (1).*

APPENDIX B. AUXILIARY RESULTS ON MODEL SELECTION VIA LASSO AND POST-LASSO

The post-double-selection estimator applies the least squares estimator to the union of variables selected for equations (2.6) and (2.7) via feasible Lasso. Therefore, the model selection properties of feasible Lasso as well as properties of least squares estimates for m and g based on the selected model play an important role in the derivation of the main result. The purpose of this appendix is to describe these properties.

Note that each of the regression models (2.6)-(2.7) obeys the following conditions.

Condition ASM: Approximate Sparse Model. *Let $\{P_n\}$ be a sequence of data-generating processes. For each n , we have data $\{(\tilde{y}_i, \tilde{z}_i, \tilde{x}_i = P(\tilde{z}_i)) : 1 \leq i \leq n\}$ defined on $(\Omega, \mathcal{A}, P_n)$*

consisting of i.n.i.d vectors that obey the following approximately sparse regression model for each n :

$$\begin{aligned}\tilde{y}_i &= f(\tilde{z}_i) + \epsilon_i = \tilde{x}_i' \beta_0 + r_i + \epsilon_i, \\ \mathbb{E}[\epsilon_i \mid \tilde{x}_i] &= 0, \bar{\mathbb{E}}[\epsilon_i^2] = \sigma^2, \\ \|\beta_0\|_0 &\leq s, \bar{\mathbb{E}}[r_i^2] \lesssim \sigma^2 s/n.\end{aligned}$$

Let \hat{T} denote the model selected by the feasible Lasso estimator $\hat{\beta}$:

$$\hat{T} = \text{support}(\hat{\beta}) = \{j \in \{1, \dots, p\} : |\hat{\beta}_j| > 0\},$$

The Post-Lasso estimator $\tilde{\beta}$ is ordinary least squares applied to the data after removing the regressors that were not selected by the feasible Lasso:

$$\tilde{\beta} \in \arg \min_{\beta \in \mathbb{R}^p} \mathbb{E}_n[(\tilde{y}_i - \tilde{x}_i' \beta)^2] \quad : \quad \beta_j = 0 \text{ for each } j \notin \hat{T}. \quad (\text{B.54})$$

The following regularity conditions are imposed to deal with non-Gaussian, heteroscedastic errors.

Condition RF: Reduced Form. *In addition to ASTE, we have*

- (i) $\log^3 p/n \rightarrow 0$ and $s \log(p \vee n)/n \rightarrow 0$,
- (ii) $\bar{\mathbb{E}}[\tilde{y}_i^2] + \max_{1 \leq j \leq p} \{\bar{\mathbb{E}}[\tilde{x}_{ij}^2 \tilde{y}_i^2] + \bar{\mathbb{E}}[\tilde{x}_{ij}^3 \epsilon_i^3] + 1/\bar{\mathbb{E}}[\tilde{x}_{ij}^2 \epsilon_i^2]\} \lesssim 1$,
- (iii) $\max_{1 \leq j \leq p} \{ |(\mathbb{E}_n - \bar{\mathbb{E}})[\tilde{x}_{ij}^2 \epsilon_i^2]| + |(\mathbb{E}_n - \bar{\mathbb{E}})[\tilde{x}_{ij}^2 \tilde{y}_i^2]| \} + \max_{1 \leq i \leq n} \|\tilde{x}_i\|_\infty^2 \frac{s \log(n \vee p)}{n} = o_P(1)$.

The main auxiliary result that we use in proving the main result is as follows.

Lemma 1 (Model Selection Properties of Lasso and Properties of Post-Lasso). *Let $\{P_n\}$ be a sequence of data-generating processes. Suppose that conditions ASM and RF hold, and that Condition SE (P_n) holds for $\mathbb{E}_n[\tilde{x}_i \tilde{x}_i']$. Consider a feasible Lasso estimator with penalty level and loadings specified as in Section 3.3.*

(i) *Then the data-dependent model \hat{T} selected by a feasible Lasso estimator satisfies with probability approaching 1:*

$$\hat{s} = |\hat{T}| \lesssim s \quad (\text{B.55})$$

and

$$\min_{\beta \in \mathbb{R}^p: \beta_j=0 \ \forall j \notin \hat{T}} \sqrt{\mathbb{E}_n[f(\tilde{z}_i) - \tilde{x}_i' \beta]^2} \lesssim \sigma \sqrt{\frac{s \log(p \vee n)}{n}}. \quad (\text{B.56})$$

(ii) *The Post-Lasso estimator obeys*

$$\sqrt{\mathbb{E}_n[f(\tilde{z}_i) - \tilde{x}_i' \tilde{\beta}]^2} \lesssim_P \sigma \sqrt{\frac{s \log(p \vee n)}{n}}.$$

and

$$\|\tilde{\beta} - \beta_0\| \lesssim_P \sqrt{\mathbb{E}_n[\{\tilde{x}_i' \tilde{\beta} - \tilde{x}_i' \beta_0\}^2]} \lesssim_P \sigma \sqrt{\frac{s \log(p \vee n)}{n}}. \quad (\text{B.57})$$

Lemma 1 was derived in Belloni, Chen, Chernozhukov, and Hansen (2012) for Iterated Lasso and by Belloni, Chernozhukov, and Wang (2010) for Square-root Lasso. These analyses build on the rate analysis of infeasible Lasso by Bickel, Ritov, and Tsybakov (2009) and on sparsity analysis and rate analysis of Post-Lasso by Belloni and Chernozhukov (2013).

APPENDIX C. PROOF OF THEOREM 1

The proof proceeds under given sequence of probability measures $\{P_n\}$, as $n \rightarrow \infty$.

Let $Y = [y_1, \dots, y_n]'$, $X = [x_1, \dots, x_n]'$, $D = [d_1, \dots, d_n]'$, $V = [v_1, \dots, v_n]'$, $\zeta = [\zeta_1, \dots, \zeta_n]'$, $m = [m_1, \dots, m_n]'$, $R_m = [r_{m1}, \dots, r_{mn}]'$, $g = [g_1, \dots, g_n]'$, $R_g = [r_{g1}, \dots, r_{gn}]'$, and so on. For $A \subset \{1, \dots, p\}$, let $X[A] = \{X_j, j \in A\}$, where $\{X_j, j = 1, \dots, p\}$ are the columns of X . Let

$$\mathcal{P}_A = X[A](X[A]'X[A])^{-1}X[A]'$$

be the projection operator sending vectors in \mathbb{R}^n onto $\text{span}[X[A]]$, and let $\mathcal{M}_A = I_n - \mathcal{P}_A$ be the projection onto the subspace that is orthogonal to $\text{span}[X[A]]$. For a vector $W \in \mathbb{R}^n$, let

$$\tilde{\beta}_W(A) := \arg \min_{b \in \mathbb{R}^p} \|W - Xb\|^2 : b_j = 0, \forall j \notin A,$$

be the coefficient of linear projection of W onto $\text{span}[X[A]]$. If $A = \emptyset$, interpret $\mathcal{P}_A = 0_n$, and $\tilde{\beta}_W = 0_p$.

Finally, denote $\phi_{\min}(m) = \phi_{\min}(m)[\mathbb{E}_n[x_i x_i']]$ and $\phi_{\max}(m) = \phi_{\max}(m)[\mathbb{E}_n[x_i x_i']]$.

Step 1.(Main) Write $\tilde{\alpha} = [D' \mathcal{M}_{\hat{I}} D/n]^{-1} [D' \mathcal{M}_{\hat{I}} Y/n]$ so that

$$\sqrt{n}(\tilde{\alpha} - \alpha_0) = [D' \mathcal{M}_{\hat{I}} D/n]^{-1} [D' \mathcal{M}_{\hat{I}} (g + \zeta)/\sqrt{n}] =: ii^{-1} \cdot i.$$

By Steps 2 and 3,

$$ii = V'V/n + o_P(1) \text{ and } i = V'\zeta/\sqrt{n} + o_P(1).$$

Next note that $V'V/n = E[V'V/n] + o_P(1)$ by Chebyshev inequality, and because $E[V'V/n]$ is bounded away from zero and from above uniformly in n by Condition SM, we have $ii^{-1} = E[V'V/n]^{-1} + o_P(1)$.

By Condition SM we have $\sigma_n^2 = \bar{E}[v_i^2]^{-1} \bar{E}[\zeta_i^2 v_i^2] \bar{E}[v_i^2]^{-1}$ is bounded away from zero and from above, uniformly in n . Hence

$$Z_n = \sigma_n^{-1} \sqrt{n}(\tilde{\alpha} - \alpha_0) = n^{-1/2} \sum_{i=1}^n z_{i,n} + o_P(1),$$

where $z_{i,n} := \sigma_n^{-1} \bar{E}[v_i^2]^{-1} v_i \zeta_i$ are i.n.i.d. with mean zero. For $\delta > 0$ such that $4 + 2\delta \leq q$

$$\bar{E}|z_{i,n}|^{2+\delta} \lesssim \bar{E} \left[|v_i|^{2+\delta} |\zeta_i|^{2+\delta} \right] \lesssim \sqrt{\bar{E}|v_i|^{4+2\delta}} \sqrt{\bar{E}|\zeta_i|^{4+2\delta}} \lesssim 1,$$

by Condition SM. This condition verifies the Lyapunov condition and thus application of the Lyapunov CLT for i.n.i.d. triangular arrays implies that

$$Z_n \rightsquigarrow N(0, 1).$$

Step 2. (Behavior of i .) Decompose, using $D = m + V$,

$$i = V'\zeta/\sqrt{n} + \underbrace{m'\mathcal{M}_{\hat{I}}g/\sqrt{n}}_{=:i_a} + \underbrace{m'\mathcal{M}_{\hat{I}}\zeta/\sqrt{n}}_{=:i_b} + \underbrace{V'\mathcal{M}_{\hat{I}}g/\sqrt{n}}_{=:i_c} - \underbrace{V'\mathcal{P}_{\hat{I}}\zeta/\sqrt{n}}_{=:i_d}.$$

First, by Step 5 and 6 below we have

$$|i_a| = |m'\mathcal{M}_{\hat{I}}g/\sqrt{n}| \leq \sqrt{n} \|\mathcal{M}_{\hat{I}}g/\sqrt{n}\| \|\mathcal{M}_{\hat{I}}m/\sqrt{n}\| \lesssim_P \sqrt{[s \log(p \vee n)]^2/n} = o(1),$$

where the last bound follows from the assumed growth condition $s^2 \log^2(p \vee n) = o(n)$.

Second, using that $m = X\beta_{m0} + R_m$ and $m'\mathcal{M}_{\hat{I}}\zeta = R'_m\zeta - (\tilde{\beta}_m(\hat{I}) - \beta_{m0})'X'\zeta$, conclude

$$|i_b| \leq |R'_m\zeta/\sqrt{n}| + |(\tilde{\beta}_m(\hat{I}) - \beta_{m0})'X'\zeta/\sqrt{n}| \lesssim_P \sqrt{[s \log(p \vee n)]^2/n} = o_P(1).$$

This follows since

$$|R'_m\zeta/\sqrt{n}| \lesssim_P \sqrt{R'_m R_m/n} \lesssim_P \sqrt{s/n},$$

holding by Chebyshev inequality and Conditions SM and ASTE(iii), and

$$|(\tilde{\beta}_m(\hat{I}) - \beta_{m0})'X'\zeta/\sqrt{n}| \leq \|\tilde{\beta}_m(\hat{I}) - \beta_{m0}\|_1 \|X'\zeta/\sqrt{n}\|_\infty \lesssim_P \sqrt{[s^2 \log(p \vee n)]/n} \sqrt{\log(p \vee n)}.$$

The latter bound follows by (a)

$$\|\tilde{\beta}_m(\hat{I}) - \beta_{m0}\|_1 \leq \sqrt{\hat{s} + s} \|\tilde{\beta}_m(\hat{I}) - \beta_{m0}\| \lesssim_P \sqrt{[s^2 \log(p \vee n)]/n}$$

holding by Step 5 and by $\hat{s} \lesssim_P s$ implied by Lemma 1, and (b) by

$$\|X'\zeta/\sqrt{n}\|_\infty \lesssim_P \sqrt{\log(p \vee n)}$$

holding by Step 4 under Condition SM.

Third, using similar reasoning, decomposition $g = X\beta_{g0} + R_g$, and Steps 4 and 6, conclude

$$|i_c| \leq |R'_g V/\sqrt{n}| + |(\tilde{\beta}_g(\hat{I}) - \beta_{g0})'X'V/\sqrt{n}| \lesssim_P \sqrt{[s \log(p \vee n)]^2/n} = o_P(1).$$

Fourth, we have

$$|i_d| \leq |\tilde{\beta}_V(\hat{I})'X'\zeta/\sqrt{n}| \leq \|\tilde{\beta}_V(\hat{I})\|_1 \|X'\zeta/\sqrt{n}\|_\infty \lesssim_P \sqrt{[s \log(p \vee n)]^2/n} = o_P(1),$$

since by Step 4 below $\|X'\zeta/\sqrt{n}\|_\infty \lesssim_P \sqrt{\log(p \vee n)}$, and

$$\begin{aligned} \|\tilde{\beta}_V(\hat{I})\|_1 &\leq \sqrt{\hat{s}} \|\tilde{\beta}_V(\hat{I})\| \leq \sqrt{\hat{s}} \|(X[\hat{I}]'X[\hat{I}]/n)^{-1}X[\hat{I}]'V/n\| \\ &\leq \sqrt{\hat{s}}\phi_{\min}^{-1}(\hat{s})\sqrt{\hat{s}}\|X'V/\sqrt{n}\|_\infty/\sqrt{n} \lesssim_P s\sqrt{[\log(p \vee n)]/n}. \end{aligned}$$

The latter bound follows from $\hat{s} \lesssim_P s$, holding by Lemma 1, so that $\phi_{\min}^{-1}(\hat{s}) \lesssim_P 1$ by Condition SE, and from $\|X'V/\sqrt{n}\|_\infty \lesssim_P \sqrt{\log(p \vee n)}$ holding by Step 4.

Step 3. (Behavior of ii .) Decompose

$$ii = (m + V)' \mathcal{M}_{\hat{I}}(m + V)/n = V'V/n + \underbrace{m' \mathcal{M}_{\hat{I}} m/n}_{=: ii_a} + \underbrace{2m' \mathcal{M}_{\hat{I}} V/n}_{=: ii_b} - \underbrace{V' \mathcal{P}_{\hat{I}} V/n}_{=: ii_c}.$$

Then $|ii_a| \lesssim_P [s \log(p \vee n)]/n = o_P(1)$ by Step 5, $|ii_b| \lesssim_P [s \log(p \vee n)]/n = o_P(1)$ by reasoning similar to deriving the bound for $|i_b|$, and $|ii_c| \lesssim_P [s \log(p \vee n)]/n = o_P(1)$ by reasoning similar to deriving the bound for $|i_d|$.

Step 4. (Auxiliary: Bounds on $\|X'\zeta/\sqrt{n}\|_\infty$ and $\|X'V/\sqrt{n}\|_\infty$) Here we show that

$$(a) \|X'\zeta/\sqrt{n}\|_\infty \lesssim_P \sqrt{\log(p \vee n)} \quad \text{and} \quad (b) \|X'V/\sqrt{n}\|_\infty \lesssim_P \sqrt{\log(p \vee n)}.$$

To show (a), we use Lemma 5 stated in Appendix I on the tail bound for self-normalized deviations to deduce the bound. Indeed, we have that $\text{wp} \rightarrow 1$ for some $\ell_n \rightarrow \infty$ but so slowly that $1/\gamma = \ell_n \lesssim \log n$, with probability $1 - o(1)$

$$\max_{1 \leq j \leq p} \left| \frac{n^{-1/2} \sum_{i=1}^n x_{ij} \zeta_i}{\sqrt{\mathbb{E}_n[x_{ij}^2 \zeta_i^2]}} \right| \leq \Phi^{-1} \left(1 - \frac{1}{2\ell_n p} \right) \lesssim \sqrt{2 \log(2\ell_n p)} \lesssim \sqrt{\log(p \vee n)}. \quad (\text{C.58})$$

By Lemma 5 the first inequality in (C.58) holds, provided that for all n sufficiently large the following holds,

$$\Phi^{-1} \left(1 - \frac{1}{2\ell_n p} \right) \leq \frac{n^{1/6}}{\ell_n} \min_{1 \leq j \leq p} M_j^2 - 1, \quad M_j := \frac{\bar{\mathbb{E}}[x_{ij}^2 \zeta_i^2]^{1/2}}{\bar{\mathbb{E}}[|x_{ij}^3| |\zeta_i^3|]^{1/3}}.$$

Since we can choose ℓ_n to grow as slowly as needed, a sufficient condition for this are the conditions:

$$\log p = o(n^{1/3}) \quad \text{and} \quad \min_{1 \leq j \leq p} M_j \gtrsim 1,$$

which both hold by Condition SM. Finally,

$$\max_{1 \leq j \leq p} \mathbb{E}_n[x_{ij}^2 \zeta_i^2] \lesssim_P 1, \quad (\text{C.59})$$

by Condition SM. Therefore (a) follows from the bounds (C.58) and (C.59). Claim (b) follows similarly.

Step 5. (Auxiliary: Bound on $\|\mathcal{M}_{\hat{I}} m\|$ and related quantities.) This step shows that

$$(a) \|\mathcal{M}_{\hat{I}} m/\sqrt{n}\| \lesssim_P \sqrt{[s \log(p \vee n)]/n} \quad \text{and} \quad (b) \|\tilde{\beta}_m(\hat{I}) - \beta_{m0}\| \lesssim_P \sqrt{[s \log(p \vee n)]/n}.$$

Observe that

$$\sqrt{[s \log(p \vee n)]/n} \underset{(1)}{\gtrsim_P} \|\mathcal{M}_{\hat{I}_1} m/\sqrt{n}\| \underset{(2)}{\gtrsim_P} \|\mathcal{M}_{\hat{I}} m/\sqrt{n}\|$$

where inequality (1) holds since by Lemma 1 $\|\mathcal{M}_{\hat{I}_1} m/\sqrt{n}\| \leq \|(X\tilde{\beta}_D(\hat{I}_1) - m)/\sqrt{n}\| \lesssim_P \sqrt{[s \log(p \vee n)]/n}$, and (2) holds by $\hat{I}_1 \subseteq \hat{I}$ by construction. This shows claim (a). To show claim (b) note that

$$\|\mathcal{M}_{\hat{I}} m/\sqrt{n}\| \underset{(3)}{\geq} \|X(\tilde{\beta}_m(\hat{I}) - \beta_{m0})/\sqrt{n}\| - \|R_m/\sqrt{n}\|$$

where (3) holds by the triangle inequality. Since $\|R_m/\sqrt{n}\| \lesssim_P \sqrt{s/n}$ by Chebyshev and Condition ASTE(iii), conclude that

$$\begin{aligned} \sqrt{[s \log(p \vee n)]/n} &\gtrsim_P \|X(\tilde{\beta}_m(\hat{I}) - \beta_{m0})/\sqrt{n}\| \\ &\geq \sqrt{\phi_{\min}(\hat{s} + s)} \|\tilde{\beta}_m(\hat{I}) - \beta_{m0}\| \gtrsim_P \|\tilde{\beta}_m(\hat{I}) - \beta_{m0}\|, \end{aligned}$$

since $\hat{s} \lesssim_P s$ by Lemma 1 so that $1/\phi_{\min}(\hat{s} + s) \lesssim_P 1$ by condition SE. This shows claim (b).

