

# Uniform post selection inference for LAD regression and other z-estimation problems

---

Alexandre Belloni  
Victor Chernozhukov  
Kengo Kato

The Institute for Fiscal Studies  
Department of Economics, UCL

**cemmap** working paper CWP74/13

# UNIFORM POST SELECTION INFERENCE FOR LAD REGRESSION AND OTHER Z-ESTIMATION PROBLEMS

A. BELLONI, V. CHERNOZHUKOV, AND K. KATO

ABSTRACT. We develop uniformly valid confidence regions for regression coefficients in a high-dimensional sparse least absolute deviation/median regression model. The setting is one where the number of regressors  $p$  could be large in comparison to the sample size  $n$ , but only  $s \ll n$  of them are needed to accurately describe the regression function. Our new methods are based on the instrumental median regression estimator that assembles the optimal estimating equation from the output of the post  $\ell_1$ -penalized median regression and post  $\ell_1$ -penalized least squares in an auxiliary equation. The estimating equation is immunized against non-regular estimation of nuisance part of the median regression function, in the sense of Neyman. We establish that in a homoscedastic regression model, the instrumental median regression estimator of a single regression coefficient is asymptotically root- $n$  normal uniformly with respect to the underlying sparse model. The resulting confidence regions are valid uniformly with respect to the underlying model. We illustrate the value of uniformity with Monte-Carlo experiments which demonstrate that standard/naive post-selection inference breaks down over large parts of the parameter space, and the proposed method does not. We then generalize our method to the case where  $p_1 \gg n$  regression coefficients are of interest in a non-smooth Z-estimation framework with approximately sparse nuisance functions, containing median regression with a single target regression coefficient as a very special case. We construct simultaneous confidence bands on all  $p_1$  coefficients, and establish their uniform validity over the underlying approximately sparse model.

Key words: uniformly valid inference, instruments, Neymanization, optimality, sparsity, model selection

## 1. INTRODUCTION

We consider the following regression model

$$y_i = d_i \alpha_0 + x_i' \beta_0 + \epsilon_i, \quad i = 1, \dots, n, \quad (1.1)$$

where  $d_i$  is the “main regressor” of interest, whose coefficient  $\alpha_0$  we would like to estimate and perform (robust) inference on. The  $(x_i)_{i=1}^n$  are other high-dimensional regressors or “controls.” The regression error  $\epsilon_i$  is independent of  $d_i$  and  $x_i$  and has median 0. The errors  $(\epsilon_i)_{i=1}^n$  are i.i.d. with distribution

---

*Date:* First version: May 2012, this version December 30, 2013. We would like to thank the participants of Luminy conference on Nonparametric and high-dimensional statistics (December 2012), Oberwolfach workshop on Frontiers in Quantile Regression (November 2012), 8th World Congress in Probability and Statistics (August 2012), and seminar at the University of Michigan (October 2012). We are grateful to Sara van de Geer, Xuming He, Richard Nickl, Roger Koenker, Vladimir Koltchinskii, Steve Portnoy, Philippe Rigollet, and Bin Yu for useful comments and discussions.

function  $F(\cdot)$  and probability density function  $f_\epsilon(\cdot)$  such that  $F(0) = 1/2$  and  $f_\epsilon = f_\epsilon(0) > 0$ . The assumption on the error term motivates the use of the least absolute deviation (LAD) or median regression, suitably adjusted for use in high-dimensional settings.

The dimension  $p$  of “controls”  $x_i$  is large, potentially much larger than  $n$ , which creates a challenge for inference on  $\alpha_0$ . Although the unknown true parameter  $\beta_0$  lies in this large space, the key assumption that will make estimation possible is its sparsity, namely  $T = \text{supp}(\beta_0)$  has  $s < n$  elements (where  $s$  can depend on  $n$ ; we shall use array asymptotics). This in turn motivates the use of regularization or model selection methods.

A standard (non-robust) approach towards inference in this setting would be first to perform model selection via the  $\ell_1$ -penalized LAD regression estimator

$$(\hat{\alpha}, \hat{\beta}) \in \arg \min_{\alpha, \beta} \mathbb{E}_n[|y - d\alpha - x'\beta|] + \frac{\lambda}{n} \|\Psi(\alpha, \beta)'\|_1, \quad (1.2)$$

where  $\lambda$  is a penalty parameter and  $\Psi^2 = \text{diag}(\mathbb{E}_n[d^2], \mathbb{E}_n[x_1^2], \dots, \mathbb{E}_n[x_p^2])$  is a diagonal matrix with normalization weights, and then to use the post-model selection estimator

$$(\tilde{\alpha}, \tilde{\beta}) \in \arg \min_{\alpha, \beta} \left\{ \mathbb{E}_n[|y - d\alpha - x'\beta|] : \beta_j = 0 \text{ if } \hat{\beta}_j = 0 \right\} \quad (1.3)$$

to perform “usual” inference for  $\alpha_0$ . (The notation  $\mathbb{E}_n[\cdot]$  denotes the average over index  $1 \leq i \leq n$ .)

This standard approach is justified if (1.2) achieves perfect model selection with probability approaching 1, so that the estimator (1.3) has the “oracle” property with probability approaching 1. However conditions for “perfect selection” are very restrictive in this model, in particular, requiring significant separation of non-zero coefficients away from zero. If these conditions do not hold, the estimator  $\tilde{\alpha}$  does not converge to  $\alpha_0$  at the  $\sqrt{n}$ -rate – uniformly with respect to the underlying model – which implies that “usual” inference breaks down and is not valid. (The statements continue to apply if  $\alpha$  is not penalized in (1.2),  $\alpha$  is restricted in (1.3), or if thresholding is applied.) We shall demonstrate the breakdown of such naive inference in the Monte-Carlo experiments where non-zero coefficients in  $\theta_0$  are not significantly separated from zero.

Note that the breakdown of inference does not mean that the aforementioned procedures are not suitable for prediction purposes. Indeed, the  $\ell_1$ -LAD estimator (1.2) and post  $\ell_1$ -LAD estimator (1.3) attain (essentially) optimal rates  $\sqrt{(s \log p)/n}$  of convergence for estimating the entire median regression function, as has been shown in [28, 3, 15, 31] and in [3]. This property means that while these procedures will not deliver perfect model recovery, they will only make “moderate” model selection mistakes (omitting only controls with coefficients local to zero).

To achieve uniformly valid inferential performance we propose a procedure whose performance does not require perfect model selection and allows potential “moderate” model selection mistakes. The latter feature is critical in achieving uniformity over a large class of data generating processes, similarly to the results for instrumental regression and mean regression studied in [2], [7], [32], [6]. This allows us to overcome the impact of (moderate) model selection mistakes on inference, avoiding (in part) the criticisms in [19], who prove that the “oracle property” sometime achieved by the naive estimators necessarily implies the failure of uniform validity of inference and their semiparametric inefficiency [20].

In order to achieve robustness with respect to moderate model selection mistakes, it will be necessary to achieve the proper orthogonality condition between the main regressors and the control variables. Towards that goal the following auxiliary equation plays a key role (in the homoscedastic case):

$$d_i = x_i' \theta_0 + v_i, \quad \mathbb{E}[v_i \mid x_i] = 0, \quad i = 1, \dots, n; \quad (1.4)$$

describing the relevant dependence of the regressor of interest  $d_i$  to the other controls  $x_i$ . We shall assume the sparsity of  $\theta_0$ , namely  $T_d = \text{supp}(\theta_0)$  has at most  $s < n$  elements, and estimate the relation (1.4) via Lasso or post-Lasso methods described below.

Given  $v_i$ , which “partials out” the effect of  $x_i$  from  $d_i$ , we shall use it as an instrument in the following estimating equations for  $\alpha_0$ :

$$\mathbb{E}[\varphi(y_i - d_i \alpha_0 - x_i' \beta_0) v_i] = 0, \quad i = 1, \dots, n,$$

where  $\varphi(t) = 1/2 - 1(t < 1/2)$ . We shall use the empirical analog of this equation to form an instrumental median regression estimator of  $\alpha_0$ , using a plug-in estimator for  $x_i' \beta_0$ . The estimating equation above has the following feature:

$$\left. \frac{\partial}{\partial \beta} \mathbb{E}[\varphi(y_i - d_i \alpha_0 - x_i' \beta) v_i] \right|_{\beta = \beta_0} = 0, \quad i = 1, \dots, n, \quad (1.5)$$

As a result, the estimator of  $\alpha_0$  will be “immunized” against “crude” estimation of  $x_i' \beta_0$ , for example, via a post-selection procedure or some regularization procedure. As we explain in Section 5, such immunization ideas can be traced back to Neyman ([21, 22]).

Our estimation procedure has the following three steps.

- Step 1: Estimation of the confounding function  $x_i' \beta_0$  in (1.1).
- Step 2: Estimation of the instruments (residuals)  $v_i$  in (1.4).
- Step 3: Estimation of the main effect  $\alpha_0$  based on the instrumental median regression using  $v_i$  as instruments for  $d_i$ .