Step 6. (Auxiliary: Bound on $\|\mathcal{M}_{\hat{I}}g\|$ and related quantities.) This step shows that

$$(a) \|\mathcal{M}_{\hat{I}}g/\sqrt{n}\| \lesssim_P \sqrt{[s \log(p \vee n)]/n} \text{ and } (b) \|\tilde{\beta}_g(\hat{I}) - \beta_{g0}\| \lesssim_P \sqrt{[s \log(p \vee n)]/n}.$$

Observe that

$$\begin{aligned} \sqrt{[s \log(p \vee n)]/n} &\stackrel{(1)}{\gtrsim_P} \|\mathcal{M}_{\hat{I}_2}(\alpha_0 m + g)/\sqrt{n}\| \\ &\stackrel{(2)}{\gtrsim_P} \|\mathcal{M}_{\hat{I}}(\alpha_0 m + g)/\sqrt{n}\| \\ &\stackrel{(3)}{\gtrsim_P} \|\mathcal{M}_{\hat{I}}g/\sqrt{n}\| - \|\mathcal{M}_{\hat{I}}\alpha_0 m/\sqrt{n}\| \end{aligned}$$

where inequality (1) holds since by Lemma 1 $\|\mathcal{M}_{\hat{I}_2}(\alpha_0 m + g)/\sqrt{n}\| \leq \|(X\tilde{\beta}_Y(\hat{I}_2) - \alpha_0 m - g)/\sqrt{n}\| \lesssim_P \sqrt{[s \log(p \vee n)]/n}$, (2) holds by $\hat{I}_2 \subseteq \hat{I}$, and (3) by the triangle inequality. Since $|\alpha_0|$ is bounded uniformly in n by assumption, by Step 5, $\|\mathcal{M}_{\hat{I}}\alpha_0 m/\sqrt{n}\| \lesssim_P \sqrt{[s \log(p \vee n)]/n}$. Hence claim (a) follows by the triangle inequality:

$$\sqrt{[s \log(p \vee n)]/n} \gtrsim_P \|\mathcal{M}_{\hat{I}}g/\sqrt{n}\|$$

To show claim (b) we note that

$$\|\mathcal{M}_{\hat{I}}g/\sqrt{n}\| \geq \|X(\tilde{\beta}_g(\hat{I}) - \beta_{g0})/\sqrt{n}\| - \|R_g/\sqrt{n}\|$$

where $\|R_g/\sqrt{n}\| \lesssim_P \sqrt{s/n}$ by Condition ASTE(iii). Then conclude similarly to Step 5 that

$$\begin{aligned} \sqrt{[s \log(p \vee n)]/n} &\gtrsim_P \|X(\tilde{\beta}_g(\hat{I}) - \beta_{g0})/\sqrt{n}\| \\ &\geq \sqrt{\phi_{\min}(\hat{s} + s)} \|\tilde{\beta}_g(\hat{I}) - \beta_{g0}\| \gtrsim_P \|\tilde{\beta}_g(\hat{I}) - \beta_{g0}\|. \end{aligned}$$

Step 7. (Variance Estimation.) Since $\hat{s} \lesssim_P s = o(n)$, $(n - \hat{s} - 1)/n = o_P(1)$, and since $\bar{\mathbb{E}}[v_i^2 \zeta_i^2]$ and $\bar{\mathbb{E}}[v_i^2]$ are bounded away from zero and from above uniformly in n by Condition SM, it suffices to show that

$$\mathbb{E}_n[\hat{v}_i^2 \hat{\zeta}_i^2] - \bar{\mathbb{E}}[v_i^2 \zeta_i^2] \rightarrow_P 0, \quad \mathbb{E}_n[\hat{v}_i^2] - \bar{\mathbb{E}}[v_i^2] \rightarrow_P 0,$$

The second relation was shown in Step 3, so it remains to show the first relation.

Let $\tilde{v}_i = v_i + r_{mi}$ and $\tilde{\zeta}_i = \zeta_i + r_{gi}$. Recall that by Condition ASTE(v) we have $\bar{\mathbb{E}}[\tilde{v}_i^2 \tilde{\zeta}_i^2] - \bar{\mathbb{E}}[v_i^2 \zeta_i^2] \rightarrow 0$, and $\mathbb{E}_n[\tilde{v}_i^2 \tilde{\zeta}_i^2] - \bar{\mathbb{E}}[\tilde{v}_i^2 \tilde{\zeta}_i^2] \rightarrow_P 0$ by Vonbahr-Esseen's inequality in von Bahr and Esseen (1965) since $\bar{\mathbb{E}}[|\tilde{v}_i \tilde{\zeta}_i|^{2+\delta}] \leq (\bar{\mathbb{E}}[|\tilde{v}_i|^{4+2\delta}] \bar{\mathbb{E}}[|\tilde{\zeta}_i|^{4+2\delta}])^{1/2}$ is uniformly bounded for $4 + 2\delta \leq q$. Thus it suffices to show that $\mathbb{E}_n[\hat{v}_i^2 \hat{\zeta}_i^2] - \mathbb{E}_n[\tilde{v}_i^2 \tilde{\zeta}_i^2] \rightarrow_P 0$.

By the triangular inequality

$$|\mathbb{E}_n[\widehat{v}_i^2 \widehat{\zeta}_i^2 - \widetilde{v}_i^2 \widetilde{\zeta}_i^2]| \leq |\mathbb{E}_n[(\widehat{v}_i^2 - \widetilde{v}_i^2) \widetilde{\zeta}_i^2]| + |\mathbb{E}_n[\widehat{v}_i^2 (\widehat{\zeta}_i^2 - \widetilde{\zeta}_i^2)]|.$$

$\quad \quad \quad =: iv \quad \quad \quad =: iii$

Then, expanding $\widehat{\zeta}_i^2 - \widetilde{\zeta}_i^2$ we have

$$\begin{aligned} iii &\leq 2\mathbb{E}_n[\{d_i(\alpha_0 - \check{\alpha})\}^2 \widehat{v}_i^2] + 2\mathbb{E}_n[\{x'_i(\check{\beta} - \beta_{g0})\}^2 \widehat{v}_i^2] \\ &\quad + |2\mathbb{E}_n[\widetilde{\zeta}_i d_i(\alpha_0 - \check{\alpha}) \widehat{v}_i^2]| + |2\mathbb{E}_n[\widetilde{\zeta}_i x'_i(\check{\beta} - \beta_{g0}) \widehat{v}_i^2]| \\ &=: iii_a + iii_b + iii_c + iii_d = o_P(1) \end{aligned}$$

where the last bound follows by the relations derived below.

First, we note

$$iii_a \leq 2 \max_{i \leq n} d_i^2 |\alpha_0 - \check{\alpha}|^2 \mathbb{E}_n[\widehat{v}_i^2] \lesssim_P n^{(2/q)-1} = o(1) \quad (\text{C.60})$$

$$iii_c \leq 2 \max_{i \leq n} \{|\widetilde{\zeta}_i| |d_i|\} \mathbb{E}_n[\widehat{v}_i^2] |\alpha_0 - \check{\alpha}| \lesssim_P n^{(2/q)-(1/2)} = o(1) \quad (\text{C.61})$$

which holds by the following argument. Condition SM assumes that $\mathbb{E}[|d_i|^q]$ which in turn implies that $\mathbb{E}[\max_{i \leq n} d_i^2] \lesssim n^{2/q}$. Similarly Condition ASTE implies that $\mathbb{E}[\max_{i \leq n} \widetilde{\zeta}_i^2] \lesssim n^{2/q}$ and $\mathbb{E}[\max_{i \leq n} \widetilde{v}_i^2] \lesssim n^{2/q}$. Thus by Markov inequality

$$\max_{i \leq n} |d_i| + |\widetilde{\zeta}_i| + |\widetilde{v}_i| \lesssim_P n^{1/q}. \quad (\text{C.62})$$

Moreover, $\mathbb{E}_n[\widehat{v}_i^2] \lesssim_P 1$ and $|\check{\alpha} - \alpha_0| \lesssim_P n^{-1/2}$ by the previous steps. These bounds and $q > 4$ imposed in Condition SM imply (C.60)-(C.61).

Next we bound,

$$\begin{aligned} iii_d &\leq 2 \max_{i \leq n} |\widetilde{\zeta}_i| \max_{i \leq n} |x'_i(\check{\beta} - \beta_{g0})| \mathbb{E}_n[\widehat{v}_i^2] \\ &\lesssim_P n^{1/q} \max_{i \leq n} \|x_i\|_\infty \sqrt{\frac{s}{\sqrt{n}} \frac{s \log(p \vee n)}{\sqrt{n}}} = o_P(1), \end{aligned} \quad (\text{C.63})$$

using (C.62) and that for $\widehat{T}_g = \text{support}(\beta_{g0}) \cup \widehat{I}$, we have

$$\max_{i \leq n} \{x'_i(\check{\beta} - \beta_{g0})\}^2 \leq \max_{i \leq n} \|x_{i\widehat{T}_g}\|^2 \|\check{\beta} - \beta_{g0}\|^2,$$

where

$$\max_{i \leq n} \|x_{i\widehat{T}_g}\|^2 \leq |\widehat{T}_g| \max_{i \leq n} \|x_i\|_\infty^2 \lesssim_P s \max_{i \leq n} \|x_i\|_\infty^2$$

by the sparsity assumption in ASTE and the sparsity bound in Lemma 1, and since $\check{\beta}[\widehat{I}] = (X[\widehat{I}]' X[\widehat{I}])^{-1} X[\widehat{I}]' (\zeta + g - (\check{\alpha} - \alpha_0)D)$ we have

$$\|\check{\beta} - \beta_{g0}\| \leq \|\check{\beta}_g(\widehat{I}) - \beta_{g0}\| + \|\check{\beta}_\zeta(\widehat{I})\| + |\check{\alpha} - \alpha_0| \cdot \|\check{\beta}_D(\widehat{I})\| \lesssim_P \sqrt{s \log(p \vee n)/n}$$

by Step 6(b), by

$$\|\check{\beta}_\zeta(\widehat{I})\| \leq \sqrt{\widehat{s}} \phi_{\min}^{-1}(\widehat{s}) \|X' \zeta / n\|_\infty \lesssim_P \sqrt{s \log(p \vee n)/n}$$

holding by Condition SE and by $\hat{s} \lesssim_P s$ from Lemma 1, and by Step 4, $|\check{\alpha} - \alpha_0| \lesssim_P 1/\sqrt{n}$ by Step 1, and

$$\|\tilde{\beta}_D(\hat{I})\| \leq \phi_{\min}^{-1}(\hat{s})\sqrt{\hat{s}} \max_{1 \leq j \leq p} |\mathbb{E}_n[x_{ij}d_i]| \leq \phi_{\min}^{-1}(\hat{s})\sqrt{\hat{s}} \max_{1 \leq j \leq p} \sqrt{\mathbb{E}_n[x_{ij}^2 d_i^2]} \lesssim_P \sqrt{s}$$

by Condition SE, $\hat{s} \lesssim_P s$ by the sparsity bound in Lemma 1, and Condition SM.

The final conclusion in (C.63) then follows by condition ASTE (iv) and (v).

Next, using the relations above and condition ASTE (iv) and (v), we also conclude that

$$\begin{aligned} iii_b &\leq 2 \max_{i \leq n} \{x'_i(\check{\beta} - \beta_{g0})\}^2 \mathbb{E}_n[\hat{v}_i^2] \\ &\lesssim_P \max_{i \leq n} \|x_i\|_\infty^2 \frac{s}{\sqrt{n}} \frac{s \log(p \vee n)}{\sqrt{n}} = o_P(1). \end{aligned} \quad (\text{C.64})$$

Finally, the argument for $iv = o_P(1)$ follows similarly to the argument for $iii = o_P(1)$ and the result follows. \square

APPENDIX D. PROOF OF COROLLARY 1

Let \mathbf{P}_n be a collection of probability measures \mathbf{P} for which conditions ASTE (P), SM (P), SE (P), and R (P) hold for the given n . Consider any sequence $\{\mathbf{P}_n\}$, with index $n \in \{n_0, n_0 + 1, \dots\}$, with $\mathbf{P}_n \in \mathbf{P}_n$ for each $n \in \{n_0, n_0 + 1, \dots\}$. By Theorem 1 we have that, for $c = \Phi^{-1}(1 - \gamma/2)$, $\lim_{n \rightarrow \infty} \mathbf{P}_n(\alpha_0 \in [\check{\alpha} \pm c\hat{\sigma}_n/\sqrt{n}]) = \Phi(c) - \Phi(-c) = 1 - \gamma$. This means that for every further subsequence $\{\mathbf{P}_{n_k}\}$ with $\mathbf{P}_{n_k} \in \mathbf{P}_{n_k}$ for each $k \in \{1, 2, \dots\}$

$$\lim_{k \rightarrow \infty} \mathbf{P}_{n_k}(\alpha_0 \in [\check{\alpha} \pm c\hat{\sigma}_{n_k}/\sqrt{n_k}]) = 1 - \gamma. \quad (\text{D.65})$$

Suppose that the claim of corollary does not hold, i.e.

$$\limsup_{n \rightarrow \infty} \sup_{\mathbf{P} \in \mathbf{P}_n} \left| \mathbf{P}(\alpha_0 \in [\check{\alpha} \pm c\hat{\sigma}_n/\sqrt{n}]) - (1 - \gamma) \right| > 0.$$

Hence there is a subsequence $\{\mathbf{P}_{n_k}\}$ with $\mathbf{P}_{n_k} \in \mathbf{P}_{n_k}$ for each $k \in \{1, 2, \dots\}$ such that:

$$\lim_{k \rightarrow \infty} \mathbf{P}_{n_k}(\alpha_0 \in [\check{\alpha} \pm c\hat{\sigma}_{n_k}/\sqrt{n_k}]) \neq 1 - \gamma.$$

This gives a contradiction to (D.65). The claim (i) follows. Claim (ii) follows from claim (i), since $\mathbf{P} \subseteq \mathbf{P}_n$ for all $n \geq n_0$. \square

APPENDIX E. PROOF OF THEOREM 2

We use the same notation as in Theorem 1. Using that notation the approximations bounds stated in Condition HLMS are equivalent to $\|\mathcal{M}_{\hat{I}}g\| \leq \delta_n n^{1/4}$ and $\|\mathcal{M}_{\hat{I}}m\| \leq \delta_n n^{1/4}$.

Step 1. It follows the same reasoning as Step 1 in the proof of Theorem 1.

Step 2. (Behavior of i .) Decompose, using $D = m + V$

$$i = V'\zeta/\sqrt{n} + \underbrace{m'\mathcal{M}_{\hat{I}}g/\sqrt{n}}_{=:i_a} + \underbrace{m'\mathcal{M}_{\hat{I}}\zeta/\sqrt{n}}_{=:i_b} + \underbrace{V'\mathcal{M}_{\hat{I}}g/\sqrt{n}}_{=:i_c} - \underbrace{V'\mathcal{P}_{\hat{I}}\zeta/\sqrt{n}}_{=:i_d}.$$

First, by Condition HLMS we have $\|\mathcal{M}_{\hat{I}}g\| = o_P(n^{1/4})$ and $\|\mathcal{M}_{\hat{I}}m\| = o_P(n^{1/4})$. Therefore

$$|i_a| = |m'\mathcal{M}_{\hat{I}}g/\sqrt{n}| \leq \sqrt{n}\|\mathcal{M}_{\hat{I}}g/\sqrt{n}\|\|\mathcal{M}_{\hat{I}}m/\sqrt{n}\| \lesssim_P o(1).$$

Second, using that $m = X\beta_{m0} + R_m$ and $m'\mathcal{M}_{\hat{I}}\zeta = R'_m\zeta - (\tilde{\beta}_m(\hat{I}) - \beta_{m0})'X'\zeta$, we have

$$\begin{aligned} |i_b| &\leq |R'_m\zeta/\sqrt{n}| + |(\tilde{\beta}_m(\hat{I}) - \beta_{m0})'X'\zeta/\sqrt{n}| \\ &\leq |R'_m\zeta/\sqrt{n}| + \|\tilde{\beta}_m(\hat{I}) - \beta_{m0}\|_1 \|X'\zeta/\sqrt{n}\|_\infty \\ &\lesssim_P \sqrt{s/n} + \sqrt{s} \{o(n^{-1/4}) + \sqrt{s/n}\} \sqrt{\log(p \vee n)} = o(1). \end{aligned}$$

This follows because

$$|R'_m\zeta/\sqrt{n}| \lesssim_P \sqrt{R'_m R_m/n} \lesssim_P \sqrt{s/n},$$

by Chebyshev inequality and Conditions SM and ASTE(iii),

$$\|\tilde{\beta}_m(\hat{I}) - \beta_{m0}\|_1 \leq \sqrt{\hat{s} + s} \|\tilde{\beta}_m(\hat{I}) - \beta_{m0}\| \lesssim_P \sqrt{s} \{o(n^{-1/4}) + \sqrt{s/n}\},$$

by Step 4 and $\hat{s} = |\hat{I}| \lesssim_P s$ by Condition HLMS, and

$$\|X'\zeta/\sqrt{n}\|_\infty \lesssim_P \sqrt{\log(p \vee n)}$$

holding by Step 4 in the proof of Theorem 1.

Third, using similar reasoning and the decomposition $g = X\beta_{g0} + R_g$ conclude

$$\begin{aligned} |i_c| &\leq |R'_g V/\sqrt{n}| + |(\tilde{\beta}_g(\hat{I}) - \beta_{g0})'X'V/\sqrt{n}| \\ &\lesssim_P \sqrt{s/n} + \sqrt{s} \{o(n^{-1/4}) + \sqrt{s/n}\} \sqrt{\log(p \vee n)} = o_P(1). \end{aligned}$$

Fourth, we have

$$|i_d| \leq |\tilde{\beta}_V(\hat{I})'X'\zeta/\sqrt{n}| \leq \|\tilde{\beta}_V(\hat{I})\|_1 \|X'\zeta/\sqrt{n}\|_\infty \lesssim_P \sqrt{[s \log(p \vee n)]^2/n} = o_P(1),$$

since $\|X'\zeta/\sqrt{n}\|_\infty \lesssim_P \sqrt{\log(p \vee n)}$ by Step 4 of the proof of Theorem 1, and

$$\begin{aligned} \|\tilde{\beta}_V(\hat{I})\|_1 &\leq \sqrt{\hat{s}} \|\tilde{\beta}_V(\hat{I})\| \leq \sqrt{\hat{s}} \|(X[\hat{I}]'X[\hat{I}]/n)^{-1}X[\hat{I}]'V/n\| \\ &\leq \sqrt{\hat{s}}\phi_{\min}^{-1}(\hat{s})\sqrt{\hat{s}}\|X'V/\sqrt{n}\|_\infty/\sqrt{n} \lesssim_P s\sqrt{[\log(p \vee n)]/n}. \end{aligned}$$

The latter bound follows from $\hat{s} \lesssim_P s$ by condition HLMS so that $\phi_{\min}^{-1}(\hat{s}) \lesssim_P 1$ by condition SE, and again invoking Step 4 of the proof of Theorem 1 to establish $\|X'V/\sqrt{n}\|_\infty \lesssim_P \sqrt{\log(p \vee n)}$.

Step 3. (Behavior of ii .) Decompose

$$ii = (m + V)'\mathcal{M}_{\hat{I}}(m + V)/n = V'V/n + \underbrace{m'\mathcal{M}_{\hat{I}}m/n}_{=:ii_a} + \underbrace{2m'\mathcal{M}_{\hat{I}}V/n}_{=:ii_b} - \underbrace{V'\mathcal{P}_{\hat{I}}V/n}_{=:ii_c}.$$

Then $|ii_a| \lesssim_P o(n^{1/2})/n = o_P(n^{-1/2})$ by condition HLMS, $|ii_b| = o(n^{-1/2})$ by reasoning similar to deriving the bound for $|i_b|$, and $|ii_c| \lesssim_P [s \log(p \vee n)]/n = o_P(1)$ by reasoning similar to deriving the bound for $|i_d|$.