Each step is computationally tractable, involving solutions of convex problems and a one-dimensional search, and relies on a different identification condition which in turn requires a different estimation procedure:

Step 1 constructs an estimate for the nuisance function  $x_i' \beta_0$  and not an estimate for  $\alpha_0$ . Here we do not need a  $\sqrt{n}$ -rate consistency for the estimates of the nuisance function; slower rate like  $o(n^{-1/4})$  will suffice. Thus, this can be based either on the  $\ell_1$ -LAD regression estimator (1.2) or the associated post-model selection estimator (1.3).

Step 2 partials out the impact of the covariates  $x_i$  on the main regressor  $d_i$ , obtaining the estimate of the residuals  $v_i$  in the decomposition (1.4). In order to estimate these residuals we rely either on heteroscedastic Lasso [2], a version of the Lasso estimator of [27, 9]:

$$\hat{\theta} \in \arg \min_{\theta} \mathbb{E}_n[(d - x' \theta)^2] + \frac{\lambda}{n} \|\hat{\Gamma} \theta\|_1 \quad \text{and set } \hat{v}_i = d_i - x_i' \hat{\theta}, \quad i = 1, \dots, n, \quad (1.6)$$

where  $\lambda$  and  $\widehat{\Gamma}$  are the penalty level and data-driven penalty loadings described in [2] (restated in Appendix D), or the associated post-model selection estimator (Post-Lasso) [4, 2] defined as

$$\widetilde{\theta} \in \arg \min_{\theta} \left\{ \mathbb{E}_n[(d - x'\theta)^2] : \theta_j = 0 \text{ if } \widehat{\theta}_j = 0 \right\} \text{ and set } \widehat{v}_i = d_i - x_i'\widetilde{\theta}. \quad (1.7)$$

Step 3 constructs an estimator  $\check{\alpha}$  of the coefficient  $\alpha_0$  via an instrumental LAD regression proposed in [12], using  $(\widehat{v}_i)_{i=1}^n$  as instruments. Formally,  $\check{\alpha}$  is defined as

$$\check{\alpha} \in \arg \inf_{\alpha \in \mathcal{A}} L_n(\alpha), \text{ where } L_n(\alpha) = \frac{4|\mathbb{E}_n[\varphi(y - x'\widehat{\beta} - d\alpha)\widehat{v}]|^2}{\mathbb{E}_n[\widehat{v}^2]}, \quad (1.8)$$

$\varphi(t) = 1/2 - 1\{t \leq 0\}$  and  $\mathcal{A}$  is a parameter space for  $\alpha_0$ . We will analyze the choice of  $\mathcal{A} = [\widehat{\alpha} - C \log^{-1} n, \widehat{\alpha} + C \log^{-1} n]$  with a suitable constant  $C > 0$ .<sup>1</sup> Several other choices for  $\mathcal{A}$  are possible.

Our main result establishes conditions under which  $\check{\alpha}$  is root- $n$  consistent for  $\alpha_0$ , asymptotically normal, and achieves the semi-parametric efficiency bound for estimating  $\alpha_0$  in the current homoscedastic setting, provided that  $(s^3 \log^3 p)/n \rightarrow 0$  and other regularity conditions hold. Specifically, we show that, despite possible model selection mistakes in Steps 1 and 2, the estimator  $\check{\alpha}$  obeys

$$\sigma_n^{-1} \sqrt{n}(\check{\alpha} - \alpha_0) \rightsquigarrow N(0, 1), \quad (1.9)$$

where  $\sigma_n^2 := 1/(4f_\epsilon^2 \mathbb{E}[v^2])$  with  $f_\epsilon = f_\epsilon(0)$ . An alternative (and more robust) expression for  $\sigma_n^2$  is given by Huber's sandwich:

$$\sigma_n^2 = J^{-1} \Omega J^{-1}, \text{ where } \Omega := \mathbb{E}[v^2]/4 \text{ and } J := \mathbb{E}[f_\epsilon dv]. \quad (1.10)$$

We recommend to estimate  $\Omega$  by the plug-in method and to estimate  $J$  by Powell's method [23]. Furthermore, we show that the criterion function at the true value  $\alpha_0$  in Step 3 has the following pivotal behavior

$$nL_n(\alpha_0) \rightsquigarrow \chi^2(1). \quad (1.11)$$

This allows the construction of a confidence region  $\widehat{A}_{n,\xi}$  with asymptotic coverage  $1 - \xi$  based on the statistic  $L_n$ ,

$$\text{pr}(\alpha_0 \in \widehat{A}_{n,\xi}) \rightarrow 1 - \xi \text{ where } \widehat{A}_{n,\xi} = \{\alpha \in \mathcal{A} : nL_n(\alpha) \leq (1 - \xi)\text{-quantile of } \chi^2(1)\}. \quad (1.12)$$

Importantly, the robustness with respect to moderate model selection mistakes, which occurs because of (1.5), allows the results (1.9) and (1.11) to hold uniformly over a large range of data generating processes, similarly to the results for instrumental regression and partially linear mean regression model established in [6, 32, 2]. One of our proposed algorithms explicitly uses  $\ell_1$ -regularization methods, similarly to [32] and [2], while the main algorithm we propose uses post-selection methods, similarly to [6, 2].

Throughout the paper, we use array asymptotics – asymptotics where the model changes with  $n$  – to better capture some finite-sample phenomena such as “small coefficients” that are local to zero. This ensures the robustness of conclusions with respect to perturbations of the data-generating process along various model sequences. This robustness, in turn, translates into uniform validity of confidence regions over substantial regions of data-generating processes.

---

<sup>1</sup>For numerical experiments we used  $C = 10(\mathbb{E}_n[d^2])^{-1/2}$  and typically we normalize  $\mathbb{E}_n[d^2] = 1$ .

In Section 3 we generalize the LAD regression to a more general setting by (i) allowing  $p_1$  target parameters defined via Huber's Z-problems are of interest, with dimension  $p_1$  potentially much larger than the sample size, and (ii) also allowing for approximately sparse models. This framework covers many other semi-parametric models since we cover smooth and non-smooth score functions. We provide sufficient conditions to derive a uniform Bahadur representation. Finally, building up on [11], we verify the validity of a multiplier bootstrap procedure.

**1.1. Notation and convention.** Denote by  $(\Omega, \mathcal{F}, \text{pr})$  the underlying probability space. The notation  $\mathbb{E}_n[\cdot]$  denotes the average over index  $1 \leq i \leq n$ , i.e., it simply abbreviates the notation  $n^{-1} \sum_{i=1}^n [\cdot]$ . For example,  $\mathbb{E}_n[x_j^2] = n^{-1} \sum_{i=1}^n x_{ij}^2$ . For a function  $f : \mathbb{R} \times \mathbb{R} \times \mathbb{R}^p \rightarrow \mathbb{R}$ , we write  $\mathbb{G}_n(f) = n^{-1/2} \sum_{i=1}^n (f(y_i, d_i, x_i) - \mathbb{E}[f(y_i, d_i, x_i)])$ . 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. Denote by  $\|\cdot\|_\infty$  the maximal absolute element of a vector. For a sequence  $(z_i)_{i=1}^n$  of constants, we write  $\|z\|_{2,n} = \sqrt{\mathbb{E}_n[z^2]}$ . For example, for a vector  $\delta \in \mathbb{R}^p$ ,  $\|x'\delta\|_{2,n} = \sqrt{\mathbb{E}_n[(x'\delta)^2]}$  denotes the prediction norm of  $\delta$ . 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 such that  $(\delta_T)_j = \delta_j$  if  $j \in T$  and  $(\delta_T)_j = 0$  if  $j \notin T$ . Also we write the support of  $\delta$  as  $\text{supp}(\delta) = \{j \in \{1, \dots, p\} : \delta_j \neq 0\}$ . 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  $a = O_P(b)$ . The arrow  $\rightsquigarrow$  denotes convergence in distribution. We assume that the quantities such as  $p$  (the dimension of  $x_i$ ),  $s$  (a bound on the numbers of non-zero elements of  $\beta_0$  and  $\theta_0$ ), and hence  $y_i, x_i, \beta_0, \theta_0, T$  and  $T_d$  are all dependent on the sample size  $n$ , and allow for the case where  $p = p_n \rightarrow \infty$  and  $s = s_n \rightarrow \infty$  as  $n \rightarrow \infty$ . However, for the notational convenience, we shall omit the dependence of these quantities on  $n$ . For a class of measurable functions  $\mathcal{F}$  equipped with the envelope  $F = \sup_{f \in \mathcal{F}} |f|$ , let  $N(\epsilon, \mathcal{F}, \|\cdot\|_{Q,2})$  denote the  $\epsilon$ -covering number of the class of functions  $\mathcal{F}$  with respect to the  $L^2(Q)$  seminorm  $\|\cdot\|_{Q,2}$ , where  $Q$  is finitely discrete, and let  $\text{ent}(\epsilon, \mathcal{F}) = \log \sup_Q N(\epsilon \|F\|_{Q,2}, \mathcal{F}, \|\cdot\|_{Q,2})$  denote the uniform covering entropy.

## 2. THE METHODS, CONDITIONS, AND RESULTS

**2.1. The methods.** Each of the steps outlined before uses a different identification condition. Several combinations are possible to implement each step, two of which are the following.

**Algorithm 1 (Based on Post-Model Selection estimators).**

- (1) Run Post- $\ell_1$ -penalized LAD (1.3) of  $y_i$  on  $d_i$  and  $x_i$ ; keep fitted value  $x_i' \tilde{\beta}$ .
- (2) Run Post-Lasso (1.7) of  $d_i$  on  $x_i$ ; keep the residual  $\hat{v}_i := d_i - x_i' \tilde{\theta}$ .
- (3) Run Instrumental LAD regression (1.8) of  $y_i - x_i' \tilde{\beta}$  on  $d_i$  using  $\hat{v}_i$  as the instrument for  $d_i$  to compute the estimator  $\tilde{\alpha}$ . Report  $\tilde{\alpha}$  and/or perform inference based upon (1.9) or (1.12).

**Algorithm 2 (Based on Regularized Estimators).**

- (1) Run  $\ell_1$ -penalized LAD (1.2) of  $y_i$  on  $d_i$  and  $x_i$ ; keep fitted value  $x_i' \hat{\beta}$ .
- (2) Run Lasso of (1.6)  $d_i$  on  $x_i$ ; keep the residual  $\hat{v}_i := d_i - x_i' \hat{\theta}$ .
- (3) Run Instrumental LAD regression (1.8) of  $y_i - x_i' \hat{\beta}$  on  $d_i$  using  $\hat{v}_i$  as the instrument for  $d_i$  to compute the estimator  $\tilde{\alpha}$ . Report  $\tilde{\alpha}$  and/or perform inference based upon (1.9) or (1.12).

**Comment 2.1** (Penalty Levels). In order to perform  $\ell_1$ -LAD and Lasso, one has to suitably choose the penalty levels. In the Supplementary Appendix F we provide implementation details including penalty choices for each step of the algorithm, and in all what follows we shall obey the penalty choices described in Appendix F.

**Comment 2.2** (Differences). Algorithm 1 relies on Post- $\ell_1$ -LAD and Post-Lasso while Algorithm 2 relies on  $\ell_1$ -LAD and Lasso. Since Algorithm 1 refits the non-zero coefficients without the penalty term it has a smaller bias. Therefore it does rely on  $\ell_1$ -LAD and Lasso obtaining sparse solutions which in turn typically relies on restricted isometry conditions [3, 2]. Algorithm 2 relies on penalized estimators. Step 3 of both algorithms relies on instrumental LAD regression with estimated data.

**Comment 2.3** (Alternative Implementations). As discussed before, the three step approach proposed here can be implemented with several different methods each with specific features. For instance, Dantzig selector, square-root Lasso or the associated post-model selection could be used instead of Lasso or Post-Lasso. Moreover, the instrumental LAD regression can be substituted by a 1-step estimator from the  $\ell_1$ -LAD estimator  $\hat{\alpha}$  of the form  $\check{\alpha} = \hat{\alpha} + (\mathbb{E}_n[f_\epsilon \hat{v}^2])^{-1} \mathbb{E}_n[\varphi(y - d\hat{\alpha} - x'\hat{\beta})\hat{v}]$  or by a LAD regression with all the covariates selected in Steps 1 and 2.

**2.2. Regularity Conditions.** Here we provide regularity conditions that are sufficient for validity of the main estimation and inference results. The behavior of the population Gram matrix  $\mathbb{E}[\tilde{x}\tilde{x}']$ , where  $\tilde{x}_i = (d_i, x_i)'$ , plays an important role in the analysis. It suffices to have good behavior of smaller submatrices despite the fact that whenever  $p + 1 > n$ , the empirical Gram matrix  $\mathbb{E}_n[\tilde{x}\tilde{x}']$  does not have full rank and in principle is not well-behaved. Define the minimal and maximal  $m$ -sparse eigenvalue of a matrix  $M$  as

$$\phi_{\min}(m, M) := \min_{1 \leq \|\delta\|_0 \leq m} \frac{\delta' M \delta}{\|\delta\|^2} \quad \text{and} \quad \phi_{\max}(m, M) := \max_{1 \leq \|\delta\|_0 \leq m} \frac{\delta' M \delta}{\|\delta\|^2}. \quad (2.13)$$

We denote  $\bar{\phi}_{\min}(m) := \phi_{\max}(m, \mathbb{E}[\tilde{x}\tilde{x}'])$  and  $\bar{\phi}_{\max}(m) := \phi_{\min}(m, \mathbb{E}[\tilde{x}\tilde{x}'])$ . To assume that  $\bar{\phi}_{\min}(m) > 0$  requires that all population Gram submatrices formed by any  $m$  components of  $\tilde{x}_i$  are positive definite.

Next we state our main condition, which contains the previously defined approximate sparsity as well as other more technical assumptions. Throughout the paper, let  $c$  and  $C$  be positive constants independent of  $n$ , and let  $\ell_n \nearrow \infty$ ,  $\delta_n \searrow 0$ , and  $\Delta_n \searrow 0$  be sequences of positive constants.

**Condition 1.** (i)  $(\epsilon_i)_{i=1}^n$  is a sequence of i.i.d. random variables with common distribution function  $F$  such that  $F(0) = 1/2$ , independent of the random vectors  $\{(d_i, x_i)'\}_{i=1}^n$ .  $\{(y_i, d_i, x_i)'\}_{i=1}^n$  is a sequence of i.i.d. random vectors generated according to models (1.1) and (1.4). (ii)  $c \leq \mathbb{E}[v^2 | x]$  and  $\mathbb{E}[|v|^3 | x] \leq C$ , a.s., and  $\mathbb{E}[d^4] + \mathbb{E}[v^4] + \max_{1 \leq j \leq p} (\mathbb{E}[x_j^2 d^2] + \mathbb{E}[|x_j v|^3]) \leq C$ . (iii) There exists  $s = s_n \geq 1$  such that  $\|\beta_0\|_0 \leq s$  and  $\|\theta_0\|_0 \leq s$ . (iv) The error distribution  $F$  is absolutely continuous with continuously differentiable density  $f_\epsilon(\cdot)$  such that  $f_\epsilon(0) \geq c > 0$  and  $f_\epsilon(t) \vee |f'_\epsilon(t)| \leq C$  for all  $t \in \mathbb{R}$ , (v) with probability  $1 - \Delta_n$  we have  $K_x \geq \max_{1 \leq i \leq n} \|x_i\|_\infty$  and  $(K_x^4 + K_x^2 s^2 + s^3) \log^3(p \vee n) \leq n \delta_n$ . (vi) We have  $c \leq \bar{\phi}_{\min}(\ell_n s) \leq \bar{\phi}_{\max}(\ell_n s) \leq C$ .

**Comment 2.4.** Condition 1(i) imposes the setting discussed in the previous section with the zero conditional median of the error distribution. Condition 1(ii) imposes moment conditions on the structural

errors and regressors to ensure good model selection performance of Lasso applied to equation (1.4). The approximate sparsity 1 (iii) imposes sparsity of the high-dimensional vectors  $\beta_0$  and  $\theta_0$ . In the theorems below we provide the required technical conditions on the growth of  $s \log p$  since it is dependent on the choice of algorithm. Condition 1(iv) is a set of standard assumptions in the LAD literature (see [16]) and in the instrumental quantile regression literature [12]. Condition 1(v) restricts the sparsity index, so that  $s^3 \log^3(p \vee n) = o(n)$  is required; this is analogous to the restriction  $p^3(\log p)^2 = o(n)$  made in [13] in the problem without selection. Most importantly, no assumptions on the separation from zero of the non-zero coefficients of  $\theta_0$  and  $\beta_0$  are made.

**Comment 2.5.** Condition 1(vi) is quite plausible for many designs of interest. Combined with Condition 1(v), an equivalence between the norms induced by the empirical Gram matrix and the population Gram matrix over  $s$ -sparse vectors follows. Other examples of such equivalence are: Theorem 3.2 in [25] (see also [33] and [1]) for i.i.d. zero-mean sub-Gaussian regressors and  $s \log^2(n \vee p) \leq \delta_n n$ ; Theorem 1.8 [25] (see also Lemma 1 in [4]) for i.i.d. uniformly bounded zero-mean regressors and  $s(\log^3 n) \log(p \vee n) \leq \delta_n n$ .