Step 4. (Auxiliary: Bounds on $\|\tilde{\beta}_m(\hat{I}) - \beta_{m0}\|$ and $\|\tilde{\beta}_g(\hat{I}) - \beta_{g0}\|$.) To establish a bound on $\|\tilde{\beta}_g(\hat{I}) - \beta_{g0}\|$ note that

$$\|\mathcal{M}_{\hat{I}}g/\sqrt{n}\| \geq \|X(\tilde{\beta}_g(\hat{I}) - \beta_{g0})/\sqrt{n}\| - \|R_g/\sqrt{n}\|$$

where $\|R_g/\sqrt{n}\| \lesssim_P \sqrt{s/n}$ holds by Chebyshev inequality and Condition ASTE(iii). Moreover, by Condition HLMS we have $\|\mathcal{M}_{\hat{I}}g/\sqrt{n}\| = o_P(n^{-1/4})$ and $\hat{s} = |\hat{I}| \lesssim_P s$. Thus

$$\begin{aligned} o(n^{-1/4}) + \sqrt{s/n} &\gtrsim_P \|X(\tilde{\beta}_g(\hat{I}) - \beta_{g0})/\sqrt{n}\| \\ &\geq \sqrt{\phi_{\min}(s + \hat{s})} \|\tilde{\beta}_g(\hat{I}) - \beta_{g0}\| \\ &\gtrsim_P \|\tilde{\beta}_g(\hat{I}) - \beta_{g0}\| \end{aligned}$$

since $\sqrt{\phi_{\min}(s + \hat{s})} \gtrsim_P 1$ by Condition SE.

The same logic yields $\|\tilde{\beta}_m(\hat{I}) - \beta_{m0}\| \lesssim_P \sqrt{s/n} + o(n^{-1/4})$.

Step 5. (Variance Estimation.) It follows similarly to Step 7 in the proof of Theorem 1 but using Condition HLMS instead of Lemma 1. □

APPENDIX F. PROOF OF COROLLARY 2

The proof is similar to the proof of Corollary 1.

APPENDIX G. VERIFICATION OF CONDITIONS FOR THE EXAMPLES

G.1. Proof of Corollary 3 (Verification for Example 1). Let \mathbf{P} be the collection of all regression models P that obey the conditions set forth above for all n for the given constants (p, b, B, q_x, q) . Below we provide explicit bounds for κ' , κ'' , c , C , δ_n and Δ_n that appear in Conditions ASTE, SE and SM that depend only on (p, b, B, q_x, q) and n which in turn establish these conditions for any $P \in \mathbf{P}$.

Condition ASTE(i) is assumed. Condition ASTE(ii) holds with $|\alpha_0| \leq C_1^{ASTE} = B$. Condition ASTE(iii) holds with $s = p$ and $r_{gi} = r_{mi} = 0$.

Condition ASTE(iv) holds with $\delta_{1n}^{ASTE} := p^2 \log^2(p \vee n)/n \rightarrow 0$ since $s = p$ is fixed. Finally, we verify ASTE(v). Because $\tilde{v}_i = v_i$, $\tilde{\zeta}_i = \zeta_i$ and the moment condition $E[|v_i^q|] + E[|\zeta_i^q|] \leq C_2^{ASTE} = 2B$ with $q > 4$, the first two requirements follow. To show the last requirement, note that because

$E[\|x_i\|^{q_x}] \leq B$ we have

$$P\left(\max_{1 \leq i \leq n} \|x_i\|_\infty > t_{1n}\right) \leq P\left(\left[\sum_{i=1}^n \|x_i\|^{q_x}\right]^{1/q_x} > t_{1n}\right) \leq nE[\|x_i\|^{q_x}]/t_{1n}^{q_x} \leq nB/t_{1n}^{q_x} =: \Delta_{1n}^{ASTE}. \quad (\text{G.66})$$

Let $t_{1n} = (n \log n)^{1/q_x} B^{1/q_x}$ so that $\Delta_{1n}^{ASTE} = 1/\log n$. Thus we have with probability $1 - \Delta_{1n}^{ASTE}$

$$\max_{1 \leq i \leq n} \|x_i\|_\infty^2 sn^{-1/2+2/q} \leq (n \log n)^{2/q_x} B^{2/q_x} pn^{-1/2+2/q} =: \delta_{2n}^{ASTE}.$$

It follows that $\delta_{2n}^{ASTE} \rightarrow 0$ by the assumption that $4/q_x + 4/q < 1$.

To verify Condition SE note that

$$\begin{aligned} P(\|\mathbb{E}_n[x_i x'_i] - \mathbb{E}[x_i x'_i]\| > t_{2n}) &\leq \sum_{k=1}^p \sum_{j=1}^p \frac{E[x_{ij}^2 x_{ik}^2]}{nt_{2n}^2} \leq \sum_{k=1}^p \sum_{j=1}^p \frac{E[x_{ij}^4] + E[x_{ik}^4]}{2nt_{2n}^2} \\ &\leq \frac{pE[\|x_i\|^4]}{nt_{2n}^2} \leq \frac{pB^{4/q_x}}{nt_{2n}^2} =: \Delta_{1n}^{SE}. \end{aligned}$$

Setting $t_{2n} := b/2$ we have $\Delta_{1n}^{SE} = (2/b)^2 B^{4/q_x} p/n \rightarrow 0$ since p is fixed. Then, with probability $1 - \Delta_{1n}^{SE}$ we have

$$\begin{aligned} \lambda_{\min}(\mathbb{E}_n[x_i x'_i]) &\geq \lambda_{\min}(\mathbb{E}[x_i x'_i]) - \|\mathbb{E}_n[x_i x'_i] - \mathbb{E}[x_i x'_i]\| \geq b/2 =: \kappa', \\ \lambda_{\max}(\mathbb{E}_n[x_i x'_i]) &\leq \lambda_{\max}(\mathbb{E}[x_i x'_i]) + \|\mathbb{E}_n[x_i x'_i] - \mathbb{E}[x_i x'_i]\| \leq E[\|x_i\|^2] + b/2 \leq 2B^{2/q_x} =: \kappa''. \end{aligned}$$

In the verification of Condition SM note that the second and third requirements in Condition SM(i) hold with $c_1^{SM} = b$ and $C_1^{SM} = B^{2/q}$. Condition SM(iii) holds with $\delta_{1n}^{SM} := \log^3 p/n \rightarrow 0$ since p is fixed.

The first requirement in Condition SM(i) and Condition SM(ii) hold by the stated moment assumptions, for $\epsilon_i = v_i$ and $\epsilon_i = \zeta_i$, $\tilde{y}_i = d_i$ and $\tilde{y}_i = y_i$,

$$\begin{aligned} E[\|\epsilon_i^q\|] &\leq B =: A_1 \\ E[\|d_i^q\|] &\leq 2^{q-1}E[|x'_i \beta_{m0}|^q] + 2^{q-1}E[|v_i^q|] \leq 2^{q-1}E[\|x_i\|^q] \|\beta_{m0}\|^q + 2^{q-1}E[|v_i^q|] \\ &\leq 2^{q-1}(B^{q/q_x} B^q + B) =: A_2 \\ E[d_i^4] &\leq 2^3(B^{4/q_x} B^4 + B) =: A'_2 \\ E[y_i^4] &\leq 3^3 \|\alpha_0\|^4 E[d_i^4] + 3^3 \|\beta_{g0}\|^4 E[\|x_i\|^4] + 3^3 E[\zeta_i^4] \\ &\leq 3^3 B^4 2^3 A'_2 + 3^3 B^4 B^{4/q_x} + 3^3 B^{4/q} =: A_3 \\ \max_{1 \leq j \leq p} E[x_{ij}^2 \tilde{y}_i^2] &\leq \max_{1 \leq j \leq p} (E[x_{ij}^4])^{1/2} (E[\tilde{y}_i^4])^{1/2} \leq B^{2/q_x} (E[\tilde{y}_i^4])^{1/2} \leq B^{2/q_x} (A'_2 \vee A_3)^{1/2} =: A_4 \\ \max_{1 \leq j \leq p} E[|x_{ij} \epsilon_i|^3] &= \max_{1 \leq j \leq p} E[|x_{ij}^3| E[|\epsilon_i^3| | x_i]] \leq B^{3/q} \max_{1 \leq j \leq p} E[|x_{ij}^3|] \leq B^{3/q+3/q_x} =: A_5 \\ \max_{1 \leq j \leq p} 1/E[x_{ij}^2] &\leq 1/\lambda_{\min}(E[x_i x'_i]) \leq 1/b =: A_6 \end{aligned}$$

since $4 < q \leq q_x$. Thus these conditions hold with $C_2^{SM} = A_2 \vee (A_1 + (A'_2 \vee A_3)^{1/2} + A_4 + A_5 + A_6)$.

Next we show Condition SM(iv). By (G.66) we have $\max_{1 \leq i \leq n} \|x_i\|_\infty^2 \leq (n \log n)^{2/q_x} B^{2/q_x}$ with probability $1 - \Delta_{1n}^{ASTE}$, thus with the same probability

$$\max_{i \leq n} \|x_i\|_\infty^2 \frac{s \log(n \vee p)}{n} \leq (B \log n)^{2/q_x} \frac{n^{2/q_x} p \log(p \vee n)}{n} =: \delta_{1n}^{SM} \rightarrow 0$$

since $q_x > 4$ and $s = p$ is fixed.

Next for $\epsilon_i = v_i$ and $\epsilon_i = \zeta_i$ we have

$$\mathbb{P} \left(\max_{1 \leq j \leq p} |(\mathbb{E}_n - \mathbb{E})[x_{ij}^2 \epsilon_i^2]| > \delta_{2n}^{SM} \right) \leq \sum_{j=1}^p \frac{\mathbb{E}[x_{ij}^4 \epsilon_i^4]}{n(\delta_{2n}^{SM})^2} \leq \frac{p B^{4/q_x + 4/q_x}}{n(\delta_{2n}^{SM})^2} =: \Delta_{1n}^{SM}$$

by the union bound, Chebyshev inequality and by $\mathbb{E}[x_{ij}^4 \epsilon_i^4] = \mathbb{E}[x_{ij}^4 \mathbb{E}[\epsilon_i^4 | x_i]] \leq B^{4/q_x + 4/q_x}$. Letting $\delta_{2n}^{SM} = B^{2/q_x + 2/q_x} n^{-1/4} \rightarrow 0$ we have $\Delta_{1n}^{SM} = p/n^{1/2} \rightarrow 0$ since p, B, q and q_x are fixed.

Next for $\tilde{y}_i = d_i$ and $\tilde{y}_i = y_i$ we have

$$\mathbb{P} \left(\max_{1 \leq j \leq p} |(\mathbb{E}_n - \mathbb{E})[x_{ij}^2 \tilde{y}_i^2]| > \delta_{3n}^{SM} \right) \leq \sum_{j=1}^p \frac{\mathbb{E}[x_{ij}^4 \tilde{y}_i^4]}{n(\delta_{3n}^{SM})^2} \leq \frac{p B^{4/q_x} A_8^{4/q}}{n(\delta_{3n}^{SM})^2} =: \Delta_{2n}^{SM}$$

by the union bound, Chebyshev inequality and by

$$\mathbb{E}[x_{ij}^4 \tilde{y}_i^4] \leq \mathbb{E}[x_{ij}^{\tilde{q}}]^{4/\tilde{q}} \mathbb{E}[\tilde{y}_i^q]^{4/q} \leq \mathbb{E}[x_{ij}^{q_x}]^{4/q_x} \mathbb{E}[\tilde{y}_i^q]^{4/q} \leq B^{4/q_x} A_8^{4/q}$$

holding by Hölder inequality where $4 < \tilde{q} \leq q_x$ such that $4/q + 4/\tilde{q} = 1$, and

$$\begin{aligned} \mathbb{E}[\tilde{y}_i^q] &\leq (1 + 3^{q-1} |\alpha_0|^q) \mathbb{E}[d_i^q] + 3^{q-1} \|\beta_{g0}\|^q \mathbb{E}[\|x_i\|^q] + 3^{q-1} \mathbb{E}[\zeta_i^q] \\ &\leq 3^q (A_2 + B^q A_2 + B^q B^{q/q_x} + B) =: A_8. \end{aligned}$$

Letting $\delta_{3n}^{SM} = B^{4/q_x} A_8^{4/q} n^{-1/4} \rightarrow 0$ we have $\Delta_{2n}^{SM} = p/n^{1/2} \rightarrow 0$ since p, B, q and q_x are fixed.

Finally, we set $c = c_1^{SM}$, $C = \max\{C_1^{ASTE}, C_2^{ASTE}, C_1^{SM}, C_2^{SM}\}$, $\delta_n = \max\{\delta_{1n}^{ASTE}, \delta_{2n}^{ASTE}, \delta_{1n}^{SM} + \delta_{2n}^{SM} + \delta_{3n}^{SM}\} \rightarrow 0$, and $\Delta_n = \max\{\Delta_{1n}^{ASTE} + \Delta_{1n}^{SM} + \Delta_{2n}^{SM}, \Delta_{1n}^{SE}\} \rightarrow 0$. \square

G.2. Proofs of Corollaries 4 and 3 (Examples 2 and 3). We will make use of the following technical lemma in the verification of examples 2 and 3.

Lemma 2 (Uniform Approximation). *Let $h_i = x_i' \theta_h + \rho_i$ be a function whose coefficients $\theta_h \in S_A^a(p)$, and $\underline{\kappa} \leq \lambda_{\min}(\mathbb{E}[x_i x_i']) \leq \lambda_{\max}(\mathbb{E}[x_i x_i']) \leq \bar{\kappa}$. For $s = A^{1/a} n^{1/2a}$, $a > 1$, define β_{h0} as in (4.31), $r_{hi} = h_i - x_i' \beta_{h0}$, for $i = 1, \dots, n$. Then we have*

$$|r_{hi}| \leq \|x_i\|_\infty (\bar{\kappa}/\underline{\kappa})^{3/2} \left\{ \frac{2a-1}{a-1} \sqrt{s^2/n} + 5 \sqrt{s \mathbb{E}[\rho_i^2]/\underline{\kappa}} \right\} + |\rho_i|.$$

Proof. Let T_h denote the support of β_{h0} and S denote the support of the s largest components of θ_h . Note that $|T_h| = |S| = s$. First we establish some auxiliary bounds on the $\|\theta_h[T_h^c]\|$ and $\|\theta_h[T_h^c]\|_1$. By the optimality of T_h and β_{h0} we have that

$$\sqrt{\mathbb{E}[(h_i - x_i' \beta_{h0})^2]} \leq \sqrt{\mathbb{E}[(x_i[S^c]' \theta_h[S^c] + \rho_i)^2]} \leq \sqrt{\bar{\kappa}} \|\theta_h[S^c]\| + \sqrt{\mathbb{E}[\rho_i^2]} \quad \text{and}$$

$$\sqrt{\mathbb{E}[(h_i - x_i' \beta_{h0})^2]} = \sqrt{\mathbb{E}[\{x_i'(\theta_h - \beta_{h0}) + \rho_i\}^2]} \geq \sqrt{\underline{\kappa}} \|\theta_h[T_h^c]\| - \sqrt{\mathbb{E}[\rho_i^2]}.$$

Thus we have $\|\theta_h[T_h^c]\| \leq \sqrt{\bar{\kappa}/\underline{\kappa}} \|\theta_h[S^c]\| + 2\sqrt{\mathbb{E}[\rho_i^2]/\underline{\kappa}}$. Moreover, since $\theta_h \in S_A^a(p)$, we have

$$\|\theta_h[S^c]\|^2 = \sum_{j=s+1}^{\infty} \theta_{h(j)}^2 \leq A^2 \sum_{j=s+1}^{\infty} j^{-2a} \leq A^2 s^{-2a+1} / [2a-1] \leq A^2 s^{-2a+1}$$

since $a > 1$. Combining these relations we have

$$\begin{aligned} \|\theta_h[T_h^c]\| &\leq \sqrt{\bar{\kappa}/\underline{\kappa}} A s^{-a+1/2} + 2\sqrt{\mathbb{E}[\rho_i^2]/\underline{\kappa}} \\ &= \sqrt{\bar{\kappa}/\underline{\kappa}} \sqrt{s/n} + 2\sqrt{\mathbb{E}[\rho_i^2]/\underline{\kappa}}. \end{aligned}$$

The second bound follows by observing that

$$\begin{aligned} \|\theta_h[T_h^c]\|_1 &\leq \sqrt{s} \|\theta_h[T_h^c \cap S]\| + \|\theta_h[S^c]\|_1 \leq \sqrt{s} \|\theta_h[T_h^c]\| + A s^{-a+1} / [a-1] \\ &\leq \sqrt{s^2/n} \sqrt{\bar{\kappa}/\underline{\kappa}} + 2\sqrt{s \mathbb{E}[\rho_i^2]/\underline{\kappa}} + (s/\sqrt{n}) / [a-1] \\ &\leq \sqrt{s^2/n} \sqrt{\bar{\kappa}/\underline{\kappa}} a / [a-1] + 2\sqrt{s \mathbb{E}[\rho_i^2]/\underline{\kappa}}. \end{aligned}$$

By the first-order optimality condition of the problem (4.31) that defines β_{h0} , we have

$$\mathbb{E}[x_i[T_h] x_i[T_h]'] (\beta_{h0}[T_h] - \theta_h[T_h]) = \mathbb{E}[x_i[T_h] x_i[T_h^c]'] \theta_h[T_h^c] + \mathbb{E}[x_i[T_h] \rho_i].$$

Thus, since $\|\mathbb{E}[x_i[T_h] \rho_i]\| = \sup_{\|\eta\|=1} \mathbb{E}[\eta' x_i[T_h] \rho_i] \leq \sup_{\|\eta\|=1} \sqrt{\mathbb{E}[(\eta' x_i[T_h])^2]} \sqrt{\mathbb{E}[\rho_i^2]}$ we have

$$\begin{aligned} \underline{\kappa} \|\beta_{h0} - \theta_h[T_h]\| &\leq \bar{\kappa} \|\theta_h[T_h^c]\| + \sqrt{\bar{\kappa} \mathbb{E}[\rho_i^2]} \\ &\leq \sqrt{s/n} (\bar{\kappa}^{3/2} / \sqrt{\underline{\kappa}}) + \sqrt{\mathbb{E}[\rho_i^2]} \sqrt{\bar{\kappa}} (1 + 2\sqrt{\bar{\kappa}/\underline{\kappa}}) \end{aligned}$$

where the last inequality follows from the definition of $s = A^{1/a} n^{1/2a}$. Therefore

$$\begin{aligned} |r_{hi}| &= |h_i - x_i' \beta_{h0}| = |x_i'(\theta_h - \beta_{h0})| + |\rho_i| \\ &\leq \|x_i\|_{\infty} \|\theta_h - \beta_{h0}\|_1 + |\rho_i| \\ &\leq \sqrt{s} \|x_i\|_{\infty} \|\theta_h[T_h] - \beta_{h0}\| + \|x_i\|_{\infty} \|\theta_h[T_h^c]\|_1 + |\rho_i| \\ &\leq \|x_i\|_{\infty} \{ \sqrt{s^2/n} (\bar{\kappa}/\underline{\kappa})^{3/2} + \sqrt{s \mathbb{E}[\rho_i^2]/\underline{\kappa}} \sqrt{\bar{\kappa}/\underline{\kappa}} (1 + 2\sqrt{\bar{\kappa}/\underline{\kappa}}) \} + \\ &\quad + \|x_i\|_{\infty} (\sqrt{s^2/n} \sqrt{\bar{\kappa}/\underline{\kappa}} a / [a-1] + 2\sqrt{s \mathbb{E}[\rho_i^2]/\underline{\kappa}}) + |\rho_i| \\ &\leq \|x_i\|_{\infty} (\bar{\kappa}/\underline{\kappa})^{3/2} \{ \frac{2a-1}{a-1} \sqrt{s^2/n} + 5\sqrt{s \mathbb{E}[\rho_i^2]/\underline{\kappa}} \} + |\rho_i|. \end{aligned}$$

□

Proof of Corollary 4 (Example 2). Let \mathbf{P} be the collection of all regression models P that obey the conditions set forth above for all n for the given constants $(\underline{\kappa}, \bar{\kappa}, a, A, B, \chi)$ and sequences p_n and $\bar{\delta}_n$. Below we provide explicit bounds for κ' , κ'' , c , C , δ_n and Δ_n that appear in Conditions ASTE, SE and SM that depend only on $(\underline{\kappa}, \bar{\kappa}, a, A, B, \chi)$, p , $\bar{\delta}_n$ and n which in turn establish these conditions for any $P \in \mathbf{P}$. In what follows we exploit Gaussianity of w_i and use that $(\mathbb{E}[|\eta' w_i|^k])^{1/k} \leq G_k (\mathbb{E}[|\eta' w_i|^2])^{1/2}$ for any vector η , $\|\eta\| < \infty$, where the constant G_k depends on k only.