**2.3. Results.** We begin with considering the estimators generated by Algorithms 1 and Algorithm 2.

**Theorem 1** (Robust Estimation and Inference). *Let  $\tilde{\alpha}$  and  $L_n$  be obtained by Algorithm 1 or Algorithm 2. Suppose that Condition 1 is satisfied for all  $n \geq 1$ . Moreover, suppose that with probability at least  $1 - \Delta_n$ ,  $\|\widehat{\beta}\|_0 \leq Cs$ . Then, as  $n \rightarrow \infty$  and for  $\sigma_n^2 = 1/(4f_\epsilon^2 \mathbb{E}[v^2])$ ,*

$$\sigma_n^{-1} \sqrt{n}(\tilde{\alpha} - \alpha_0) \rightsquigarrow N(0, 1) \text{ and } nL_n(\alpha_0) \rightsquigarrow \chi^2(1).$$

Theorem 1 establishes the first main result of the paper. Algorithm 1 relies on the post model selection estimators which in turn hinge on achieving sufficiently sparse estimates  $\widehat{\beta}$  and  $\widehat{\theta}$ . Sparsity of the former can be directly achieved under sharp penalty choices for optimal rates as discussed in the Supplementary Appendix F.2. The sparsity for the latter potentially requires heavier penalty as shown in [3]. Alternatively, sparsity for the estimator in Step 1 can also be achieved by truncating the smallest components of estimate  $\widehat{\beta}$ .<sup>2</sup> Algorithm 2 relies on the regularized estimators instead of the post-model selection estimators. Theorem 1 establishes that Algorithm 2 achieves the same inferential guarantees as Algorithm 1.

An important consequence of these results is the following corollary. Here  $\mathcal{P}_n$  denotes a collection of distributions for  $\{(y_i, d_i, x_i')'\}_{i=1}^n$  and for  $P_n \in \mathcal{P}_n$  the notation  $\text{pr}_{P_n}$  means that under  $\text{pr}_{P_n}$ ,  $\{(y_i, d_i, x_i')'\}_{i=1}^n$  is distributed according to the law determined by  $P_n$ .

**Corollary 1 (Uniformly Valid Confidence Intervals).** *Let  $\tilde{\alpha}$  be the estimator of  $\alpha_0$  constructed according to Algorithm 1 or Algorithm 2, and let  $\mathcal{P}_n$  be the collection of all distributions of  $\{(y_i, d_i, x_i')'\}_{i=1}^n$  for which Condition 1 and  $\|\widehat{\beta}\|_0 \leq Cs$  holds with with probability at least  $1 - \Delta_n$  for given  $n \geq 1$ . Then as  $n \rightarrow \infty$ ,*

$$\sup_{P_n \in \mathcal{P}_n} \left| \text{pr}_{P_n}(\alpha_0 \in [\tilde{\alpha} \pm \sigma_n z_{\xi/2} / \sqrt{n}]) - (1 - \xi) \right| = o(1), \quad \sup_{P_n \in \mathcal{P}_n} \left| \text{pr}_{P_n}(\alpha_0 \in \widehat{A}_{n,\xi}) - (1 - \xi) \right| = o(1),$$

<sup>2</sup>Lemma 3 in Appendix D.2 formally shows that a suitable truncation preserves the rate of convergence under our conditions.



where  $z_{\xi/2} = \Phi^{-1}(1 - \xi/2)$  and  $\widehat{A}_{n,\xi} = \{\alpha \in \mathcal{A} : nL_n(\alpha) \leq (1 - \xi)\text{-quantile of } \chi^2(1)\}$ .

Corollary 1 establishes the second main result of the paper; it highlights the uniformity nature of the results. As long as the overall sparsity requirements hold, imperfect model selection in Steps 1 and 2 do not compromise the results. The robustness of the approach is also apparent from the fact that Corollary 1 allows for the data-generating process to change with  $n$ . This result is new even under the traditional case of fixed- $p$  asymptotics. Condition 1 explicitly characterizes regions of data-generating processes for which the uniformity result holds. Simulations results discussed next also provide an additional evidence that these regions are substantial.

**2.4. Generalization to Many Target Coefficients and Approximate Sparsity.** We consider the following generalization to the previous model:

$$y = \sum_{\ell=1}^{p_1} d_\ell \alpha_\ell + g(u) + \epsilon, \quad \epsilon \sim F, \quad F(0) = 1/2.$$

where  $(d, u)$  are regressors, and  $\epsilon$  is the noise with distribution function  $F$  that is independent of regressors, and has median 0, i.e.  $F(0) = 1/2$ . The coefficients  $\alpha_\ell$  for each  $\ell \in \mathcal{L} = \{1, \dots, p_1\}$  are now the high-dimensional parameter of interest.

We can rewrite this model as  $p_1$  models of the previous form:

$$y = \alpha_\ell d_\ell + g_\ell(z_\ell) + \epsilon, \quad d_\ell = m_\ell(z_\ell) + v_\ell, \quad \mathbb{E}[v_\ell | z_\ell] = 0, \quad (\ell \in \mathcal{L}), \quad (2.14)$$

where  $\alpha_\ell$  is the target coefficient,  $g_\ell(z_\ell) = \sum_{l \neq \ell}^{p_1} d_l \alpha_l + g(u)$ ,  $m_\ell(z_\ell) = \mathbb{E}[d_\ell | z_\ell]$  and  $z_\ell = (d, u) \setminus d_\ell$ . We would like to estimate and perform inference on each of the  $p_1$  coefficients  $\alpha_\ell$  simultaneously. Moreover, we would like to allow regression functions  $g_\ell$  and  $m_\ell$  that are approximately sparse in terms of some dictionary of technical regressors, generalizing the previous exact sparsity. By approximate sparsity we mean that we can decompose the regression function into a sum of a sparse approximation and an approximation error,

$$g_\ell(z_\ell) = \sum_{k=1}^p f_{\ell k}(z_\ell) \theta_{\ell k} + r_{g_\ell}(z_\ell), \quad m_\ell = \sum_{k=1}^p f_{\ell k}(z_\ell) \vartheta_{\ell k} + r_{m_\ell}(z_\ell) \quad (2.15)$$

where the sparse approximation is formulated in terms of a dictionary  $\{f_{\ell k}\}_{k=1}^K$  of technical regressors containing  $(d_l)_{l \neq \ell}$  as a subvector. Moreover, we require that the sparse approximations have dimension  $1 \leq s < n$  and the resulting squared approximation errors are small in expectation, namely for each  $\ell \in \mathcal{L}$

$$\|(\theta_{\ell k})_{k=1}^p\|_0 \leq s, \quad \|(\vartheta_{\ell k})_{k=1}^p\|_0 \leq s, \quad \mathbb{E}[r_{g_\ell}^2(z_\ell)] \leq Cs/n, \quad \mathbb{E}[r_{m_\ell}^2(z_\ell)] \leq Cs/n.$$

The approximately sparse framework above is quite general, in particular it contains traditional linear sieve/series framework which uses  $s = o(n)$  dictionary terms to approximate and estimate regression functions. Here we allow for the same possibility, except that we assume no a priori knowledge of the most important dictionary terms. It is important to note that we allow for the regression functions to have non-zero Fourier coefficients associated with each term in the dictionary, but we do require that keeping the largest  $s$  coefficients while setting to zero the rest does produce a good approximation to the target regression function.

Given the setting, we can apply our previous instrumental median regression to estimate the sparse approximations to the function  $h_\ell = (g_\ell, m_\ell)$  appearing in equation (2.14), and each of the target parameters  $(\alpha_\ell)_{\ell \in \mathcal{L}}$  can be identified and estimated by working with the system of “immunized” equations:

$$\mathbb{E}[\psi_\ell\{w_\ell, \alpha_\ell, h_\ell(z_\ell)\}] = 0,$$

where  $\psi_\ell(w_\ell, \alpha, t) = \varphi(y - d_\ell \alpha - t_1)(d_\ell - t_2)$  and  $w_\ell = (y, d_\ell, z_\ell)$ , and  $\partial_t \mathbb{E}[\psi_\ell(w_\ell, \alpha_\ell, t)|z_\ell]|_{t=h_\ell(z_\ell)} = 0$  is the “immunization” property. We will treat this generalization as a special case of a more general, unified framework of the next section.

### 3. INFERENCE ON MANY TARGET PARAMETERS IN Z-PROBLEMS WITH APPROXIMATELY SPARSE NUISANCE FUNCTIONS

In this section we generalize the previous example to a more general setting, where  $p_1$  target parameters defined via Huber’s Z-problems are of interest, with dimension  $p_1$  potentially much larger than the sample size. This framework covers the median regression example, its generalization discussed above, as well many other semi-parametric models.