Conditions ASTE(i) is assumed. Condition ASTE(ii) holds with $|\alpha_0| \leq B =: C_1^{ASTE}$. Because $\theta_m, \theta_g \in S_A^a(p)$, Condition ASTE(iii) holds with

$$s = A^{1/a} n^{1/2a}, \quad r_{mi} = m(z_i) - \sum_{j=1}^p z_{ij} \beta_{m0j}, \quad \text{and} \quad r_{gi} = g(z_i) - \sum_{j=1}^p z_{ij} \beta_{g0j}$$

where $\|\beta_{m0}\|_0 \leq s$ and $\|\beta_{g0}\|_0 \leq s$. Indeed, we have

$$\mathbb{E}[r_{mi}^2] \leq \mathbb{E} \left[\left(\sum_{j \geq s+1} \theta_{m(j)} z_{i(j)} \right)^2 \right] \leq \bar{\kappa} \sum_{j \geq s+1} \theta_{m(j)}^2 \leq \bar{\kappa} A^2 s^{-2a+1} / [2a-1] \leq \bar{\kappa} s / n$$

where the first inequality follows by the definition of β_{m0} in (4.31), the second inequality follows from $\theta_m \in S_A^a(p)$, and the last inequality because $s = A^{1/a} n^{1/2a}$. Similarly we have $\mathbb{E}[r_{gi}^2] \leq \mathbb{E}[(\sum_{j \geq s+1} \theta_{g(j)} z_{i(j)})^2] \leq \bar{\kappa} A^2 s^{-2a+1} / [2a-1] \leq \bar{\kappa} s / n$. Thus let $C_2^{ASTE} := \sqrt{f}$.

Condition ASTE(iv) holds with $\delta_{1n}^{ASTE} := A^{2/a} n^{1/a-1} \log^2(p \vee n) \rightarrow 0$ since $s = A^{1/a} n^{1/2a}$, A is fixed, and the assumed condition $n^{(1-a)/a} \log^2(p \vee n) \log^2 n \leq \bar{\delta}_n \rightarrow 0$.

The moment restrictions in Condition ASTE(v) are satisfied by the Gaussianity. Indeed, we have for $q = 4/\chi$ (where $\chi < 1$ by assumption)

$$\begin{aligned} \mathbb{E}[|\tilde{\zeta}_i|^q] &\leq 2^{q-1} \mathbb{E}[|\zeta_i^q|] + 2^{q-1} \mathbb{E}[|r_{gi}^q|] \leq 2^{q-1} G_q^q (\mathbb{E}[\zeta_i^2]^{q/2} + \mathbb{E}[r_{gi}^2]^{q/2}) \\ &\leq 2^{q-1} G_q^q \{ \bar{\kappa}^{q/2} + \bar{\kappa}^{q/2} (s/n)^{q/2} \} \\ &\leq 2^q G_q^q \bar{\kappa}^{q/2} =: C_3^{ASTE} \end{aligned}$$

for $s \leq n$, i.e., $n \geq n_{01}^{ASTE} := A^{2/[2a-1]}$. Similarly, $\mathbb{E}[|\tilde{v}_i|^q] \leq C_3^{ASTE}$. Moreover,

$$\begin{aligned} |\mathbb{E}[\tilde{\zeta}_i^2 \tilde{v}_i^2] - \mathbb{E}[\zeta_i^2 v_i^2]| &\leq \mathbb{E}[\zeta_i^2 r_{mi}^2] + \mathbb{E}[r_{gi}^2 v_i^2] + \mathbb{E}[r_{mi}^2 r_{gi}^2] \\ &\leq \sqrt{\mathbb{E}[\zeta_i^4] \mathbb{E}[r_{mi}^4]} + \sqrt{\mathbb{E}[r_{gi}^4] \mathbb{E}[v_i^4]} + \sqrt{\mathbb{E}[r_{mi}^4] \mathbb{E}[r_{gi}^4]} \\ &\leq G_4^2 \bar{\kappa} \mathbb{E}[r_{mi}^2] + G_4^2 \bar{\kappa} \mathbb{E}[r_{gi}^2] + G_4^2 \mathbb{E}[r_{mi}^2] \mathbb{E}[r_{gi}^2] \\ &\leq G_4^2 \bar{\kappa}^2 \{2 + \bar{\kappa} s / n\} s / n =: \delta_{2n}^{ASTE} \rightarrow 0. \end{aligned}$$

Next note that by Gaussian tail bounds and $\lambda_{\max}(\mathbb{E}[w_i w_i']) \leq \bar{\kappa}$ we have

$$\begin{aligned} \max_{i \leq n} \|x_i\|_\infty &\leq \|\mathbb{E}[x_i]\|_\infty + \max_{i \leq n} \|x_i - \mathbb{E}[x_i]\|_\infty \\ &\leq \sqrt{\bar{\kappa}} + \sqrt{2\bar{\kappa} \log(pn)} \quad \text{with probability at least } 1 - \Delta_{1n}^{ASTE} \end{aligned} \tag{G.67}$$

where $\Delta_{1n}^{ASTE} = 1/\sqrt{2\bar{\kappa} \log(pn)}$. The last requirement in Condition ASTE(v) holds with $q = 4/\chi$

$$\max_{i \leq n} \|x_i\|_\infty^2 s n^{-1/2+2/q} \leq 6\bar{\kappa} \log(pn) A^{1/a} n^{\frac{1}{2a}-\frac{1}{2}+\chi/2} =: \delta_{3n}^{ASTE}$$

with probability $1 - \Delta_{1n}^{ASTE}$. By the assumption on a, p, χ , and n , $\delta_{3n}^{ASTE} \rightarrow 0$.

To verify Condition SE with $\ell_n = \log n$ note that the minimal and maximal eigenvalues of $\mathbb{E}[x_i x_i']$ are bounded away from zero by $\underline{\kappa} > 0$ and from above by $\bar{\kappa} < \infty$ uniformly in n . Also, let $\mu = \mathbb{E}[x_i]$ so that $x_i = \tilde{x}_i + \mu$ where \tilde{x}_i is zero mean. By constricton $\mathbb{E}[x_i x_i'] = \mathbb{E}[\tilde{x}_i \tilde{x}_i'] + \mu \mu'$ and $\|\mu\| \leq \sqrt{\bar{\kappa}}$.

For any $\eta \in \mathbb{R}^p$, $\|\eta\|_0 \leq k := s \log n$ and $\|\eta\| = 1$, we have that

$$\mathbb{E}_n[(\eta' x_i)^2] - \mathbb{E}[(\eta' x_i)^2] = \mathbb{E}_n[(\eta' \tilde{x}_i)^2] - \mathbb{E}[(\eta' \tilde{x}_i)^2] + 2\eta' \mathbb{E}_n[\tilde{x}_i] \cdot \eta' \mu.$$

Moreover, by Gaussianity of x_i , with probability $1 - \Delta_{1n}^{SE}$, where $\Delta_{1n}^{SE} = 1/\sqrt{2\bar{\kappa} \log(pn)}$,

$$\begin{aligned} |\eta' \mathbb{E}_n[\tilde{x}_i]| &\leq \|\eta\|_1 \|\mathbb{E}_n[\tilde{x}_i]\|_\infty \leq \sqrt{k} \sqrt{2\bar{\kappa} \log(pn)} / \sqrt{n} \\ |\eta' \mu| &\leq \|\eta\| \|\mu\| \leq \sqrt{\bar{\kappa}}. \end{aligned}$$

By the sub-Gaussianity of $\tilde{x}_i = (\mathbb{E}[x_i x_i'] - \mu \mu')^{-1/2} \Psi_i$, where $\Psi_i \sim N(0, I_p)$, by Theorem 3.2 in Rudelson and Zhou (2011) (restated in Lemma 11 in Appendix G) with $\tau = 1/6$, $k = s \log n$, $\alpha = \sqrt{8/3}$, provided that

$$n \geq N_n := 80(\alpha^4/\tau^2)(s \log n \log(12ep/[\tau s \log n])),$$

we have

$$(1 - \tau)^2 \mathbb{E}[(\eta' \tilde{x}_i)^2] \leq \mathbb{E}_n[(\eta' \tilde{x}_i)^2] \leq (1 + \tau)^2 \mathbb{E}[(\eta' \tilde{x}_i)^2]$$

with probability $1 - \Delta_{1n}^{SE}$, where $\Delta_{1n}^{SE} = 2\exp(-\tau^2 n/80\alpha^4)$. Note that under ASTE(iv) we have $\Delta_{1n}^{SE} \rightarrow 0$ and

$$n_{01}^{SE} := \max\{n : n \leq N_n\} \leq \max\{(12e/\tau)^{2a} A^{-2}, 80^2(\alpha^8/\tau^4)A^{2/a}, n^*\}$$

where n^* is the smallest n such that $\bar{\delta}_n < 1$.

Therefore, with probability $1 - \Delta_{1n}^{SE}$ and $n \geq n_{01}^{SE}$, we have for any $\eta \in \mathbb{R}^p$, $\|\eta\|_0 \leq k$ and $\|\eta\| = 1$,

$$\begin{aligned} \mathbb{E}_n[(\eta' x_i)^2] &\geq \mathbb{E}[(\eta' x_i)^2] - |\mathbb{E}_n[(\eta' x_i)^2] - \mathbb{E}[(\eta' x_i)^2]| \\ &\geq \mathbb{E}[(\eta' x_i)^2] - |\mathbb{E}_n[(\eta' \tilde{x}_i)^2] - \mathbb{E}[(\eta' \tilde{x}_i)^2]| - 2|\eta' \mathbb{E}_n[\tilde{x}_i]| \cdot |\eta' \mu| \\ &\geq \mathbb{E}[(\eta' x_i)^2] \{1 - 2\tau - \tau^2\} - 2\bar{\kappa} \sqrt{2k \log(pn)} / \sqrt{n} \\ &\geq \mathbb{E}[(\eta' x_i)^2] / 2 - 2\bar{\kappa} \sqrt{2k \log(pn)} / \sqrt{n} \end{aligned}$$

since $\tau = 1/6$ and $\mathbb{E}[(\eta' \tilde{x}_i)^2] \leq \mathbb{E}[(\eta' x_i)^2]$. So for $n \geq n_{02}^{SE} := 288k(\bar{\kappa}/\underline{\kappa})^2 \log(pn)$ we have

$$\phi_{\min}(s \log n) [\mathbb{E}_n[x_i x_i']] \geq \underline{\kappa}/3 =: \kappa'.$$

Similarly, we have

$$\begin{aligned} \mathbb{E}_n[(\eta' x_i)^2] &\leq \mathbb{E}[(\eta' x_i)^2] + |\mathbb{E}_n[(\eta' x_i)^2] - \mathbb{E}[(\eta' x_i)^2]| \\ &\leq \mathbb{E}[(\eta' x_i)^2] + |\mathbb{E}_n[(\eta' \tilde{x}_i)^2] - \mathbb{E}[(\eta' \tilde{x}_i)^2]| + 2|\eta' \mathbb{E}_n[\tilde{x}_i]| \cdot |\eta' \mu| \\ &\leq \mathbb{E}[(\eta' x_i)^2] \{1 + 2\tau + \tau^2\} + 2\bar{\kappa} \sqrt{2k \log(pn)} / \sqrt{n} \\ &\leq 2\mathbb{E}[(\eta' x_i)^2] + 2\bar{\kappa} \sqrt{2k \log(pn)} / \sqrt{n} \end{aligned}$$

since $\tau = 1/6$ and $\mathbb{E}[(\eta' \tilde{x}_i)^2] \leq \mathbb{E}[(\eta' x_i)^2]$. So for $n \geq n_{03}^{SE} := 2k \log(pn)$ we have

$$\phi_{\max}(s \log n) [\mathbb{E}_n[x_i x_i']] \leq 4\bar{\kappa} =: \kappa''.$$

The second and third requirements in Conditions SM(i) holds by the Gaussianity of w_i , $E[\zeta_i | x_i, v_i] = 0$, $E[v_i | x_i] = 0$, and the assumption that the minimal and maximum eigenvalues of the covariance matrix (operator) $E[w_i w_i']$ are bounded below and above by positive absolute constants.

The first requirement in Condition SM(i) and Condition SM(ii) also hold by Gaussianity. Indeed, we have for $\epsilon_i = v_i$ and $\epsilon_i = \zeta_i$, $\tilde{y}_i = d_i$ and $\tilde{y}_i = y_i$

$$\begin{aligned}
E[|v_i^q|] + E[|\zeta_i^q|] &\leq 2^{q-1} G_q^q \{ (E[v_i^2])^{q/2} + (E[\zeta_i^2])^{q/2} \} \leq 2^q G_q^q \bar{\kappa}^{q/2} =: A_1 \\
E[|d_i^q|] &\leq 2^{q-1} E[|\theta'_m z|^q] + 2^{q-1} E[|v_i^q|] \leq 2^{q-1} G_q^q (E[|\theta'_m z|^2])^{q/2} + 2^{q-1} G_q^q (E[v_i^2])^{q/2} \\
&\leq 2^{q-1} G_q^q \|\theta_m\|^q \bar{\kappa}^{q/2} + 2^{q-1} G_q^q \bar{\kappa}^{q/2} \leq 2^q G_q^q \bar{\kappa}^{q/2} (1 + (2A)^q) =: A_2 \\
E[d_i^2] &\leq 2E[|\theta'_m z|^2] + 2E[v_i^2] \leq 2\bar{\kappa} \|\theta_m\|^2 + 2\bar{\kappa} \leq 2\bar{\kappa}(4A^2 + 1) =: A'_2 \\
E[y_i^2] &\leq 3|\alpha_0|^2 E[d_i^2] + 3E[|\theta'_m z|^2] + 3E[\zeta_i^2] \leq 3B^2 A'_2 + 3A'_2 + 3\bar{\kappa} =: A_3 \\
\max_{1 \leq j \leq p} E[x_{ij}^2 \tilde{y}_i^2] &\leq \max_{1 \leq j \leq p} (E[x_{ij}^4])^{1/2} (E[\tilde{y}_i^4])^{1/2} \leq G_4^4 \max_{1 \leq j \leq p} E[x_{ij}^2] E[\tilde{y}_i^2] \\
&\leq G_4^4 \bar{\kappa} (A'_2 \vee A_3) =: A_4 \\
\max_{1 \leq j \leq p} E[|x_{ij} \epsilon_i|^3] &\leq \max_{1 \leq j \leq p} (E[x_{ij}^6])^{1/2} (E[\epsilon_i^6])^{1/2} \leq G_6^6 \max_{1 \leq j \leq p} (E[x_{ij}^2])^{3/2} (E[\epsilon_i^2])^{3/2} \\
&\leq G_6^6 \bar{\kappa}^3 =: A_5 \\
\max_{1 \leq j \leq p} 1/E[x_{ij}^2] &\leq 1/\lambda_{\min}(E[w_i w_i']) \leq 1/\underline{\kappa} =: A_6
\end{aligned}$$

because $\|\theta_m\| \leq 2A$ and $\|\theta_g\| \leq 2A$ since $\theta_m, \theta_g \in S_A^g(p)$. Thus the first requirement in Condition SM(i) holds with $C_2^{SM} = A_2$. Condition SM(ii) holds with $C_3^{SM} = A_1 + (A'_2 \vee A_3) + A_4 + A_5 + A_6$.

Condition SM(iii) is assumed.

To verify Condition SM(iv) note that for $\epsilon_i = v_i$ and $\epsilon_i = \zeta_i$, by (G.67), with probability $1 - \Delta_{1n}^{ASTE}$,

$$\begin{aligned}
\max_{j \leq p} \sqrt{E_n[x_{ij}^4 \epsilon_i^4]} &\leq \max_{j \leq p} \sqrt[4]{E_n[x_{ij}^8]} \sqrt[4]{E_n[\epsilon_i^8]} \\
&\leq \{\sqrt{\bar{\kappa}} + \sqrt{2\bar{\kappa} \log(pn)}\} \max_{j \leq p} \sqrt[4]{E_n[x_{ij}^4]} \sqrt[4]{E_n[\epsilon_i^8]}.
\end{aligned} \tag{G.68}$$

By Lemma 3 with $k = 4$ we have with probability $1 - \Delta_{1n}^{SM}$, where $\Delta_{1n}^{SM} = 1/n$

$$\begin{aligned}
\max_{j \leq p} \sqrt[4]{E_n[x_{ij}^4]} &\leq \|E[x_i]\|_\infty + \max_{j \leq p} \sqrt[4]{E_n[(x_{ij} - E[x_{ij}])^4]} \\
&\leq \sqrt{\bar{\kappa}} + \sqrt{\bar{\kappa} 2\bar{C}} + \sqrt{\bar{\kappa} n^{-1/4}} \sqrt{2 \log(2pn)} \leq 4\bar{C} \sqrt{\bar{\kappa}}
\end{aligned} \tag{G.69}$$

for $n \geq n_{01}^{SM} = 4 \log^2(2pn)$. Also, Lemma 3 with $k = 8$ and $p = 1$ we have with probability $1 - \Delta_{1n}^{SM}$ that

$$\sqrt[4]{E_n[\epsilon_i^8]} \leq 2\bar{\kappa} 8\bar{C}^2 + 2\bar{\kappa} n^{-1/4} 2 \log(2n) \leq 20\bar{C}^2 \bar{\kappa} \tag{G.70}$$

for $n \geq n_{02}^{SM} = 16 \log^4(2n)$. Moreover, we have

$$\max_{1 \leq j \leq p} \sqrt{E[x_{ij}^4 \epsilon_i^4]} \leq \max_{1 \leq j \leq p} \sqrt[4]{E[x_{ij}^8]} \sqrt[4]{E[\epsilon_i^8]} \leq G_8^4 \bar{\kappa}^2.$$

Applying Lemma 7, for $\tau = 2\Delta_{1n}^{ASTE} + \Delta_{1n}^{SM}$, with probability $1 - 8\tau$ we have

$$\max_{j \leq p} |(\mathbb{E}_n - \bar{\mathbb{E}})[x_{ij}^2 \epsilon_i^2]| \leq 4\sqrt{\frac{2\log(2p/\tau)}{n}} \sqrt{Q(\max_{1 \leq j \leq p} \mathbb{E}_n[x_{ij}^4 \epsilon_i^4], 1 - \tau)} \vee \frac{2\sqrt{2}G_8^4 \bar{\kappa}^2}{\sqrt{n}}$$

where by (G.68), (G.69) and (G.70) we have

$$Q(\max_{1 \leq j \leq p} \sqrt{\mathbb{E}_n[x_{ij}^4 \epsilon_i^4]}, 1 - \tau) \leq \bar{\kappa}^2 \sqrt{2\log(pn)} 80\bar{C}^3.$$

So we let $\delta_{1n}^{SM} = 640\bar{C}^3 \bar{\kappa}^2 \sqrt{\frac{\log(2p/\tau)}{n}} \sqrt{\log(pn)} \vee 2\sqrt{2} \frac{G_8^4 \bar{\kappa}^2}{\sqrt{n}} \rightarrow 0$ under the condition that $\log^2(p \vee n)/n \leq \bar{\delta}_n$.

Similarly for $\tilde{y}_i = d_i$ and $\tilde{y}_i = y_i$, by Lemma 3, we have with probability $1 - \Delta_{1n}^{SM}$, for $n \geq n_{02}^{SM}$ we have

$$\begin{aligned} \sqrt[8]{\mathbb{E}_n[\tilde{y}_i^8]} &\leq |\mathbb{E}[\tilde{y}_i]| + \sqrt[8]{\mathbb{E}_n[(\tilde{y}_i - \mathbb{E}[\tilde{y}_i])^8]} \\ &\leq [A'_2 \vee A_3]^{1/2} + (20\bar{C}^2 \mathbb{E}[\tilde{y}_i^2])^{1/2} \leq 6\bar{C}[A'_2 \vee A_3]^{1/2}. \end{aligned} \quad (\text{G.71})$$

Moreover, $\sqrt[4]{\mathbb{E}[\tilde{y}_i^8]} \leq G_8^2 \mathbb{E}[\tilde{y}_i^2] \leq G_8^2 [A'_2 \vee A_3]$. Therefore by Lemma 7, for $\tau = 2\Delta_{1n}^{ASTE} + \Delta_{2n}^{SM}$, with probability $1 - 8\tau$ we have by the arguments in (G.68), (G.69), and (G.71)

$$\max_{j \leq p} |(\mathbb{E}_n - \bar{\mathbb{E}})[x_{ij}^2 \tilde{y}_i^2]| \leq 4\sqrt{\frac{2\log(2p/\tau)}{n}} \sqrt{6\bar{\kappa} \log(pn)} 4\bar{C} \sqrt{\bar{\kappa}} (36\bar{C}^2 [A'_2 \vee A_3]) \vee \frac{2\sqrt{2}G_8^4 \bar{\kappa} [A'_2 \vee A_3]}{\sqrt{n}} =: \delta_{2n}^{SM}$$

where $\delta_{2n}^{SM} \rightarrow 0$ under the condition $\log^2(p \vee n)/n \leq \bar{\delta}_n \rightarrow 0$.

We have that the last term in Condition SM(iv) satisfies with probability $1 - \Delta_{1n}^{ASTE}$

$$\max \|x_i\|_\infty^2 \frac{s \log(p \vee n)}{n} \leq 6\bar{\kappa} \log(pn) A^{1/a} n^{-1+1/2a} \log(p \vee n) =: \delta_{3n}^{SM}.$$

Under ASTE(iv) and $s = A^{1/a} n^{1/2a}$ we have $\delta_{3n}^{SM} \rightarrow 0$.

Finally, we set $n_0 = \max\{n_{01}^{ASTE}, n_{01}^{SE}, n_{02}^{SE}, n_{03}^{SE}, n_{01}^{SM}, n_{02}^{SM}\}$, $C = \max\{C_1^{ASTE}, C_2^{ASTE}, 2C_3^{ASTE}, C_1^{SM}, C_2^{SM}\}$, $\delta_n = \max\{\bar{\delta}_n, \delta_{1n}^{ASTE}, \delta_{2n}^{ASTE}, \delta_{1n}^{SM} + \delta_{2n}^{SM} + \delta_{3n}^{SM}\} \rightarrow 0$, and $\Delta_n = \max\{33\Delta_{1n}^{ASTE} + 16\Delta_{1n}^{SM}, \Delta_{1n}^{SE}\} \rightarrow 0$.

□

Lemma 3. Let $f_{ij} \sim N(0, \sigma_j^2)$, $\sigma_j \leq \sigma$, independent across $i = 1, \dots, n$, where $j = 1, \dots, p$. Then, for some universal constant $\bar{C} \geq 1$, we have that for any $k \geq 2$ and $\gamma \in (0, 1)$

$$P\left(\max_{1 \leq j \leq p} \{\mathbb{E}_n[|f_{ij}^k|]\}^{1/k} \geq \sigma \bar{C} \sqrt{k} + \sigma n^{-1/k} \sqrt{2\log(2p/\gamma)}\right) \leq \gamma.$$

Proof. Note that $P(\mathbb{E}_n[|f_{ij}^k|] > M) = P(\|f_{\cdot j}\|_k^k > Mn) = P(\|f_{\cdot j}\|_k > (Mn)^{1/k})$.

Since $\|f\|_k - \|g\|_k \leq \|f - g\|_k \leq \|f - g\|$, we have that $\|\cdot\|_k$ is 1-Lipschitz for $k \geq 2$. Moreover,

$$\mathbb{E}[\|f_{\cdot j}\|_k] \leq (\mathbb{E}[\|f_{\cdot j}\|_k^k])^{1/k} = \left(\sum_{i=1}^n \mathbb{E}[|f_{ij}^k|]\right)^{1/k} = n^{1/k} (\mathbb{E}[|f_{1j}^k|])^{1/k}$$

$$= n^{1/k} \{\sigma_j^k 2^{k/2} \Gamma((k+1)/2) / \Gamma(1/2)\}^{1/k} \leq n^{1/k} \sigma \sqrt{k} \bar{C}.$$

By Ledoux and Talagrand (1991), page 21 equation (1.6), we have

$$P(\|f_{\cdot j}\|_k > (Mn)^{1/k}) \leq 2 \exp(-\{(Mn)^{1/k} - \mathbb{E}[\|f_{\cdot j}\|_k]\}^2 / 2\sigma_j^2).$$

Setting $M := \{\sigma \sqrt{k} \bar{C} + \sigma n^{-1/k} \sqrt{2 \log(2p/\gamma)}\}^k$, so that $(Mn)^{1/k} = n^{1/k} \sigma \sqrt{k} \bar{C} + \sigma \sqrt{2 \log(2p/\gamma)}$ we have by the union bound and $\sigma \geq \sigma_j$

$$P(\max_{1 \leq j \leq p} \mathbb{E}_n[\|f_{ij}^k\|] \geq M) \leq p \max_{1 \leq j \leq p} P(\mathbb{E}_n[\|f_{ij}^k\|] \geq M) \leq \gamma.$$

□

Proof for Corollary 5 (Example 3). Let \mathbf{P} be the collection of all regression models P that obey the conditions set forth above for all n for the given constants $(\underline{f}, \bar{f}, a, A, b, B, q)$ and the sequence $\bar{\delta}_n$. Below we provide explicit bounds for $\kappa', \kappa'', c, C, \delta_n$ and Δ_n that appear in Conditions ASTE, SE and SM that depend only on $(\underline{f}, \bar{f}, a, A, b, B, q)$ and $\bar{\delta}_n$ which in turn establish these conditions for all $P \in \mathbf{P}$.

Conditions ASTE(i) is assumed. Condition ASTE(ii) holds with $|\alpha_0| \leq B =: C_1^{ASTE}$. Because $\theta_m, \theta_g \in S_A^a(p)$, Condition ASTE(iii) holds with

$$s = A^{1/a} n^{\frac{1}{2a}}, \quad r_{mi} = m(z_i) - \sum_{j=1}^p \beta_{m0j} P_j(z_i) \quad \text{and} \quad r_{gi} = g(z_i) - \sum_{j=1}^p \beta_{g0j} P_j(z_i)$$

where $\|\beta_{m0}\|_0 \leq s$ and $\|\beta_{g0}\|_0 \leq s$. Indeed, we have

$$\mathbb{E}[r_{mi}^2] \leq \mathbb{E} \left[\left(\sum_{j \geq s+1} \theta_{m(j)} P_j(z_i) \right)^2 \right] \leq \bar{f} \sum_{j \geq s+1} \theta_{m(j)}^2 \leq \bar{f} A^2 s^{-2a+1} / [2a-1] = \bar{f} s / n$$

where the first inequality follows by the definition of β_{m0} in (4.31), the second inequality follows from the upper bound on the density and orthogonality of the basis, the third inequality follows from $\theta_m \in S_A^a(p)$, and the last inequality because $s = A^{1/a} n^{1/2a}$. Similarly we have $\mathbb{E}[r_{gi}^2] \leq \mathbb{E}[(\sum_{j \geq s+1} \theta_{g(j)} z_{i(j)})^2] \leq \bar{f} A^2 s^{-2a+1} / [2a-1] = \bar{f} s / n$. Let $C_2^{ASTE} = \sqrt{\bar{f}}$.

Condition ASTE(iv) holds with $\delta_{1n}^{ASTE} := A^{2/a} n^{1/a-1} \log^2(p \vee n) \rightarrow 0$ since $s = A^{1/a} n^{1/2a}$, A is fixed, and the assumed condition $n^{(1-a)/a} \log^2(p \vee n) \leq \bar{\delta}_n \rightarrow 0$.

Next we establish the moment restrictions in Condition ASTE(v). Because $\underline{f} \leq \lambda_{\min}(\mathbb{E}[x_i x_i']) \leq \lambda_{\max}(\mathbb{E}[x_i x_i']) \leq \bar{f}$, by the assumption on the density and orthonormal basis, and $\max_{i \leq n} \|x_i\|_\infty \leq B$, by Lemma 2 with $\rho_i = 0$ we have

$$\max_{1 \leq i \leq n} |r_{mi}| \vee |r_{gi}| \leq \max_{1 \leq i \leq n} \|x_i\|_\infty (\bar{f}/\underline{f})^{3/2} \frac{2a-1}{a-1} \sqrt{s^2/n} \leq B(\bar{f}/\underline{f})^{3/2} \frac{2a-1}{a-1} \sqrt{s^2/n} =: \delta_{2n}^{ASTE}$$

where $\delta_{2n}^{ASTE} \rightarrow 0$ under $s = A^{1/a} n^{1/2a}$ and $a > 1$.

Thus we have

$$\begin{aligned} \mathbb{E}[|\tilde{\zeta}_i|^q] &\leq 2^{q-1}\mathbb{E}[|\zeta_i^q|] + 2^{q-1}\mathbb{E}[|r_{gi}^q|] \leq 2^{q-1}B + 2^{q-1}(\delta_{2n}^{ASTE})^q \\ &\leq 2^{q-1}B + 2^{q-1}(\delta_{2n_0}^{ASTE})^q =: C_3^{ASTE}. \end{aligned}$$

Similarly, $\mathbb{E}[|\tilde{v}_i|^q] \leq C_3^{ASTE}$. Moreover, since $\delta_{2n}^{ASTE} \rightarrow 0$ we have

$$\begin{aligned} |\mathbb{E}[\tilde{\zeta}_i^2 \tilde{v}_i^2] - \mathbb{E}[\zeta_i^2 v_i^2]| &\leq \mathbb{E}[\zeta_i^2 r_{mi}^2] + \mathbb{E}[r_{gi}^2 v_i^2] + \mathbb{E}[r_{mi}^2 r_{gi}^2] \\ &\leq \sqrt{\mathbb{E}[\zeta_i^4] \mathbb{E}[r_{mi}^4]} + \sqrt{\mathbb{E}[r_{gi}^4] \mathbb{E}[v_i^4]} + \sqrt{\mathbb{E}[r_{mi}^4] \mathbb{E}[r_{gi}^4]} \\ &\leq 2B^{2/q}(\delta_{2n}^{ASTE})^2 + (\delta_{2n}^{ASTE})^4 =: \delta_{3n}^{ASTE} \rightarrow 0. \end{aligned}$$

Finally, the last requirement holds because $(1-a)/a + 4/q < 0$ implies

$$\max_{i \leq n} \|x_i\|_\infty^2 s n^{-1/2+2/q} \leq B^2 A^{1/a} n^{1/2a-1/2+2/q} =: \delta_{4n}^{ASTE} \rightarrow 0,$$

since $s = A^{1/a} n^{1/2a}$ and $\max_{i \leq n} \|x_i\|_\infty \leq B$.

To show Condition SE with $\ell_n = \log n$ note that regressors are uniformly bounded, and minimal and maximal eigenvalues of $\mathbb{E}[x_i x_i']$ are bounded below by \underline{f} and above by \bar{f} uniformly in n . Thus Condition SE follows by Corollary 4 in the supplementary material in Belloni and Chernozhukov (2013) (restated in Lemma 10 in Appendix G) which is based on Rudelson and Vershynin (2008). Let

$$\delta_{1n}^{SE} := 2\bar{C}B\sqrt{s \log n} \log(1 + s \log n) \sqrt{\log(p \vee n)} \sqrt{\log n} / \sqrt{n}$$

and $\Delta_{1n}^{SE} := (2/\underline{f})(\delta_{1n}^{SE})^2 + \delta_{1n}^{SE}(2\bar{f}/\underline{f})$, where \bar{C} is an universal constant. By this result and the Markov inequality, we have with probability $1 - \Delta_{1n}^{SE}$

$$\kappa' := \underline{f}/2 \leq \phi_{\min}(s \log n) [\mathbb{E}_n[x_i x_i']] \leq \phi_{\max}(s \log n) [\mathbb{E}_n[x_i x_i']] \leq 2\bar{f} =: \kappa''.$$

We need to show that $\Delta_{1n}^{SE} \rightarrow 0$ which follows from $\delta_{1n}^{SE} \rightarrow 0$. We have that

$$\delta_{1n}^{SE} \leq \frac{2\bar{C}B(1+A)^2 \sqrt{n^{1/2a}} \log^2(n) \sqrt{\log(p \vee n)}}{\sqrt{n}} = 2\bar{C}B(1+A)^2 \sqrt{\frac{n^{1/2a} \log^4 n}{n^{2/3}}} \sqrt{\frac{\log(p \vee n)}{n^{1/3}}}.$$

By assumption we have $\log^3 p/n \leq \bar{\delta}_n \rightarrow 0$ and $a > 1$ we have $\delta_{1n}^{SE} \rightarrow 0$.

The second and third requirements in Condition SM(i) hold with $C_1^{SM} = B^{2/q}$ and $c_1^{SM} = b$ by assumption. Condition SM(iii) is assumed.

The first requirement in Condition SM(i) and Condition SM(ii) follow by, for $\epsilon_i = v_i$ and $\epsilon_i = \zeta_i$, $\tilde{y}_i = d_i$ and $\tilde{y}_i = y_i$

$$\begin{aligned}
\mathbb{E}[|v_i^q|] + \mathbb{E}[|\zeta_i^q|] &\leq 2B =: A_1 \\
\mathbb{E}[|d_i^q|] &\leq 2^{q-1}\mathbb{E}[|\theta'_m x_i|^q] + 2^{q-1}\mathbb{E}[|v_i^q|] \leq 2^{q-1}\|\theta_m\|_1^q \mathbb{E}[\|x_i\|_\infty^q] + 2^{q-1}B \\
&\leq 2^{q-1}(2A)^q B^q + 2^{q-1}B =: A_2 \\
\mathbb{E}[d_i^2] &\leq 2\bar{f}\|\theta_m\|^2 + 2\mathbb{E}[v_i^2] \leq 8\bar{f}A^2 + 2B^{2/q} =: A'_2 \\
\mathbb{E}[y_i^2] &\leq 3|\alpha_0|^2\mathbb{E}[d_i^2] + 3\|\theta_g\|_1^2\mathbb{E}[\|x_i\|_\infty^2] + 3\mathbb{E}[\zeta_i^2] \\
&\leq 3B^2A'_2 + 12A^2B^2 + 3B^{2/q} =: A_3 \\
\max_{1 \leq j \leq p} \mathbb{E}[x_{ij}^2 \tilde{y}_i^2] &\leq B^2\mathbb{E}[\tilde{y}_i^2] \leq B^2(A'_2 \vee A_3) =: A_4 \\
\max_{1 \leq j \leq p} \mathbb{E}[|x_{ij}\epsilon_i|^3] &\leq B^3\mathbb{E}[|\epsilon_i^3|] \leq B^3B^{3/q} =: A_5 \\
\max_{1 \leq j \leq p} 1/\mathbb{E}[x_{ij}^2] &\leq 1/\lambda_{\min}(\mathbb{E}[x_i x_i']) \leq 1/\underline{f} =: A_6
\end{aligned}$$

where we used that $\max_{i \leq n} \|x_i\|_\infty \leq B$, the moment assumptions of the disturbances, $\|\theta_m\| \leq \|\theta_m\|_1 \leq 2A$, $\|\theta_g\|_1 \leq 2A$ since $\theta_m, \theta_g \in S_A^a(p)$ for $a > 1$. Thus the first requirement in Condition SM(i) holds with $C_2^{SM} = A_2$. Condition SM(ii) holds with $C_3^{SM} := A_1 + (A'_2 \vee A_3) + A_4 + A_5 + A_6$.

To verify Condition SM(iv) note that for $\epsilon_i = v_i$ and $\epsilon_i = \zeta_i$ we have by Lemma 7 with probability $1 - 8\tau$, where $\tau = 1/\log n$,

$$\begin{aligned}
\max_{1 \leq j \leq p} |(\mathbb{E}_n - \bar{\mathbb{E}})[x_{ij}^2 \epsilon_i^2]| &\leq 4\sqrt{\frac{2\log(2p/\tau)}{n}} Q(\max_{1 \leq j \leq p} \sqrt{\mathbb{E}_n[x_{ij}^4 \epsilon_i^4]}, 1 - \tau) \vee \frac{2 \max_{1 \leq j \leq p} \sqrt{2\mathbb{E}[x_{ij}^4 \epsilon_i^4]}}{\sqrt{n}} \\
&\leq 4\sqrt{\frac{2\log(2p/\tau)}{n}} B^2 Q(\sqrt{\mathbb{E}_n[\epsilon_i^4]}, 1 - \tau) \vee \frac{2B^2 \sqrt{2\mathbb{E}[\epsilon_i^4]}}{\sqrt{n}} \\
&\leq 4\sqrt{\frac{2\log(2p \log n)}{n}} B^2 B^{2/q} \log n =: \delta_{1n}^{SM}
\end{aligned}$$

where we used $\mathbb{E}[\epsilon_i^4] \leq B^{4/q}$ and the Markov inequality. By the definition of τ and the assumed rate $\log^3(p \vee n)/n \leq \bar{\delta}_n \rightarrow 0$, we have $\delta_{1n}^{SM} \rightarrow 0$.

Similarly, we have for $\tilde{y}_i = d_i$ and $\tilde{y}_i = y_i$, with probability $1 - 8\tau$

$$\begin{aligned}
\max_{1 \leq j \leq p} |(\mathbb{E}_n - \bar{\mathbb{E}})[x_{ij}^2 \tilde{y}_i^2]| &\leq 4\sqrt{\frac{2\log(2p/\tau)}{n}} Q(\max_{1 \leq j \leq p} \sqrt{\mathbb{E}_n[x_{ij}^4 \tilde{y}_i^4]}, 1 - \tau) \vee \frac{2 \max_{1 \leq j \leq p} \sqrt{2\mathbb{E}[x_{ij}^4 \tilde{y}_i^4]}}{\sqrt{n}} \\
&\leq 4\sqrt{\frac{2\log(2p/\tau)}{n}} B^2 Q(\sqrt{\mathbb{E}_n[\tilde{y}_i^4]}, 1 - \tau) \vee \frac{2B^2 \sqrt{2\mathbb{E}[\tilde{y}_i^4]}}{\sqrt{n}} \\
&\leq 4\sqrt{\frac{2\log(2p \log n)}{n}} B^2 A_7 \log n =: \delta_{2n}^{SM}
\end{aligned}$$

where we used the Markov inequality and

$$\begin{aligned}
\mathbb{E}[\tilde{y}_i^4] &\leq \mathbb{E}[d_i^4] + 3^3|\alpha_0|^4\mathbb{E}[d_i^4] + 3^3\|\theta_g\|_1^4\mathbb{E}[\|x_i\|_\infty^4] + 3^3\mathbb{E}[\zeta_i^4] \\
&\leq A_2^{4/q} + 3^3B^4A_2^{4/q} + 3^3(2A)^4B^4 + 3^3B^{4/q} =: A_7.
\end{aligned}$$

By the definition of τ and the assumed rate $\log^3(p \vee n)/n \leq \bar{\delta}_n \rightarrow 0$, we have $\delta_{2n}^{SM} \rightarrow 0$.