The interest lies in  $p_1$  real-valued target parameters  $\alpha_\ell$  indexed by  $\ell \in \mathcal{L} = \{1, \dots, p_1\}$ . We assume that  $\alpha_\ell \in \mathcal{A}_\ell \subset \mathcal{A}$  for each  $\ell$ , where  $\mathcal{A}$  is a fixed compact interval in  $\mathbb{R}$ . For each  $\ell \in \mathcal{L}$  the true value  $\alpha_\ell$  is identified as a unique solution of the following moment condition:

$$\mathbb{E}[\psi_\ell\{w_\ell, \alpha_\ell, h_\ell(z_\ell)\}] = 0. \tag{3.16}$$

Here for each  $\ell \in \mathcal{L}$ , vector  $w_\ell$  is a random vector taking values in  $\mathcal{W}_\ell \subset \mathbb{R}^{d_w}$ , containing vector  $z_\ell$  taking values in  $\mathcal{Z}_\ell$  as a subcomponent; the function  $(w, \alpha, t) \mapsto \psi_\ell(w, \alpha, t)$  is a measurable map from an open neighborhood of  $\mathcal{W}_\ell \times \mathcal{A}_\ell \times T_\ell$  to  $\mathbb{R}$ , and  $z \mapsto h_\ell(z) = \{h_{\ell m}(z)\}_{m=1}^M$  is a measurable map from  $\mathcal{Z}_\ell$  to  $T_\ell \subset \mathbb{R}^M$ , where  $M$  is fixed. The latter map is the nuisance parameter, possibly infinite-dimensional.

We assume that the nuisance functions  $(h_\ell)_{\ell \in \mathcal{L}}$  are approximately sparse in the sense of Condition 3 given below. We also assume that these functions can be estimated via sparse estimators, generated by the use of the post-selection or  $\ell_1$ -penalized methods; with examples being given in the previous section. We let  $\hat{h}_\ell = (\hat{h}_{\ell m})_{m=1}^M$  denote the estimator of  $h_\ell$ , which obeys Condition 3 stated below. The estimator  $\hat{\alpha}_\ell$  of  $\alpha_\ell$  is constructed as an Z-estimator, which solves the sample analogue of the equation (3.16):

$$|\mathbb{E}_n[\psi_\ell\{w_\ell, \hat{\alpha}_\ell, \hat{h}_\ell(z)\}]| \leq \inf_{\alpha \in \mathcal{A}_\ell} |\mathbb{E}_n[\psi\{w_\ell, \alpha, \hat{h}_\ell(z)\}]| + \epsilon_n, \tag{3.17}$$

where  $\epsilon_n = o(b_n^{-1}n^{-1/2})$  is the numerical tolerance parameter, and  $b_n = \sqrt{\log(ep_1)}$ .

In order to achieve robust inference results, we shall need to rely on the condition of orthogonality (“immunity”) of the scores with respect to small perturbations in the value of the nuisance parameters, which we can express in the following condition:

$$\partial_t \mathbb{E}[\psi_\ell(w_\ell, \alpha_\ell, t)|z_\ell]|_{t=h_\ell(z_\ell)} = 0, \text{ a.s.}, \tag{3.18}$$

where here and below we use the symbol  $\partial_t$  to abbreviate  $\partial/\partial t$ . It is important to construct the scores  $\psi_\ell$  to have property (3.18). Generally, we can construct the scores  $\psi_\ell$  that obey (3.18) by projecting some

initial non-orthogonal scores onto the orthocomplement of the tangent space for the nuisance parameter (see [30, 29, 17]). Sometimes the resulting construction generates additional nuisance parameters, for example, the auxiliary regression function in the case of the median regression problem in Section 2.

In what follows, we shall denote by  $\varsigma$ ,  $c_0$ ,  $c$ ,  $n_0$ , and  $C$  some positive constants.

**Condition 2.** For each  $n$ , we observe independent and identically distributed copies of  $(w_i)_{i=1}^n$  of the random vector  $w = (w_\ell)_{\ell \in \mathcal{L}}$ , whose law is determined by the probability measure  $P \in \mathcal{P}_n$ . Uniformly for all  $n \geq n_0$  and  $P \in \mathcal{P}_n$  and  $\ell \in \mathcal{L}$ , the following conditions hold. (i) The true parameter values  $\alpha_\ell$  obeys (3.16); there is an interval of fixed positive radius centered at  $\alpha_\ell$  contained in  $\mathcal{A}_\ell \subset \mathcal{A}$ , where  $\mathcal{A}$  is a fixed compact set in  $\mathbb{R}$ . (ii) For each  $\nu = (\nu_k)_{k=1}^{1+M} = (\alpha, t) \in \mathcal{A}_\ell \times T_\ell$ , the map  $\nu \mapsto \mathbb{E}\{\psi_\ell(w_\ell, \nu) | z_\ell\}$  is twice continuously differentiable a.s., and  $\mathbb{E}[\sup_{\nu \in \mathcal{A}_\ell \times T_\ell} |\partial_{\nu_r} \mathbb{E}\{\psi_\ell(w_\ell, \nu) | z_\ell\}|^2] \leq C$ ,  $\sup_{\nu \in \mathcal{A}_\ell \times T_\ell} |\partial_{\nu_k} \partial_{\nu_r} \mathbb{E}\{\psi_\ell(w_\ell, \nu) | z_\ell\}| \leq C$  for each  $r = 1, \dots, K$  and  $k = 1, \dots, K$ ; moreover,  $\sup_{(\nu, \bar{\nu}) \in (\mathcal{A}_\ell \times T_\ell)^2} \mathbb{E}[\{\psi_\ell(w_\ell, \nu) - \psi_\ell(w_\ell, \bar{\nu})\}^2 | z_\ell] \leq C \|\nu - \bar{\nu}\|^\varsigma$  almost surely. (iii) The orthogonality condition (3.18) holds. (iv) The following global and local identifiability condition holds:  $2|\mathbb{E}[\psi_\ell\{w_\ell, \alpha, h_\ell(z_\ell)\}]| \geq |\Gamma_\ell(\alpha - \alpha_\ell)| \vee c_0$  for all  $\alpha \in \mathcal{A}_\ell$ , where  $\Gamma_\ell = \partial_\alpha \mathbb{E}[\psi_\ell\{w_\ell, \alpha, h_\ell(z_\ell)\}]$ , with  $c \leq |\Gamma_\ell| \leq C$  for all  $\ell \in \mathcal{L}$ . (v) The moments of the scores are well-behaved:  $c \leq (\mathbb{E}[\psi_\ell^2\{w_\ell, \alpha, h_\ell(z_\ell)\}])^{1/2} \leq (\mathbb{E}[\psi_\ell^q\{w_\ell, \alpha, h_\ell(z_\ell)\}])^{1/q} \leq C$ , for  $q \geq 4$ .

These conditions impose rather mild assumptions for Z-estimation problems, in particular, allowing for non-smooth scores  $\psi_\ell$  such as those arising in median regression. These conditions are analogous to assumptions imposed in the  $p = o(n)$  setting, e.g., in [13]. In what follows, let  $\delta_n \searrow 0$  and  $\rho_n \searrow 0$  be a sequence of constants approaching zero from above. Let  $a_n = \max(p_1, p, n, e)$  and recall that  $b_n = \sqrt{\log(ep_1)}$ .

**Condition 3.** Uniformly for all  $n \geq n_0$  and  $P \in \mathcal{P}_n$  and  $\ell \in \mathcal{L}$ , the following conditions hold: (i) The nuisance functions  $h_\ell = (h_{\ell m})_{m=1}^M : \mathcal{Z}_\ell \mapsto T_\ell \subset \mathbb{R}^M$ , where  $M$  is fixed, are approximately sparse with sparsity index at most  $s = s_n \geq 1$ , namely,  $z \mapsto h_{\ell m}(z) = \sum_{k=1}^p f_{\ell mk}(z) \beta_{\ell km} + r_{\ell m}(z)$ , where  $f_{\ell mk} : \mathcal{Z}_\ell \mapsto \mathbb{R}$  are approximating functions,  $\beta_{\ell m} = (\beta_{\ell mk})_{k=1}^p$  obeys  $|\text{supp}(\beta_{\ell m})| \leq s$ , and the approximation errors  $(r_{\ell m})_{m=1}^M : \mathcal{Z}_\ell \rightarrow \mathbb{R}$  obey  $\|r_{\ell m}\|_{P,2} \leq C\{s \log(a_n)/n\}^{1/2}$  for all  $m$  and  $\ell$ . (ii) There is a sparse estimator  $\hat{h}_\ell = (\hat{h}_{\ell m})_{m=1}^M$  of  $h_\ell$  with good sparsity and rate properties, namely with probability  $1 - \delta_n$ ,  $\hat{h}_\ell \in \mathcal{H}_\ell$ , where  $\mathcal{H}_\ell = \times_{m=1}^M \mathcal{H}_{\ell m}$  consists of functions  $(\bar{h}_{\ell m})_{m=1}^M : \mathcal{Z}_\ell \rightarrow T_\ell$ , where  $z \mapsto \bar{h}_{\ell m}(z) = \sum_{k=1}^p f_{\ell mk}(z) \bar{\beta}_{\ell k}$  is such that  $|\text{supp}(\bar{\beta}_{\ell m})| \leq Cs$  and  $\|\sum_{k=1}^p f_{\ell mk}(z_\ell) (\bar{\beta}_{\ell mk} - \beta_{\ell mk})\|_{P,2} \leq C\{s \log(a_n)/n\}^{1/2}$  for all  $m$  and  $\ell$ . (iii) Application of  $\psi_\ell$  to  $\mathcal{A}_\ell$  and  $\mathcal{H}_\ell$  does not increase the entropy too much, namely  $\mathcal{F}_\ell = [\psi_\ell\{w_\ell, \alpha, h(z_\ell)\}, \alpha \in \mathcal{A}_\ell, h \in \mathcal{H}_{\ell m} \cup \{h_\ell\}]$  has  $\text{ent}(\varepsilon, \mathcal{F}_\ell) \leq C\{\log(e/\varepsilon) + \sum_{m=1}^M \text{ent}(\varepsilon/C, \mathcal{H}_{\ell m})\}$ . (iv) For  $F_\ell$  denoting the envelope of  $\mathcal{F}_\ell$ , and  $F = \max_{\ell \in \mathcal{L}} F_\ell$ , we have  $\|F\|_{P,q} \leq C$  for  $q \geq 4$ . (v) The dimensions  $p_1, p$ , and  $s$  obey the following growth conditions with respect to  $n$ :

$$n^{-1/2} \left( \sqrt{s \log a_n} + n^{-1/2} s n^{\frac{1}{4}} \log a_n \right) \leq \rho_n, \quad \rho_n^{\varsigma/2} \sqrt{s \log a_n} \leq \delta_n b_n^{-1}. \quad (3.19)$$

Condition 2(i), (ii), (v) records a formal sense in which approximate sparsity of  $h_\ell$  is used, as well as requires reasonable behavior of sparse estimators  $\hat{h}_\ell$ . In the previous section, this type of behavior occurred in the cases where  $h_\ell$  consisted of (a part of) median regression function and a conditional

expectation function in an auxiliary equation. There are lots of conditions in the literature that imply these conditions from various primitive assumptions. Note that for the case with  $q = \infty$ , condition (v) implies the following restrictions on the sparsity indices:  $s^2 \log^3 a_n/n \rightarrow 0$  for the case where  $\varsigma = 2$  (smooth  $\psi_\ell$ ) and  $s^3 \log^5 a_n/n \rightarrow 0$  for the case where  $\varsigma = 1$  (non-smooth  $\psi_\ell$ ). Condition (iii) is a mild condition on  $\psi_\ell$  – it holds for example, when  $\psi_\ell$  is generated by applying monotone and Lipschitz transformations to its arguments, as was the case in median regression (see [30] for many other ways). The condition (iv) bounds the moments of the envelopes, and it can be relaxed to a bound that grows with  $n$ , with an appropriate strengthening of the growth condition (3.19).

Define, for each  $\ell \in \mathcal{L}$ ,

$$\sigma_\ell^2 = \mathbb{E}[\Gamma_\ell^{-2} \psi_\ell^2(w_\ell, \alpha_\ell, h_\ell(z_\ell))], \quad \phi_\ell(w) = -\sigma_\ell^{-1} \Gamma_\ell^{-1} \psi_\ell(w_\ell, \alpha_\ell, h_\ell(z_\ell)).$$

**Theorem 2** (Uniform Bahadur Representation). *Under Conditions 2 and 3, uniformly in  $P \in \mathcal{P}_n$ , with probability  $1 - o(1)$ , as  $n \rightarrow \infty$ ,*

$$\max_{\ell \in \mathcal{L}} \left| \sqrt{n} \left( \frac{\hat{\alpha}_\ell - \alpha_\ell}{\sigma_\ell} \right) - \frac{1}{\sqrt{n}} \sum_{i=1}^n \phi_\ell(w) \right| \leq o(b_n^{-1}).$$

An immediate implication is a corollary on the uniform in  $P \in \mathcal{P}_n$  and  $\ell \in \mathcal{L}$  asymptotic normality, which follows from Liapunov's central limit theorem for triangular arrays.

**Corollary 2** (Uni-Dimensional Central Limit Theorem). *Under conditions of Theorem 2, as  $n \rightarrow \infty$ ,*

$$\max_{\ell \in \mathcal{L}} \sup_{P \in \mathcal{P}_n} \sup_{t \in \mathbb{R}} \left| \Pr_P \left\{ \sqrt{n} \left( \frac{\hat{\alpha}_\ell - \alpha_\ell}{\sigma_\ell} \right) \leq t \right\} - \Pr_P \{N(0, 1) \leq t\} \right| = o(1).$$

*This implies, in particular, that for  $\bar{c}_{1-a} = (1 - a/2)$ -quantile of  $N(0, 1)$  variable,*

$$\max_{\ell \in \mathcal{L}} \sup_{P \in \mathcal{P}_n} \left| \Pr_P \left\{ \alpha_\ell \in [\hat{\alpha}_\ell \pm \hat{\sigma}_\ell \bar{c}_{1-a/2} n^{-1/2}] \right\} - (1 - a) \right| = o(1),$$

*provided  $\max_{\ell \in \mathcal{L}} |\hat{\sigma}_\ell - \sigma_\ell| = o_P(1)$  uniformly in  $P \in \mathcal{P}_n$ .*

This result construct pointwise confidence bands for  $\alpha_\ell$ , and shows that they are valid uniformly in  $P \in \mathcal{P}$  and in  $\ell \in \mathcal{L}$ .

Another useful implication is the *high-dimensional* central limit theorem uniformly over rectangles in  $\mathbb{R}^{p_1}$ , provided that  $(\log p_1)^7 = o(n)$ , which follows from [11]'s central limit theorem for  $p_1$ -dimensional (approximate) sample means, with  $p_1$  potentially much larger than  $n$ . Let

$$\mathcal{N} = (\mathcal{N}_\ell)_{\ell \in \mathcal{L}} = N(0, \Omega)$$

be a random vector with normal distribution with mean zero and variance matrix  $\Omega$ , where  $\Omega_{\ell\bar{\ell}} = \mathbb{E}\{\bar{\psi}_\ell(w)\bar{\psi}_{\bar{\ell}}(w)\}$  for  $(\ell, \bar{\ell}) \in \mathcal{L}^2$ . Let  $\mathcal{R}$  be a collection of rectangles  $R$  in  $\mathbb{R}^{p_1}$  of the form

$$R = \left\{ z \in \mathbb{R}^{p_1} : \max_{\ell \in A} z_\ell \leq t, \max_{\ell \in B} -z_\ell \leq t \right\} \quad (t \in \mathbb{R}, A \subset \mathcal{L}, B \subset \mathcal{L}).$$

**Corollary 3** (High-Dimensional Central Limit Theorem over Rectangles). *Under conditions of Theorem 2, provided that  $(\log p_1)^7 = o(n)$ ,*

$$\sup_{P \in \mathcal{P}_n} \sup_{R \in \mathcal{R}} \left| \Pr_P \left\{ \sqrt{n} \left( \frac{\hat{\alpha}_\ell - \alpha_\ell}{\sigma_\ell} \right)_{\ell \in \mathcal{L}} \in R \right\} - \Pr_P \{ \mathcal{N} \in R \} \right| = o(1). \quad (3.20)$$

*This implies, in particular, that for  $c_{1-a} = (1-a)$ -quantile of  $\max_{\ell \in \mathcal{L}} |\mathcal{N}_\ell|$ ,*

$$\sup_{P \in \mathcal{P}_n} \left| \Pr_P \left\{ \alpha_\ell \in [\hat{\alpha}_\ell \pm c_{1-a} \sigma_\ell n^{-1/2}], \text{ for all } \ell \in \mathcal{L} \right\} - (1-a) \right| = o(1).$$

The result provides simultaneous confidence bands for  $(\alpha_\ell)_{\ell \in \mathcal{L}}$ , which are valid uniformly in  $P \in \mathcal{P}_n$ . Moreover, the result (3.20) is immediately useful for performing *multiple hypotheses testing* about  $(\alpha_\ell)_{\ell \in \mathcal{L}}$  via the step-down methods of [24] which control the family-wise error rates— we refer the reader to [11] for further discussion and details of multiple testing with  $p_1 \gg n$ .