The last term in the requirement of Condition SM(iv), because $\max_{i \leq n} \|x_i\|_\infty \leq B$ and Condition ASTE(iv) holds, is bounded by $\delta_{3n}^{SM} := B^2 A^{1/a} n^{1/2a} \log(p \vee n)/n \rightarrow 0$.

Finally, we set $c = c_1^{SM}$, $C = \max\{C_1^{ASTE}, C_2^{ASTE}, 2C_3^{ASTE}, C_1^{SM}, C_2^{SM}, C_3^{SM}\}$,

$$\delta_n = \max\{\bar{\delta}_n, \delta_{1n}^{ASTE}, \delta_{2n}^{ASTE}, \delta_{3n}^{ASTE}, \delta_{4n}^{ASTE}, \delta_{1n}^{SM} + \delta_{2n}^{SM} + \delta_{3n}^{SM}\} \rightarrow 0,$$

$$\Delta_n = \max\{16/\log n, \Delta_{1n}^{SE}\} \rightarrow 0. \quad \square$$

APPENDIX H. PROOF OF THEOREMS 3 AND 4.

The two results have identical structure and have nearly the same proof, and so we present the proof of Theorem 3 only.

In the proof $a \lesssim b$ means that $a \leq Ab$, where the constant A depends on the constants in Condition HT only, but not on n once $n \geq n_0 = \min\{j : \delta_j \leq 1/2\}$, and not on $P \in \mathbf{P}_n$. For the proof of claims (1) and (2) we consider a sequence P_n in \mathbf{P}_n , but for simplicity, we write P throughout the proof, omitting the index n . Since the argument is asymptotic, we can just assume that $n \geq n_0$ in what follows.

Step 1. In this step we establish claim (1).

(a) We begin with a preliminary observation. Define, for $t = (t_1, t_2, t_3)$,

$$\psi(y, d, t) = \frac{d(y - t_2)}{t_3} - \frac{(1 - d)(y - t_1)}{1 - t_3} + t_2 - t_1.$$

The derivatives of this function with respect to t obey for all $k = (k_j)_{j=1}^3 \in \mathbb{N} : 0 \leq |k| \leq 3$,

$$|\partial_t^k \psi(y, d, t)| \leq L, \quad \forall (y, d, t) : |y| \leq C, |t_1| \leq C, |t_2| \leq C, c'/2 \leq |t_3| \leq 1 - c'/2, \quad (\text{H.72})$$

where L depends only on c' and C , $|k| = \sum_{j=1}^3 k_j$, and

$$\partial_t^k := \partial_{t_1}^{k_1} \partial_{t_2}^{k_2} \partial_{t_3}^{k_3}.$$

(b). Let

$$\hat{h}(z_i) := (\hat{g}(0, z_i), \hat{g}(1, z_i), \hat{m}(z_i))', \quad h_0(z_i) := (g(0, z_i), g(1, z_i), m(z_i))',$$

$$f_{\hat{h}}(y_i, d_i, z_i) := \psi(y_i, d_i, \hat{h}(z_i)), \quad f_{h_0}(y_i, d_i, z_i) := \psi(y_i, d_i, h_0(z_i)).$$

We observe that with probability no less than $1 - \Delta_n$,

$$\hat{g}(0, \cdot) \in \mathcal{G}_0, \quad \hat{g}(1, \cdot) \in \mathcal{G}_1 \text{ and } \hat{m} \in \mathcal{M},$$

$$\mathcal{G}_d := \{z \mapsto x' \beta : \|\beta\|_0 \leq sC, \|x'_i \beta - g(d, z_i)\|_{P,2} \lesssim \delta_n n^{-1/4}, \|x'_i \beta - g(d, z_i)\|_{P,\infty} \lesssim \delta_n\},$$

$$\mathcal{M} := \{z \mapsto \Lambda(x' \beta) : \|\beta\|_0 \leq sC, \|\Lambda(x'_i \beta) - m(z_i)\|_{P,2} \lesssim \delta_n n^{-1/4}, \|\Lambda(x'_i \beta) - m(z_i)\|_{P,\infty} \lesssim \delta_n\}.$$

To see this note, that under assumption HT (P), under condition (i)-(ii), under the event occurring under condition (ii) of that assumption: for $n \geq n_0 = \min\{j : \delta_j \leq 1/2\}$:

$$\begin{aligned} \|\tilde{x}'_i \beta - g(d_i, z_i)\|_{P,2} &\leq \|\tilde{x}'_i(\beta - \beta_{g0})\|_{P,2} + \|r_{gi}\|_{P,2} \leq 2\|\tilde{x}'_i(\beta - \beta_{g0})\|_{\mathbb{P}_{n,2}} + \|r_{gi}\|_{P,2} \leq 4\delta_n n^{-1/4}, \\ \|\tilde{x}'_i \beta - g(d_i, z_i)\|_{P,\infty} &\leq \|\tilde{x}'_i(\beta - \beta_{g0})\|_{P,\infty} + \|r_{gi}\|_{P,\infty} \leq K_n \|\beta - \beta_{g0}\|_1 + \delta_n \leq 2\delta_n, \end{aligned}$$

for $\beta = \hat{\beta}_g$, with evaluation after computing the norms, and noting that for any β

$$\|x'_i \beta - g(1, z_i)\|_{P,2} \vee \|x'_i \beta - g(0, z_i)\|_{P,2} \lesssim \|\tilde{x}'_i \beta - g(d_i, z_i)\|_{P,2}$$

under condition (iii). Furthermore, for $n \geq n_0 = \min\{j : \delta_j \leq 1/2\}$:

$$\begin{aligned} \|\Lambda(x'_i \beta) - m(z_i)\|_{P,2} &\leq \|\Lambda(x'_i \beta) - \Lambda(x'_i \beta_{m0})\|_{P,2} + \|r_{mi}\|_{P,2} \\ &\lesssim \|\partial \Lambda\|_{\infty} \|\tilde{x}'_i(\beta - \beta_{m0})\|_{P,2} + \|r_{mi}\|_{P,2} \\ &\lesssim \|\partial \Lambda\|_{\infty} \|\tilde{x}'_i(\beta - \beta_{m0})\|_{\mathbb{P}_{n,2}} + \|r_{mi}\|_{P,2} \lesssim \delta_n n^{-1/4} \\ \|\Lambda(x'_i \beta) - m(z_i)\|_{P,\infty} &\leq \|\partial \Lambda\|_{\infty} \|\tilde{x}'_i(\beta - \beta_{g0})\|_{P,\infty} + \|r_{mi}\|_{P,\infty} \\ &\lesssim K_n \|\beta - \beta_{m0}\|_1 + \delta_n \leq 2\delta_n, \end{aligned}$$

for $\beta = \hat{\beta}_{m0}$, with evaluation after computing the norms.

Hence with probability at least $1 - \Delta_n$,

$$\hat{h} \in \mathcal{H}_n := \{h = (\bar{g}(0, z), \bar{g}(1, z), \bar{m}(z)) \in \mathcal{G}_0 \times \mathcal{G}_1 \times \mathcal{M}\}.$$

(c) We have that

$$\alpha_0 = \mathbb{E}[f_{h_0}] \text{ and } \check{\alpha} = \mathbb{E}_n[f_{\hat{h}}],$$

so that

$$\sqrt{n}(\check{\alpha} - \alpha_0) = \underbrace{\mathbb{G}_n[f_{h_0}]}_i + \underbrace{(\mathbb{G}_n[f_{\hat{h}}] - \mathbb{G}_n[f_{h_0}])}_{ii} + \underbrace{\sqrt{n}(\mathbb{E}[f_{\hat{h}}] - f_{h_0})}_{iii},$$

with h evaluated at $h = \hat{h}$. By the Lyapunov central limit theorem,

$$\sigma_n^{-1} i \rightsquigarrow N(0, 1).$$

(d) Note that for $\Delta_i = h(z_i) - h_0(z_i)$,

$$\begin{aligned} iii &= \sqrt{n} \sum_{|k|=1} \mathbb{E}[\partial_t^k \psi(y_i, d_i, h_0(z_i)) \Delta_i^k] \\ &+ \sqrt{n} \sum_{|k|=2} 2^{-1} \mathbb{E}[\partial_t^k \psi(y_i, d_i, h_0(z_i)) \Delta_i^k] \\ &+ \sqrt{n} \sum_{|k|=3} \int_0^1 6^{-1} \mathbb{E}[\partial_t^k \psi(y_i, d_i, h_0(z_i) + \lambda \Delta_i) \Delta_i^k] d\lambda, \\ &=: iii_a + iii_b + iii_c, \end{aligned}$$

(with h evaluated at $h = \hat{h}$). By the law of iterated expectations and because

$$\mathbb{E}[\partial_t^k \psi(y_i, d_i, h_0(d_i, z_i)) | d_i, z_i] = 0 \quad \forall m \in \mathbb{N}^3 : |k| = 1,$$

we have that

$$iii_a = 0.$$

Moreover, uniformly for any $h \in \mathcal{H}_n$ we have that

$$|iii_b| \lesssim \sqrt{n} \|h - h_0\|_{\mathbb{P},2}^2 \lesssim \sqrt{n} (\delta_n n^{-1/4})^2 \leq \delta_n^2,$$

$$|iii_c| \lesssim \sqrt{n} \|h - h_0\|_{\mathbb{P},2}^2 \|h - h_0\|_{\mathbb{P},\infty} \lesssim \sqrt{n} (\delta_n n^{-1/4})^2 \delta_n \leq \delta_n^3.$$

Since $\hat{h} \in \mathcal{H}_n$ with probability $1 - \Delta_n$, we have that once $n \geq n_0$,

$$\mathbb{P}(|iii| \lesssim \delta_n^2) \geq 1 - \Delta_n.$$

(e). Furthermore, we have that

$$|ii| \leq \sup_{h \in \mathcal{H}_n} |\mathbb{G}_n[f_h] - \mathbb{G}_n[f_{h_0}]|.$$

The class of functions \mathcal{G}_d for $d \in \{0, 1\}$ is a union of at most $\binom{p}{C_s}$ VC-subgraph classes of functions with VC indices bounded by $C's$. The class of functions \mathcal{M} is a union of at most $\binom{p}{C_s}$ VC-subgraph classes of functions with VC indices bounded by $C's$ (monotone transformation Λ preserve the VC-subgraph property). These classes are uniformly bounded and their entropies therefore satisfy

$$\log N(\varepsilon, \mathcal{M}, \|\cdot\|_{\mathbb{P}_{n,2}}) + \log N(\varepsilon, \mathcal{G}_0, \|\cdot\|_{\mathbb{P}_{n,2}}) + \log N(\varepsilon, \mathcal{G}_1, \|\cdot\|_{\mathbb{P}_{n,2}}) \lesssim s \log p + s \log(1/\varepsilon).$$

Finally, the class $\mathcal{F}_n = \{f_h - f_{h_0} : h \in \mathcal{H}_n\}$ is a Lipschitz transform of \mathcal{H}_n with bounded Lipschitz coefficients and with a constant envelope. Therefore, we have that

$$\log N(\varepsilon, \mathcal{F}_n, \|\cdot\|_{\mathbb{P}_{n,2}}) \lesssim s \log p + s \log(1/\varepsilon).$$

We shall invoke the following lemma derived in Belloni and Chernozhukov (2011a).

Lemma 4 (A Self-Normalized Maximal inequality). *Let \mathcal{F} be a measurable function class on a sample space. Let $F = \sup_{f \in \mathcal{F}} |f|$, and suppose that there exist some constants $\omega_n > 3$ and $v > 1$, such that*

$$\log N(\epsilon \|F\|_{\mathbb{P}_{n,2}}, \mathcal{F}, \|\cdot\|_{\mathbb{P}_{n,2}}) \leq vm(\log(n \vee \omega_n) + \log(1/\epsilon)), \quad 0 < \epsilon < 1.$$

Then for every $\delta \in (0, 1/6)$ we have

$$\sup_{f \in \mathcal{F}} |\mathbb{G}_n(f)| \leq C_v \sqrt{2/\delta} \sqrt{m \log(n \vee \omega_n)} (\sup_{f \in \mathcal{F}} \|f\|_{\mathbb{P},2} \vee \sup_{f \in \mathcal{F}} \|f\|_{\mathbb{P}_{n,2}}),$$

with probability at least $1 - \delta$ for some constant that C_v .

Then by Lemma 4 together and some simple calculations, we have that

$$\begin{aligned} |ii| &\leq \sup_{f \in \mathcal{F}_n} |\mathbb{G}_n(f)| = O_P(1) \sqrt{s \log(p \vee n)} \left(\sup_{f \in \mathcal{F}_n} \|f\|_{\mathbb{P}_{n,2}} \vee \sup_{f \in \mathcal{F}_n} \|f\|_{P,2} \right) \\ &\leq O_P(1) \sqrt{s \log(p \vee n)} \left(\sup_{h \in \mathcal{H}_n} \|h - h_0\|_{\mathbb{P}_{n,2}} \vee \sup_{h \in \mathcal{H}_n} \|h - h_0\|_{P,2} \right) = o_P(1). \end{aligned}$$

The last conclusion follows because $\sup_{h \in \mathcal{H}_n} \|h - h_0\|_{P,2} \lesssim \delta_n n^{-1/4}$ by definition of the norm, and

$$\sup_{h \in \mathcal{H}_n} \|h - h_0\|_{\mathbb{P}_{n,2}} \leq O_P(1) \cdot \left(\sup_{h \in \mathcal{H}_n} \|h - h_0\|_{P,2} + \|r_{gi}\|_{P,2} + \|r_{mi}\|_{P,2} \right),$$

where the last conclusion follows from the same argument as in step (b) but in a reverse order, switching from empirical norms to population norms, using equivalence of norms over sparse sets imposed in condition (ii), and also using an application of Markov inequality to argue that $\|r_{gi}\|_{\mathbb{P}_{n,2}} + \|r_{mi}\|_{\mathbb{P}_{n,2}} = O_P(1)(\|r_{gi}\|_{P,2} + \|r_{mi}\|_{P,2})$.

Step 2. Claim (2) follows from consistency: $\hat{\sigma}_n/\sigma_n = 1 + o_P(1)$, which follows from $\hat{\sigma}_n$ being a Lipschitz transform of \hat{h} with respect to $\|\cdot\|_{\mathbb{P}_{n,2}}$, once $\hat{h} \in \mathcal{H}_n$ and the consistency of \hat{h} for h under $\|\cdot\|_{\mathbb{P}_{n,2}}$.

Step 3. Claim (3) is immediate from claims (2) and (3) by the way of contradiction. \square

APPENDIX I. TOOLS

I.0.1. Moderate Deviations for a Maximum of Self-Normalized Averages.

Lemma 5 (Moderate Deviation Inequality for Maximum of a Vector). *Suppose that*

$$\mathcal{S}_j = \frac{\sum_{i=1}^n U_{ij}}{\sqrt{\sum_{i=1}^n U_{ij}^2}},$$

where U_{ij} are independent variables across i with mean zero. We have that

$$P \left(\max_{1 \leq j \leq p} |\mathcal{S}_j| > \Phi^{-1}(1 - \gamma/2p) \right) \leq \gamma \left(1 + \frac{A}{\ell_n^3} \right),$$

where A is an absolute constant, provided for $\ell_n > 0$

$$0 \leq \Phi^{-1}(1 - \gamma/(2p)) \leq \frac{n^{1/6}}{\ell_n} \min_{1 \leq j \leq p} M_j^2 - 1, \quad M_j := \frac{\left(\frac{1}{n} \sum_{i=1}^n E[U_{ij}^2] \right)^{1/2}}{\left(\frac{1}{n} \sum_{i=1}^n E[|U_{ij}|^3] \right)^{1/3}}.$$

The proof of this result, given in Belloni, Chen, Chernozhukov, and Hansen (2012), follows from a simple combination of union bounds with the bounds in Theorem 7.4 in de la Peña, Lai, and Shao (2009) (see also (Jing, Shao, and Wang, 2003)).

I.1. Inequalities based on Symmetrization. Next we proceed to use symmetrization arguments to bound the empirical process. In what follows for a random variable Z let $Q(Z, 1 - \tau)$ denote its $(1 - \tau)$ -quantile.

Lemma 6 (Maximal inequality via symmetrization). *Let Z_1, \dots, Z_n be arbitrary independent stochastic processes and \mathcal{F} a finite set of measurable functions. For any $\tau \in (0, 1/2)$, and $\delta \in (0, 1)$ we have that with probability at least $1 - 4\tau - 4\delta$*

$$\max_{f \in \mathcal{F}} |\mathbb{G}_n(f(Z_i))| \leq \left\{ 4\sqrt{2\log(2|\mathcal{F}|/\delta)} Q\left(\max_{f \in \mathcal{F}} \sqrt{\mathbb{E}_n[f(Z_i)^2]}, 1 - \tau\right) \right\} \vee 2 \max_{f \in \mathcal{F}} Q\left(|\mathbb{G}_n(f(Z_i))|, \frac{1}{2}\right).$$

Proof. Let

$$e_{1n} = \sqrt{2\log(2|\mathcal{F}|/\delta)} Q\left(\max_{f \in \mathcal{F}} \sqrt{\mathbb{E}_n[f(Z_i)^2]}, 1 - \tau\right), \quad e_{2n} = \max_{f \in \mathcal{F}} Q\left(|\mathbb{G}_n(f(Z_i))|, \frac{1}{2}\right)$$

and the event $\mathcal{E} = \{\max_{f \in \mathcal{F}} \sqrt{\mathbb{E}_n[f^2(Z_i)]} \leq Q(\max_{f \in \mathcal{F}} \sqrt{\mathbb{E}_n[f^2(Z_i)]}, 1 - \tau)\}$ which satisfies $P(\mathcal{E}) \geq 1 - \tau$. By the symmetrization Lemma 2.3.7 of van der Vaart and Wellner (1996) (by definition of e_{2n} we have $\beta_n(x) \geq 1/2$ in Lemma 2.3.7) we obtain

$$\begin{aligned} \mathbb{P}\{\max_{f \in \mathcal{F}} |\mathbb{G}_n(f(Z_i))| > 4e_{1n} \vee 2e_{2n}\} &\leq 4\mathbb{P}\{\max_{f \in \mathcal{F}} |\mathbb{G}_n(\varepsilon_i f(Z_i))| > e_{1n}\} \\ &\leq 4\mathbb{P}\{\max_{f \in \mathcal{F}} |\mathbb{G}_n(\varepsilon_i f(Z_i))| > e_{1n} | \mathcal{E}\} + 4\tau \end{aligned}$$

where ε_i are independent Rademacher random variables, $P(\varepsilon_i = 1) = P(\varepsilon_i = -1) = 1/2$.