In practice the distribution of  $\mathcal{N}$  is unknown due to the unknown covariance matrix, but it can be approximated by the Gaussian multiplier bootstrap procedure, which generates a vector  $\mathcal{N}^*$  as follows:

$$\mathcal{N}^* = (\mathcal{N}^*_\ell)_{\ell \in \mathcal{L}} = \left\{ \frac{1}{\sqrt{n}} \sum_{i=1}^n \xi_i \hat{\phi}_\ell(w_i) \right\}_{\ell \in \mathcal{L}}, \quad (3.21)$$

where  $(\xi_i)_{i=1}^n$  are i.i.d. draws of  $N(0, 1)$  variables, which are independently distributed of the data  $(w_i)_{i=1}^n$ , and  $\hat{\phi}_\ell$  are any estimators of  $\bar{\psi}_\ell$ , such that  $\max_{(\ell, \bar{\ell}) \in \mathcal{L}} |\mathbb{E}_n[\hat{\phi}_\ell(w) \hat{\phi}_{\bar{\ell}}(w)] - \mathbb{E}_n[\phi_\ell(w) \phi_{\bar{\ell}}(w)]| = o_P(b_n^{-4})$  uniformly in  $P \in \mathcal{P}_n$ . Let  $\hat{\sigma}_\ell^2 = \mathbb{E}_n[\hat{\phi}_\ell^2(w)]$ .

[11]’s results for multiplier bootstrap then imply the following theorem.

**Corollary 4** (Validity of Multiplier Bootstrap). *Under conditions of Theorem 2, provided that  $(\log p_1)^7 = o(n)$ , uniformly in  $P \in \mathcal{P}_n$  with probability  $1 - o(1)$ ,*

$$\sup_{P \in \mathcal{P}_n} \sup_{R \in \mathcal{R}} \left| \Pr_P \{ \mathcal{N}^* \in R \mid (w_i)_{i=1}^n \} - \Pr_P \{ \mathcal{N} \in R \} \right| = o(1). \quad (3.22)$$

*This implies, in particular, that for  $\hat{c}_{1-a} = (1-a)$ -quantile of  $\max_{\ell \in \mathcal{L}} |\mathcal{N}_\ell^*|$ ,*

$$\sup_{P \in \mathcal{P}_n} \left| \Pr_P \left\{ \alpha_\ell \in [\hat{\alpha}_\ell \pm \hat{c}_{1-a} \hat{\sigma}_\ell n^{-1/2}], \text{ for all } \ell \in \mathcal{L} \right\} - (1-a) \right| = o(1).$$

#### 4. MONTE-CARLO EXPERIMENTS

In this section we examine the finite sample performance of the proposed estimators. We focus on the estimator associated with Algorithm 1 based on post-model selection methods.

We considered the following regression model:

$$y = d\alpha_0 + x'(c_y \theta_0) + \epsilon, \quad d = x'(c_d \theta_0) + v, \quad (4.23)$$

where  $\alpha_0 = 1/2$ ,  $\theta_{0j} = 1/j^2$ ,  $j = 1, \dots, 10$ , and  $\theta_{0j} = 0$  otherwise,  $x = (1, z')'$  consists of an intercept and covariates  $z \sim N(0, \Sigma)$ , and the errors  $\epsilon$  and  $v$  are independently and identically distributed as  $N(0, 1)$ . The dimension  $p$  of the covariates  $x$  is 300, and the sample size  $n$  is 250. The regressors are correlated with  $\Sigma_{ij} = \rho^{|i-j|}$  and  $\rho = 0.5$ . The coefficients  $c_y$  and  $c_d$  are used to control the  $R^2$  of the reduce

form equation. For each equation, we consider the following values for the  $R^2$ :  $\{0, 0.1, 0.2, \dots, 0.8, 0.9\}$ . Therefore we have 100 different designs and results are based on 500 repetitions for each design. For each repetition we draw new vectors  $x_i$ 's and errors  $\epsilon_i$ 's and  $v_i$ 's.

The design above with  $x'(c_y\theta_0)$  is a sparse model. However, the decay of the components of  $\theta_0$  rules out typical “separation from zero” assumptions of the coefficients of “important” covariates (since the last component is of the order of  $1/n$ ), unless  $c_y$  is very large. Thus, we anticipate that “standard” post-selection inference procedures – which rely on model selection of the outcome equation only – work poorly in the simulation study. In contrast, based upon the prior theoretical arguments, we anticipate that our instrumental median estimator – which works off both equations in (4.23) – to work well in the simulation study.

The simulation study focuses on Algorithm 1. Standard errors are computed using the formula (1.10). (Algorithm 2 worked similarly, though somewhat worse due to larger biases). As the main benchmark we consider the standard post-model selection estimator  $\tilde{\alpha}$  based on the post  $\ell_1$ -penalized LAD method, as defined in (1.3).

In Figure 1, we display the (empirical) rejection probability of tests of a true hypothesis  $\alpha = \alpha_0$ , with nominal size of tests equal to 0.05. The left-top plot shows the rejection frequency of the standard post-model selection inference procedure based upon  $\tilde{\alpha}$  (where the inference procedure assumes perfect recovery of the true model). The rejection frequency deviates very sharply from the ideal rejection frequency of 0.05. This confirms the anticipated failure (lack of uniform validity) of inference based upon the standard post-model selection procedure in designs where coefficients are not well separated from zero (so that perfect recovery does not happen). In sharp contrast, the right top and bottom plots show that both of our proposed procedures (based on estimator  $\check{\alpha}$  and the result (1.9) and on the statistic  $L_n$  and the result (1.12)) perform well, closely tracking the ideal level of 0.05. This is achieved uniformly over all the designs considered in the study, and this confirms our theoretical results established in Corollary 1.

In Figure 2, we compare the performance of the standard post-selection estimator  $\tilde{\alpha}$  (defined in (1.3)) and our proposed post-selection estimator  $\check{\alpha}$  (obtained via Algorithm 1). We display results in three different metrics of performance – mean bias (top row), standard deviation (middle row), and root mean square error (bottom row) of the two approaches. The significant bias for the standard post-selection procedure occurs when the indirect equation (1.4) is nontrivial, that is, when the main regressor is correlated to other controls. Such bias can be positive or negative depending on the particular design. The proposed post-selection estimator  $\check{\alpha}$  performs well in all three metrics. The root mean square error for the proposed estimator  $\check{\alpha}$  are typically much smaller than those for standard post-model selection estimators  $\tilde{\alpha}$  (as shown by bottom plots in Figure 2). This is fully consistent with our theoretical results and minimax efficiency considerations given in Section 5.

## Appendix

### APPENDIX A. GENERALIZATION OF SECTION 3 TO HETEROSCEDASTIC CASE

We emphasize that both proposed algorithms exploit the homoscedasticity of the model (1.1) with respect to the error term  $\epsilon_i$ . The generalization to the heteroscedastic case can be achieved as follows. In order to achieve the semiparametric efficiency bound we need to consider the weighted version of the auxiliary equation (1.4). Specifically, we can rely on the following of weighted decomposition:

$$f_i d_i = f_i x_i' \theta_0^* + v_i^*, \quad \mathbb{E}[f_i v_i^* | x_i] = 0, \quad i = 1, \dots, n, \quad (\text{A.24})$$

where the weights are conditional densities of error terms  $\epsilon_i$  evaluated at their medians of 0,

$$f_i = f_{\epsilon_i}(0 | d_i, x_i), \quad i = 1, \dots, n, \quad (\text{A.25})$$

which in general vary under heteroscedasticity. With that in mind it is straightforward to adapt the proposed algorithms when the weights  $(f_i)_{i=1}^n$  are known. For example Algorithm 1 becomes as follows.

**Algorithm 1' (Based on Post-Model Selection estimators).**

- (1) Run Post- $\ell_1$ -penalized LAD of  $y_i$  on  $d_i$  and  $x_i$ ; keep fitted value  $x_i' \tilde{\beta}$ .
- (2) Run Post-Lasso of  $f_i d_i$  on  $f_i x_i$ ; keep the residual  $\tilde{v}_i^* := f_i(d_i - x_i' \tilde{\theta})$ .
- (3) Run Instrumental LAD regression of  $y_i - x_i' \tilde{\beta}$  on  $d_i$  using  $\tilde{v}_i^*$  as the instrument for  $d_i$  to compute the estimator  $\tilde{\alpha}$ . Report  $\tilde{\alpha}$  and/or perform inference.

An analogous generalization of Algorithm 2 based on regularized estimator results from removing the word ‘‘Post’’ in the algorithm above.

Under similar regularity conditions, uniformly over a large collection  $\mathcal{P}_n^*$  of distributions of  $\{(y_i, d_i, x_i')\}_{i=1}^n$ , the estimator  $\tilde{\alpha}$  above obeys

$$(4\mathbb{E}[v^{*2}])^{1/2} \sqrt{n}(\tilde{\alpha} - \alpha_0) \rightsquigarrow N(0, 1). \quad (\text{A.26})$$

Moreover, the criterion function at the true value  $\alpha_0$  in Step 3 also has a pivotal behavior, namely

$$nL_n(\alpha_0) \rightsquigarrow \chi^2(1), \quad (\text{A.27})$$

which can also be used to construct a confidence region  $\hat{A}_{n,\xi}$  based on the  $L_n$ -statistic as in (1.12) with coverage  $1 - \xi$  uniformly over the collection of distributions  $\mathcal{P}_n^*$ .