Thus a union bound yields

$$\mathbb{P}\left\{\max_{f \in \mathcal{F}} |\mathbb{G}_n(f(Z_i))| > 4e_{1n} \vee 2e_{2n}\right\} \leq 4\tau + 4|\mathcal{F}| \max_{f \in \mathcal{F}} \mathbb{P}\{|\mathbb{G}_n(\varepsilon_i f(Z_i))| > e_{1n} | \mathcal{E}\}. \quad (\text{I.73})$$

We then condition on the values of Z_1, \dots, Z_n and \mathcal{E} , denoting the conditional probability measure as \mathbb{P}_ε . Conditional on Z_1, \dots, Z_n , by the Hoeffding inequality the symmetrized process $\mathbb{G}_n(\varepsilon_i f(Z_i))$ is sub-Gaussian for the $L_2(\mathbb{P}_n)$ norm, namely, for $f \in \mathcal{F}$, $\mathbb{P}_\varepsilon\{|\mathbb{G}_n(\varepsilon_i f(Z_i))| > x\} \leq 2\exp(-x^2/\{2\mathbb{E}_n[f^2(Z_i)]\})$. Hence, under the event \mathcal{E} , we can bound

$$\begin{aligned} \mathbb{P}_\varepsilon\{|\mathbb{G}_n(\varepsilon_i f(Z_i))| > e_{1n} | Z_1, \dots, Z_n, \mathcal{E}\} &\leq 2\exp(-e_{1n}^2/[2\mathbb{E}_n[f^2(Z_i)]]) \\ &\leq 2\exp(-\log(2|\mathcal{F}|/\delta)). \end{aligned}$$

Taking the expectation over Z_1, \dots, Z_n does not affect the right hand side bound. Plugging in this bound yields the result. \square

The following specialization will be convenient.

Lemma 7. *Let $\tau \in (0, 1)$ and $\{(x'_{ij}, \epsilon_i)' \in \mathbb{R}^p \times \mathbb{R}, i = 1, \dots, n\}$ be random vectors that are independent across i . Then with probability at least $1 - 8\tau$*

$$\max_{1 \leq j \leq p} |\mathbb{E}_n[x_{ij}^2 \epsilon_i^2] - \bar{\mathbb{E}}[x_{ij}^2 \epsilon_i^2]| \leq 4\sqrt{\frac{2\log(2p/\tau)}{n}} Q\left(\max_{1 \leq j \leq p} \mathbb{E}_n[x_{ij}^4 \epsilon_i^4], 1 - \tau\right) \vee 2 \max_{1 \leq j \leq p} \sqrt{\frac{2\bar{\mathbb{E}}[x_{ij}^4 \epsilon_i^4]}{n}}$$

Proof. Let $Z_i = x_i \epsilon_i$, $f_j(Z_i) = x_{ij}^2 \epsilon_i^2$, $\mathcal{F} = \{f_1, \dots, f_p\}$, so that $n^{-1/2} \mathbb{G}_n(f_j(Z_i)) = \mathbb{E}_n[x_{ij}^2 \epsilon_i^2] - \bar{\mathbb{E}}[x_{ij}^2 \epsilon_i^2]$. Also, for $\tau_1 \in (0, 1/2)$ and $\tau_2 \in (0, 1)$, let

$$e_{1n} = \sqrt{2 \log(2p/\tau_1)} \sqrt{Q \left(\max_{1 \leq j \leq p} \mathbb{E}_n[x_{ij}^4 \epsilon_i^4], 1 - \tau_2 \right)} \quad \text{and} \quad e_{2n} = \max_{1 \leq j \leq p} Q(|\mathbb{G}_n(x_{ij}^2 \epsilon_i^2)|, 1/2)$$

where we have $e_{2n} \leq \max_{1 \leq j \leq p} \sqrt{2 \bar{\mathbb{E}}[x_{ij}^4 \epsilon_i^4]}$ by Chebyshev.

By Lemma 6 we have

$$P \left(\max_{1 \leq j \leq p} |\mathbb{E}_n[x_{ij}^2 \epsilon_i^2] - \bar{\mathbb{E}}[x_{ij}^2 \epsilon_i^2]| > \frac{4e_{1n} \vee 2e_{2n}}{\sqrt{n}} \right) \leq 4\tau_1 + 4\tau_2.$$

The result follows by setting $\tau_1 = \tau_2 = \tau < 1/2$. Note that for $\tau \geq 1/2$ the result is trivial. \square

I.2. Moment Inequality. We shall be using the following result, which is based on Markov inequality and (von Bahr and Esseen, 1965).

Lemma 8 (Vonbahr-Esseen's LLN). *Let $r \in [1, 2]$, and independent zero-mean random variables X_i with $\bar{\mathbb{E}}[|X_i|^r] \leq C$. Then for any $\ell_n > 0$*

$$Pr \left(\left| \frac{\sum_{i=1}^n X_i}{n} \right| > \ell_n n^{-(1-1/r)} \right) \leq \frac{2C}{\ell_n^r}.$$

I.3. Matrices Deviation Bounds. In this section we collect matrices deviation bounds. We begin with a bound due to Rudelson (1999) for the case that $p < n$.

Lemma 9 (Essentially in Rudelson (1999)). *Let x_i , $i = 1, \dots, n$, be independent random vectors in \mathbb{R}^p and set*

$$\delta_n := \bar{C} \frac{\sqrt{\log(n \wedge p)}}{\sqrt{n}} \sqrt{\mathbb{E}[\max_{1 \leq i \leq n} \|x_i\|^2]}.$$

for some universal constant \bar{C} . Then, we have

$$\mathbb{E} \left[\sup_{\|\alpha\|=1} |\mathbb{E}_n[(\alpha' x_i)^2] - \mathbb{E}[(\alpha' x_i)^2]| \right] \leq \delta_n^2 + \delta_n \sup_{\|\alpha\|=1} \sqrt{\bar{\mathbb{E}}[(\alpha' x_i)^2]}.$$

Based on results in Rudelson and Vershynin (2008), the following lemma for bounded regressors was derived in the supplementary material of Belloni and Chernozhukov (2013)

Lemma 10 (Essentially in Theorem 3.6 of Rudelson and Vershynin (2008)). *Let x_i , $i = 1, \dots, n$, be independent random vectors in \mathbb{R}^p be such that $\sqrt{\mathbb{E}[\max_{1 \leq i \leq n} \|x_i\|_\infty^2]} \leq K$. Let*

$$\delta_n := 2 \left(\bar{C} K \sqrt{k} \log(1+k) \sqrt{\log(p \vee n)} \sqrt{\log n} \right) / \sqrt{n},$$

where \bar{C} is the universal constant. Then,

$$\mathbb{E} \left[\sup_{\|\alpha\|_0 \leq k, \|\alpha\|=1} |\mathbb{E}_n[(\alpha' x_i)^2] - \mathbb{E}[(\alpha' x_i)^2]| \right] \leq \delta_n^2 + \delta_n \sup_{\|\alpha\|_0 \leq k, \|\alpha\|=1} \sqrt{\bar{\mathbb{E}}[(\alpha' x_i)^2]}.$$

Proof. Let

$$V_k = \sup_{\|\alpha\|_0 \leq k, \|\alpha\|=1} \left| \mathbb{E}_n [(\alpha' x_i)^2 - \mathbb{E}[(\alpha' x_i)^2]] \right|.$$

Then, by a standard symmetrization argument (Guédon and Rudelson (2007), page 804)

$$n\mathbb{E}[V_k] \leq 2\mathbb{E}_x \mathbb{E}_\varepsilon \left[\sup_{\|\alpha\|_0 \leq k, \|\alpha\|=1} \left| \sum_{i=1}^n \varepsilon_i (\alpha' x_i)^2 \right| \right].$$

Letting

$$\phi(k) = \sup_{\|\alpha\|_0 \leq k, \|\alpha\| \leq 1} \mathbb{E}_n[(\alpha' x_i)^2] \quad \text{and} \quad \varphi(k) = \sup_{\|\alpha\|_0 \leq k, \|\alpha\|=1} \bar{\mathbb{E}}[(\alpha' x_i)^2],$$

we have $\phi(k) \leq \varphi(k) + V_k$ and by Lemma 3.8 in Rudelson and Vershynin (2008) to bound the expectation in ε ,

$$\begin{aligned} n\mathbb{E}[V_k] &\leq 2 \left(\bar{C} \sqrt{k} \log(1+k) \sqrt{\log(p \vee n)} \sqrt{\log n} \right) \sqrt{n} \mathbb{E}_x \left[\max_{i \leq n} \|x_i\|_\infty \sqrt{\phi(k)} \right] \\ &\leq 2 \left(\bar{C} \sqrt{k} \log(1+k) \sqrt{\log(p \vee n)} \sqrt{\log n} \right) \sqrt{n} \sqrt{\mathbb{E}_x [\max_{i \leq n} \|x_i\|_\infty^2]} \sqrt{\mathbb{E}_x [\phi(k)]} \\ &\leq 2 \left(\bar{C} K \sqrt{k} \log(1+k) \sqrt{\log(p \vee n)} \sqrt{\log n} \right) \sqrt{n} \sqrt{\varphi(k) + \mathbb{E}[V_k]}. \end{aligned}$$

The result follows by noting that for positive numbers v, A, B , $v \leq A(v+B)^{1/2}$ implies $v \leq A^2 + A\sqrt{B}$. \square

The following result establishes an approximation bound for sub-Gaussian regressors and was developed in Rudelson and Zhou (2011). Recall that a random vector $Z \in \mathbb{R}^p$ is isotropic if $\mathbb{E}[ZZ'] = I$, and it is called ψ_2 with a constant α if for every $w \in \mathbb{R}^p$ we have

$$\|Z'w\|_{\psi_2} := \inf \{t : \mathbb{E}[\exp(-(Z'w)^2/t^2)] \leq 2\} \leq \alpha \|w\|_2.$$

Lemma 11 (Essentially in Theorem 3.2 of Rudelson and Zhou (2011)). *Let Ψ_i , $i = 1, \dots, n$, be i.i.d. isotropic random vectors in \mathbb{R}^p that are ψ_2 with a constant α . Let $x_i = \Sigma^{1/2} \Psi_i$ so that $\Sigma = \mathbb{E}[x_i x_i']$. For $m \leq p$ and $\tau \in (0, 1)$ assume that*

$$n \geq \frac{80m\alpha^4}{\tau^2} \log \left(\frac{12ep}{m\tau} \right).$$

Then with probability at least $1 - 2\exp(-\tau^2 n / 80\alpha^4)$, for all $u \in \mathbb{R}^p$, $\|u\|_0 \leq m$, we have

$$(1 - \tau) \|\Sigma^{1/2} u\|_2 \leq \sqrt{\mathbb{E}_n[(x'_i u)^2]} \leq (1 + \tau) \|\Sigma^{1/2} u\|_2.$$

For example, Lemma 11 covers the case of $x_i \sim N(0, \Sigma)$ by setting $\Psi_i \sim N(0, I)$ which is isotropic and ψ_2 with a constant $\alpha = \sqrt{8/3}$.

REFERENCES

- ABADIE, A., AND G. W. IMBENS (2011): “Bias-corrected matching estimators for average treatment effects,” *J. Bus. Econom. Statist.*, 29(1), 1–11.
- ANDREWS, D., AND X. CHENG (2011): “Maximum Likelihood Estimation and Uniform Inference with Sporadic Identification Failure,” Cowless Foundation Discussion Paper.
- ANDREWS, D., X. CHENG, AND P. GUGGENBERGER (2011): “Generic Results for Establishing the Asymptotic Size of Confidence Sets and Tests,” Cowless Foundation Discussion Paper.
- ANGRIST, J. D., AND J.-S. PISCHKE (2008): *Mostly Harmless Econometrics: An Empiricist’s Companion*. Princeton University Press.
- BACH, F. (2010): “Self-concordant analysis for logistic regression,” *Electronic Journal of Statistics*, 4, 384–414.
- BARANIUK, R., M. DAVENPORT, R. DEVORE, AND M. WAKIN (2008): “A Simple Proof of the Restricted Isometry Property for Random Matrices,” *Constructive Approximation*, 28, 253–263.
- BELLONI, A., D. CHEN, V. CHERNOZHUKOV, AND C. HANSEN (2012): “Sparse Models and Methods for Optimal Instruments with an Application to Eminent Domain,” *Econometrica*, 80, 2369–2429.
- BELLONI, A., AND V. CHERNOZHUKOV (2011a): “ ℓ_1 -penalized quantile regression in high-dimensional sparse models,” *Ann. Statist.*, 39(1), 82–130.
- (2011b): “High Dimensional Sparse Econometric Models: An Introduction,” *Inverse problems and high dimensional estimation - Stats in the Château summer school in econometrics and statistics, 2009, Springer Lecture Notes in Statistics - Proceedings*, pp. 121–156.
- (2013): “Least Squares After Model Selection in High-dimensional Sparse Models,” *Bernoulli*, 19(2), 521–547.
- BELLONI, A., V. CHERNOZHUKOV, AND C. HANSEN (2010): “LASSO Methods for Gaussian Instrumental Variables Models,” arXiv:[math.ST], <http://arxiv.org/abs/1012.1297>.
- (2011): “Inference for High-Dimensional Sparse Econometric Models,” *Advances in Economics and Econometrics. 10th World Congress of Econometric Society*.
- BELLONI, A., V. CHERNOZHUKOV, AND K. KATO (2013): “Uniform Post Selection Inference for LAD Regression Models,” *arXiv preprint arXiv:1304.0282*.
- BELLONI, A., V. CHERNOZHUKOV, AND L. WANG (2010): “Square-Root-LASSO: Pivotal Recovery of Nonparametric Regression Functions via Conic Programming,” *Duke and MIT Working Paper*.
- (2011): “Square-Root-LASSO: Pivotal Recovery of Sparse Signals via Conic Programming,” *Biometrika*, 98(4), 791–806.
- BELLONI, A., V. CHERNOZHUKOV, AND Y. WEI (2013): “Honest Confidence Regions for Logistic Regression with a Large Number of Controls,” *arXiv preprint arXiv:1304.3969*.
- BICKEL, P. J., Y. RITOV, AND A. B. TSYBAKOV (2009): “Simultaneous analysis of Lasso and Dantzig selector,” *Annals of Statistics*, 37(4), 1705–1732.
- CANDÈS, E., AND T. TAO (2007): “The Dantzig selector: statistical estimation when p is much larger than n ,” *Ann. Statist.*, 35(6), 2313–2351.
- CATTANEO, M., M. JANSSON, AND W. NEWEY (2010): “Alternative Asymptotics and the Partially Linear Model with Many Regressors,” *Working Paper*, <http://econ-www.mit.edu/files/6204>.
- CATTANEO, M. D. (2010): “Efficient semiparametric estimation of multi-valued treatment effects under ignorability,” *Journal of Econometrics*, 155(2), 138–154.
- CHEN, X. (2007): “Large Sample Sieve Estimation of Semi-Nonparametric Models,” *Handbook of Econometrics*, 6, 5559–5632.
- CHEN, X., AND D. POUZO (2009): “Efficient estimation of semiparametric conditional moment models with possibly nonsmooth residuals,” *Journal of Econometrics*, 152(1), 46–60.

- (2012): “Estimation of nonparametric conditional moment models with possibly nonsmooth moments,” *Econometrica*, 80(1), 277–322.
- DE LA PEÑA, V. H., T. L. LAI, AND Q.-M. SHAO (2009): *Self-normalized processes*, Probability and its Applications (New York). Springer-Verlag, Berlin, Limit theory and statistical applications.
- DONALD, S. G., AND W. K. NEWBY (1994): “Series estimation of semilinear models,” *J. Multivariate Anal.*, 50(1), 30–40.
- DONOHUE III, J. J., AND S. D. LEVITT (2001): “The Impact of Legalized Abortion on Crime,” *Quarterly Journal of Economics*, 116(2), 379–420.
- (2008): “Measurement Error, Legalized Abortion, and the Decline in Crime: A Response to Foote and Goetz,” *Quarterly Journal of Economics*, 123(1), 425–440.
- DUFLO, E., R. GLENNERSTER, AND M. KREMER (2008): “Using Randomization in Development Economics Research: A Toolkit,” in *Handbook of Development Economics. Volume 4*, ed. by T. P. Schultz, and J. A. Strauss. Elsevier: North-Holland.
- FAN, J., S. GUO, AND N. HAO (2011): “Variance estimation using refitted cross-validation in ultrahigh dimensional regression,” *forthcoming Journal of the Royal Statistical Society: Series B (Statistical Methodology)*.
- FAN, J., AND R. LI (2001): “Variable selection via nonconcave penalized likelihood and its oracle properties,” *Journal of American Statistical Association*, 96(456), 1348–1360.
- FOOTE, C. L., AND C. F. GOETZ (2008): “The Impact of Legalized Abortion on Crime: Comment,” *Quarterly Journal of Economics*, 123(1), 407–423.
- FRANK, I. E., AND J. H. FRIEDMAN (1993): “A Statistical View of Some Chemometrics Regression Tools,” *Technometrics*, 35(2), 109–135.
- GAUTIER, E., AND A. TSYBAKOV (2011): “High-dimensional Instrumental Variables Regression and Confidence Sets,” arXiv:1105.2454v2 [math.ST].
- GUÉDON, O., AND M. RUDELSON (2007): “ L_p -moments of random vectors via majorizing measures,” *Advances in Mathematics*, 208, 798–823.
- HAHN, J. (1998): “On the role of the propensity score in efficient semiparametric estimation of average treatment effects,” *Econometrica*, pp. 315–331.
- HANSEN, B. E. (2005): “Challenges for Econometric Model Selection,” *Econometric Theory*, 21, 60–68.
- HÄRDLE, W., H. LIANG, AND J. GAO (2000): *Partially linear models*, Contributions to Statistics. Physica-Verlag, Heidelberg.
- HECKMAN, J., R. LALONDE, AND J. SMITH (1999): “The economics and econometrics of active labor market programs,” *Handbook of labor economics*, 3, 1865–2097.
- HECKMAN, J. J., H. ICHIMURA, AND P. TODD (1998): “Matching as an econometric evaluation estimator,” *Rev. Econom. Stud.*, 65(2), 261–294.
- HIRANO, K., G. W. IMBENS, AND G. RIDDER (2003): “Efficient estimation of average treatment effects using the estimated propensity score,” *Econometrica*, 71(4), 1161–1189.
- HUANG, J., J. L. HOROWITZ, AND S. MA (2008): “Asymptotic properties of bridge estimators in sparse high-dimensional regression models,” *The Annals of Statistics*, 36(2), 587–613.
- HUANG, J., J. L. HOROWITZ, AND F. WEI (2010): “Variable selection in nonparametric additive models,” *Ann. Statist.*, 38(4), 2282–2313.
- IMBENS, G. W. (2004): “Nonparametric Estimation of Average Treatment Effects Under Exogeneity: A Review,” *The Review of Economics and Statistics*, 86(1), 4–29.
- JING, B.-Y., Q.-M. SHAO, AND Q. WANG (2003): “Self-normalized Cramér-type large deviations for independent random variables,” *Ann. Probab.*, 31(4), 2167–2215.

- KERKYACHARIAN, G., AND D. PICARD (1992): “Density estimation in Besov spaces,” *Statist. Probab. Lett.*, 13(1), 15–24.
- KOENKER, R. (1988): “Asymptotic Theory and Econometric Practice,” *Journal of Applied Econometrics*, 3, 139–147.
- KREMER, M., AND R. GLENNERSTER (2011): “Improving Health in Developing Countries: Evidence from Randomized Evaluations,” in *Handbook of Health Economics. Volume 2*, ed. by T. G. M. M. V. Pauly, and P. P. Barros. Elsevier: North-Holland.
- LEDoux, M., AND M. TALAGRAND (1991): *Probability in Banach Spaces (Isoperimetry and processes)*. Ergebnisse der Mathematik und ihrer Grenzgebiete, Springer-Verlag.
- LEEB, H., AND B. M. PÖTSCHER (2008): “Can one estimate the unconditional distribution of post-model-selection estimators?,” *Econometric Theory*, 24(2), 338–376.
- MACKINNON, J. G., AND H. WHITE (1985): “Some heteroskedasticity consistent covariance matrix estimators with improved finite sample properties,” *Journal of Econometrics*, 29, 305–325.
- MEINSHAUSEN, N., AND B. YU (2009): “Lasso-type recovery of sparse representations for high-dimensional data,” *Annals of Statistics*, 37(1), 2246–2270.
- MIKUSHEVA, A. (2007): “Uniform inference in autoregressive models,” *Econometrica*, 75(5), 1411–1452.
- NEWBY, W. K. (1997): “Convergence Rates and Asymptotic Normality for Series Estimators,” *Journal of Econometrics*, 79, 147–168.
- NEYMAN, J. (1979): “ $C(\alpha)$ tests and their use,” *Sankhya*, 41, 1–21.
- PÖTSCHER, B. (2009a): “Confidence Sets Based on Sparse Estimators Are Necessarily Large,” *Sankhya*, 71-A, 1–18.
- PÖTSCHER, B. M. (2009b): “Confidence sets based on sparse estimators are necessarily large,” *Sankhyā*, 71(1, Ser. A), 1–18.
- ROBINS, J. M., AND A. ROTNITZKY (1995): “Semiparametric efficiency in multivariate regression models with missing data,” *J. Amer. Statist. Assoc.*, 90(429), 122–129.
- ROBINSON, P. M. (1988): “Root- N -consistent semiparametric regression,” *Econometrica*, 56(4), 931–954.
- ROMANO, J. P. (2004): “On non-parametric testing, the uniform behaviour of the t -test, and related problems,” *Scand. J. Statist.*, 31(4), 567–584.
- ROTHER, C., AND S. FIRPO (2013): “Semiparametric Estimation and Inference Using Doubly Robust Moment Conditions,” Discussion paper, NYU preprint.
- RUDELSON, M. (1999): “Random vectors in the isotropic position,” *Journal of Functional Analysis*, 164, 6072.
- RUDELSON, M., AND R. VERSHYNIN (2008): “On sparse reconstruction from Fourier and Gaussian measurements,” *Communications on Pure and Applied Mathematics*, 61, 10251045.
- RUDELSON, M., AND S. ZHOU (2011): “Reconstruction from anisotropic random measurements,” *ArXiv:1106.1151*.
- TIBSHIRANI, R. (1996): “Regression shrinkage and selection via the Lasso,” *J. Roy. Statist. Soc. Ser. B*, 58, 267–288.
- VAN DE GEER, S. A. (2008): “High-dimensional generalized linear models and the lasso,” *Annals of Statistics*, 36(2), 614–645.
- VAN DER VAART, A. W., AND J. A. WELLNER (1996): *Weak Convergence and Empirical Processes*. Springer Series in Statistics.
- VON BAHR, B., AND C.-G. ESSEEN (1965): “Inequalities for the r th absolute moment of a sum of random variables, $1 \leq r \leq 2$,” *Ann. Math. Statist.*, 36, 299–303.
- ZHOU, S. (2009): “Restricted eigenvalue conditions on subgaussian matrices,” *ArXiv:0904.4723v2*.

Table 1. Simulation Results for Selected R^2 Values

Estimation Procedure	First Stage $R^2 = .2$		First Stage $R^2 = .2$		First Stage $R^2 = .8$		First Stage $R^2 = .8$	
	Structure $R^2 = 0$		Structure $R^2 = .8$		Structure $R^2 = 0$		Structure $R^2 = .8$	
	RMSE	Rej. Rate	RMSE	Rej. Rate	RMSE	Rej. Rate	RMSE	Rej. Rate
	A. Design 1. Quadratic Decay							
Oracle	0.090	0.048	0.090	0.048	0.045	0.057	0.045	0.057
Double-Selection Oracle	0.102	0.050	0.102	0.050	0.143	0.047	0.143	0.047
Post-Lasso	0.137	0.205	0.110	0.064	0.402	0.987	0.489	0.974
Double-Selection	0.107	0.063	0.107	0.058	0.109	0.074	0.104	0.062
Double-Selection + Ridge	0.260	0.064	0.256	0.055	0.132	0.049	0.130	0.050
	B. Design 2. Quadratic Decay with Heteroscedasticity							
Oracle	0.139	0.060	0.139	0.060	0.066	0.062	0.066	0.062
Double-Selection Oracle	0.169	0.072	0.169	0.072	0.225	0.085	0.225	0.085
Post-Lasso	0.175	0.139	0.178	0.097	0.409	0.994	0.501	0.993
Double-Selection	0.165	0.098	0.167	0.081	0.162	0.082	0.165	0.083
Double-Selection + Ridge	0.308	0.060	0.290	0.058	0.183	0.064	0.185	0.075
	C. Design 3. Quadratic Decay with Random Coefficients							
Oracle	0.070	0.055	0.070	0.055	0.041	0.060	0.041	0.060
Double-Selection Oracle	0.114	0.056	0.114	0.056	0.151	0.058	0.151	0.058
Post-Lasso	0.105	0.082	0.131	0.133	0.329	0.940	0.435	0.953
Double-Selection	0.109	0.055	0.118	0.075	0.105	0.056	0.117	0.086
Double-Selection + Ridge	0.227	0.040	0.230	0.035	0.151	0.054	0.153	0.057

Note: The table reports root-mean-square-error (RMSE) rejection rates for 5% level tests (Rej. Rate) from a Monte Carlo simulation experiment. Results are based on 1000 simulation replications. Data in Panels A and B are based on models with coefficients that decay quadratically, and the data in Panel C are based on a with five quadratically decaying coefficients and 95 random coefficients. Further details about the simulation models are provided in the text as are details about the estimation procedures. Rejection rates are for t-tests of the null hypothesis that the structural coefficient is equal to the true population value and are formed using jack-knife standard errors that are robust to heteroscedasticity; see MacKinnon and White (1985).

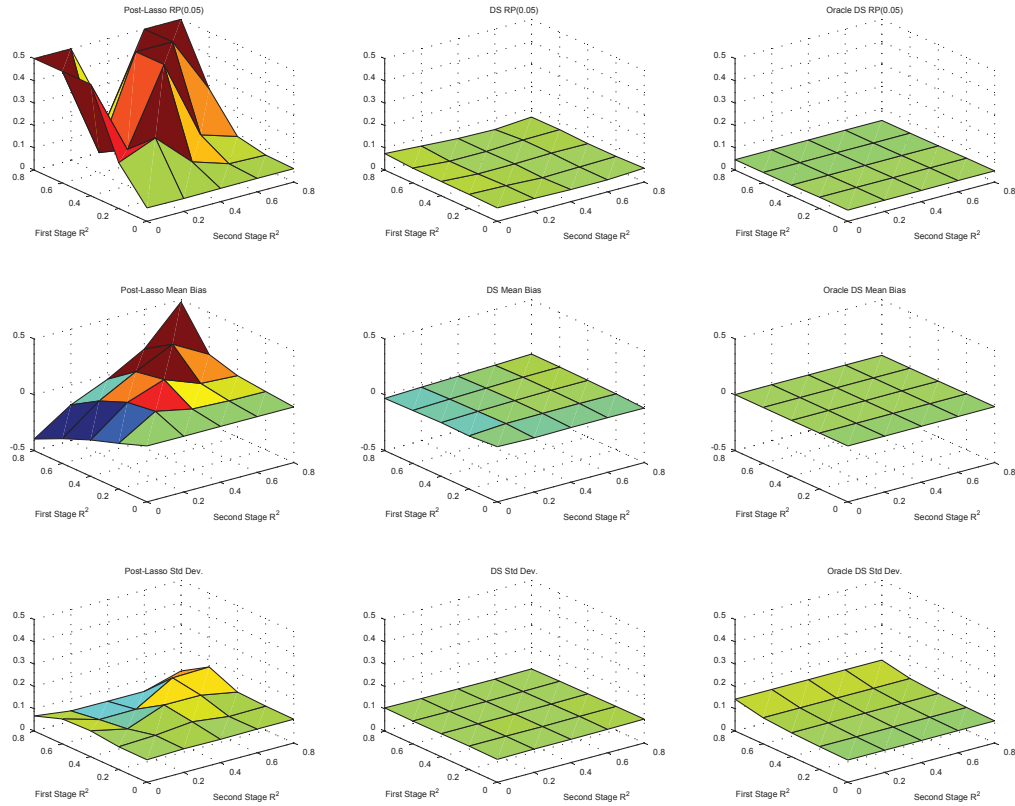


FIGURE 2. This figure presents rejection frequencies for 5% level tests, biases, and standard deviations for estimating the treatment effect from Design 1 of the simulation study which has quadratically decaying coefficients and homoscedasticity. Results are reported for a one-step Post-Lasso estimator, our proposed double selection procedure, and the infeasible OLS estimator that uses the set of variables that have coefficients larger than 0.1 in either equation (2.6) or (2.7). Reduced form and first stage R^2 correspond to the population R^2 of (2.6) and (2.7) respectively. Note that rejection frequencies are censored at 0.5.

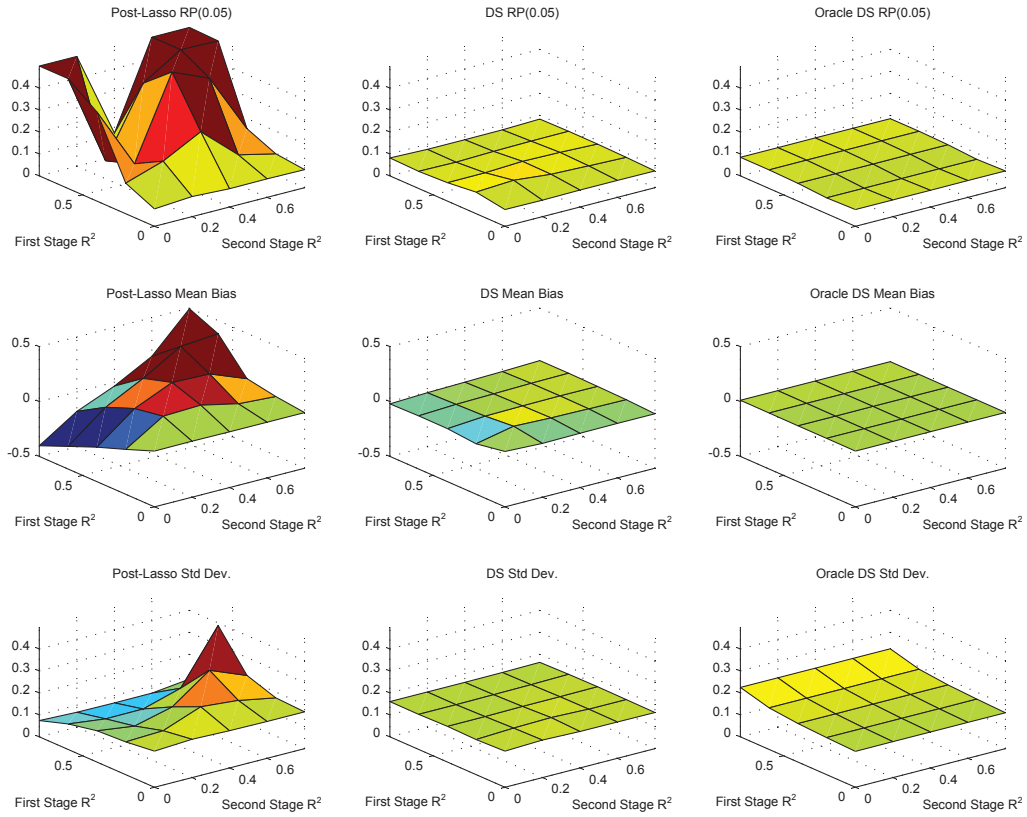


FIGURE 3. This figure presents rejection frequencies for 5% level tests, biases, and standard deviations for estimating the treatment effect from Design 2 of the simulation study which has quadratically decaying coefficients and heteroscedasticity. Results are reported for a one-step Post-Lasso estimator, our proposed double selection procedure, and the infeasible OLS estimator that uses the set of variables that have coefficients larger than 0.1 in either equation (2.6) or (2.7). Reduced form and first stage R^2 correspond to the population R^2 of (2.6) and (2.7) respectively. Note that rejection frequencies are censored at 0.5.

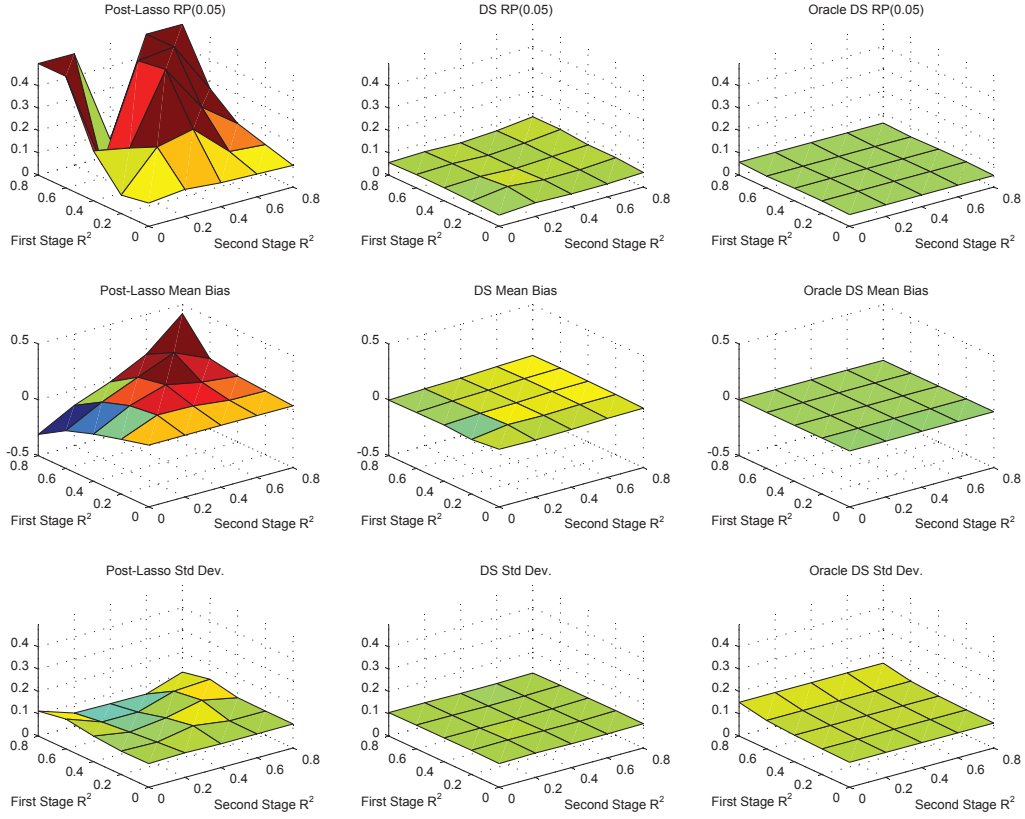


FIGURE 4. This figure presents rejection frequencies for 5% level tests, biases, and standard deviations for estimating the treatment effect from Design 3 of the simulation study which has five quadratically decaying coefficients and 95 Gaussian random coefficients. Results are reported for a one-step Post-Lasso estimator, our proposed double selection procedure, and the infeasible OLS estimator that uses the set of variables that have coefficients larger than 0.1 in either equation (2.6) or (2.7). Reduced form and first stage R^2 correspond to what would be the population R^2 of (2.6) and (2.7) if all of the random coefficients were equal to zero. Note that rejection frequencies are censored at 0.5.

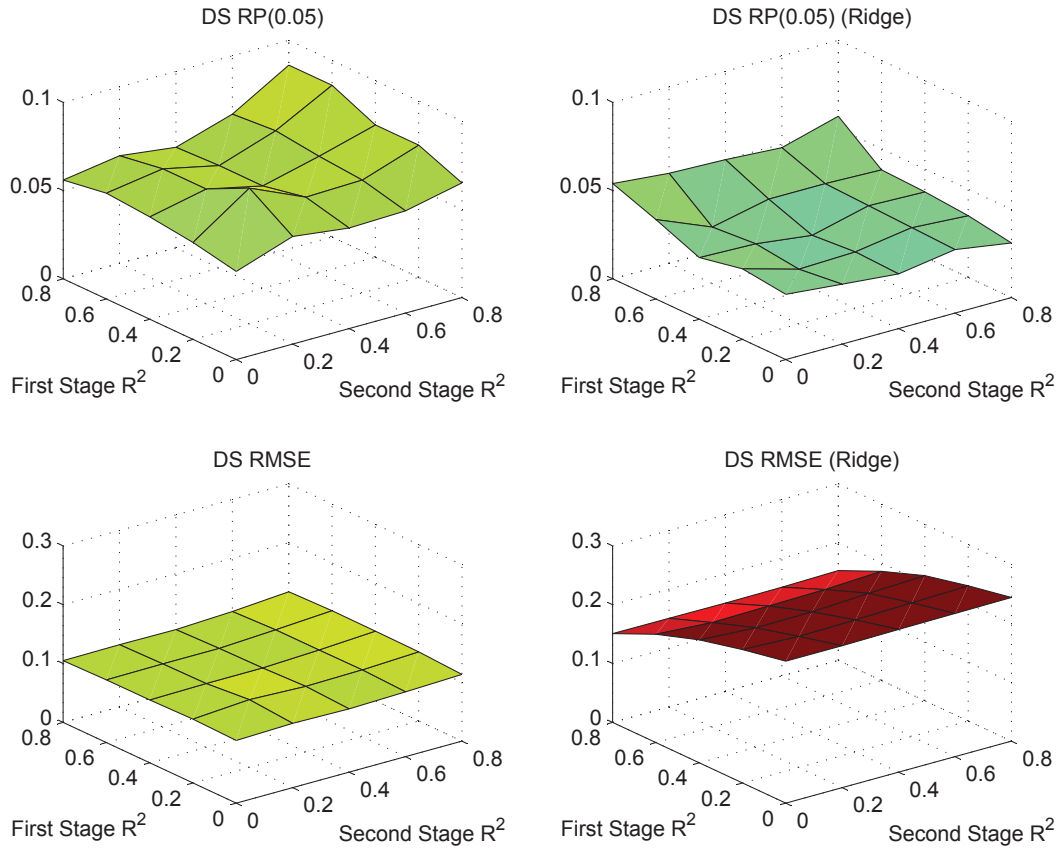


FIGURE 5. This figure presents rejection frequencies for 5% level tests and RMSE's for estimating the treatment effect from Design 3 of the simulation study which has five quadratically decaying coefficients and 95 Gaussian random coefficients. Results in the first column are for the proposed double selection procedure, and the results in the second column are for the proposed double selection procedure when the ridge fit from (2.6) is added as an additional potential control. Reduced form and first stage R^2 correspond to what would be the population R^2 of (2.6) and (2.7) if all of the random coefficients were equal to zero. Note that the vertical axis on the rejection frequency graph is from 0 to 0.1.

Table 2. Estimated Effects of Abortion on Crime Rates

	Violent Crime		Property Crime		Murder	
	Effect	Std. Err.	Effect	Std. Err.	Effect	Std. Err.
A. Donohue III and Levitt (2001) Table IV						
Donohue III and Levitt (2001) Table IV	-0.129	0.024	-0.091	0.018	-0.121	0.047
First-Difference	-0.152	0.034	-0.108	0.022	-0.204	0.068
All Controls	0.014	0.719	-0.195	0.225	2.343	2.796
Post-Double-Selection	-0.176	0.110	-0.034	0.042	0.012	0.165
Post-Double-Selection+	-0.157	0.111	-0.036	0.042	0.046	0.216

Note: The table displays the estimated coefficient on the abortion rate, "Effect," and its estimated standard error. Numbers in the first row are taken from Donohue III and Levitt (2001) Table IV, columns (2), (4), and (6). The remaining rows are estimated by first differences, include a full set of time dummies, and use standard errors clustered at the state-level. Estimates in the row labeled "First-Difference" are obtained using the same controls as in the first row. Estimates in the row labeled "All Controls" use 284 control variables as discussed in the text. Estimates in the row "Post-Double-Selection" use the variable selection technique developed in this paper to search among the set of 284 potential controls. Estimates in the row "Post-Double-Selection+" use the variables selected by the procedure of this paper augmented with the set of variables from Donohue III and Levitt (2001).