In practice the density function values  $(f_i)_{i=1}^n$  are typically unknown and need to be replaced by estimates  $(\hat{f}_i)_{i=1}^n$ . The analysis of the impact of such estimation is very delicate and is developed in the companion work [8], which considers the more general problem of uniformly valid inference for quantile regression models in approximately sparse models.

### APPENDIX B. ADDITIONAL DISCUSSION FOR SECTION 3

**B.1. Connection to Neymanization.** In this section we make some connections to Neyman’s  $C(\alpha)$  test ([21, 22]). For the sake of exposition we assume that  $(y_i, d_i, x_i)_{i=1}^n$  are i.i.d. but we shall use the

heteroscedastic setup introduced in the previous section. We consider the estimating equation for  $\alpha_0$ :

$$\mathbb{E}[\varphi(y_i - d_i\alpha_0 - x_i'\beta_0)v_i] = 0.$$

Our problem is to find useful instruments  $v_i$  such that

$$\frac{\partial}{\partial \beta} \mathbb{E}[\varphi(y_i - d_i\alpha_0 - x_i'\beta)v_i] \Big|_{\beta=\beta_0} = 0.$$

If this property holds, the estimator of  $\alpha_0$  will be “immunized” against “crude” or nonregular estimation of  $\beta_0$ , for example, via a post-selection procedure or some regularization procedure. Such immunization ideas are in fact behind Neyman’s classical construction of his  $C(\alpha)$  test, so we shall use the term “Neymanization” to describe such procedure. There will be many instruments  $v_i$  that can achieve the property stated above, and there will be one that is optimal.

The instruments can be constructed by taking  $v_i := z_i/f_i$ , where  $z_i$  is the residual in the regression equation:

$$w_i d_i = w_i m_0(x_i) + z_i, \quad \mathbb{E}[w_i z_i | x_i] = 0, \quad (\text{B.28})$$

where  $w_i$  is a nonnegative weight, a function of  $(d_i, z_i)$  only, for example  $w_i = 1$  or  $w_i = f_i$  (the latter choice will in fact be optimal). Note that function  $m_0(x_i)$  solves the least squares problem

$$\min_{h \in \mathcal{H}} \mathbb{E} [\{wd - wh(x)\}^2], \quad (\text{B.29})$$

where  $\mathcal{H}$  is the class of measurable functions  $h(x)$  such that  $\mathbb{E}[w^2 h^2(x)] < \infty$ . Our assumption is that the  $m_0(x)$  is a sparse function  $x'\theta_0$ , with  $\|\theta_0\|_0 \leq s$  so that

$$w_i d_i = w_i x_i' \theta_0 + z_i, \quad \mathbb{E}[w_i z_i | x_i] = 0. \quad (\text{B.30})$$

In finite samples, the sparsity assumption allows to employ post-Lasso and Lasso to solve the least squares problem above approximately, and estimate  $z_i$ . Of course, the use of other structured assumptions may motivate the use of other regularization methods.

Arguments similar to those in the proofs show that, for  $\sqrt{n}(\alpha - \alpha_0) = O(1)$ ,

$$\sqrt{n} \{ \mathbb{E}_n[\varphi(y - d\alpha - x'\widehat{\beta})v] - \mathbb{E}_n[\varphi(y - d\alpha - x'\beta_0)v] \} = o_P(1),$$

for  $\widehat{\beta}$  based on a sparse estimation procedure, despite the fact that  $\widehat{\beta}$  converges to  $\beta_0$  at a slower rate than  $1/\sqrt{n}$ . That is, the empirical estimating equations behave as if  $\beta_0$  is known. Hence for estimation we can use  $\widehat{\alpha}$  as a minimizer of the statistic:

$$L_n(\alpha) = c_n^{-1} |\sqrt{n} \mathbb{E}_n[\varphi(y - d\alpha - x'\widehat{\beta})v]|^2,$$

where  $c_n = \mathbb{E}_n[v^2]/4$ . Since  $L_n(\alpha_0) \rightsquigarrow \chi^2(1)$ , we can also use the statistic directly for testing hypotheses and for construction of confidence sets.

This is in fact a version of Neyman’s  $C(\alpha)$  test statistic, adapted to the present non-smooth setting. The usual expression of  $C(\alpha)$  statistic is different. To see a more familiar form, note that  $\theta_0 = \mathbb{E}[w^2 x x']^{-1} \mathbb{E}[w^2 d x']$ , where  $A^{-}$  denotes a generalized inverse of  $A$ , and write

$$v_i = (w_i/f_i)d_i - (w_i/f_i)x_i' \mathbb{E}[w^2 x x']^{-1} \mathbb{E}[w^2 d x'], \quad \text{and} \quad \widehat{\varphi}_i := \varphi(y_i - d_i\alpha - x_i'\widehat{\beta}),$$



so that,

$$L_n(\alpha) = c_n^{-1} |\sqrt{n} \{ \mathbb{E}_n[\widehat{\varphi}(w/f)d] - \mathbb{E}_n[\widehat{\varphi}(w/f)x]' \mathbb{E}[w^2 x x']^{-1} \mathbb{E}[w^2 d x'] \}|^2.$$

This is indeed a familiar form of a  $C(\alpha)$  statistic.

The estimator  $\widehat{\alpha}$  that minimizes  $L_n$  up to  $o_P(1)$ , under suitable regularity conditions,

$$\sigma_n^{-1} \sqrt{n}(\widehat{\alpha} - \alpha_0) \rightsquigarrow N(0, 1), \quad \sigma_n^2 = \frac{1}{4} \mathbb{E}[fdv]^{-2} \mathbb{E}[v^2].$$

It is easy to show that the smallest value of  $\sigma_n^2$  is achieved by using  $v_i = v_i^*$  induced by setting  $w_i = f_i$ :

$$\sigma_n^{*2} = \frac{1}{4} \mathbb{E}[v^{*2}]^{-1}. \quad (\text{B.31})$$

Thus, setting  $w_i = f_i$  gives an optimal instrument amongst all “immunizing” instruments generated by the process described above. Obviously, this improvement translates into shorter confidence intervals and better testing based on either  $\widehat{\alpha}$  or  $L_n$ . While  $w_i = f_i$  is optimal,  $f_i$  will have to be estimated in practice, resulting actually in more stringent condition than when using non-optimal, known weights, e.g.,  $w_i = 1$ . The use of known weights may also give better behavior under misspecification of the model. Under homoscedasticity,  $w_i = 1$  is an optimal weight.

**B.2. Minimax Efficiency.** There is also a clean connection to the (local) minimax efficiency analysis from the semiparametric efficiency analysis. [18] derives an efficient score function for the partially linear median regression model:

$$S_i = 2\varphi(y_i - d_i\alpha_0 - x_i'\beta_0) f_i[d_i - m_0^*(x)],$$

where  $m_0^*(x_i)$  is  $m_0(x_i)$  in (B.28) induced by the weight  $w_i = f_i$ :

$$m_0^*(x_i) = \frac{\mathbb{E}[f_i^2 d_i | x_i]}{\mathbb{E}[f_i^2 | x_i]}.$$

Using the assumption  $m_0^*(x_i) = x_i'\theta_0^*$ , where  $\|\theta_0^*\|_0 \leq s \ll n$  is sparse, we have that

$$S_i = 2\varphi(y_i - d_i\alpha_0 - x_i'\beta_0) v_i^*,$$

which is the score that was constructed using Neymanization. It follows that the estimator based on the instrument  $v_i^*$  is actually efficient in the minimax sense (see Theorem 18.4 in [17]), and inference about  $\alpha_0$  based on this estimator provides best minimax power against local alternatives (see Theorem 18.12 in [17]).

The claim above is formal as long as, given a law  $P_n$ , the least favorable submodels are permitted as deviations that lie within the overall model. Specifically, given a law  $P_n$ , we shall need to allow for a certain neighborhood  $\mathcal{P}_n^\delta$  of  $P_n$  such that  $P_n \in \mathcal{P}_n^\delta \subset \mathcal{P}_n$ , where the overall model  $\mathcal{P}_n$  is defined similarly as before, except now permitting heteroscedasticity (or we can keep homoscedasticity  $f_i = f_\epsilon$  to maintain formality). To allow for this we consider a collection of models indexed by a parameter  $t = (t_1, t_2)$ :

$$y_i = d_i(\alpha_0 + t_1) + x_i'(\beta_0 + t_2\theta_0^*) + \epsilon_i, \quad \|t\| \leq \delta, \quad (\text{B.32})$$

$$f_i d_i = f_i x_i' \theta_0^* + v_i^*, \quad \mathbb{E}[f_i v_i^* | x_i] = 0, \quad (\text{B.33})$$

where  $\|\beta_0\|_0 \vee \|\theta_0^*\|_0 \leq s/2$  and conditions as in Section 2 hold. The case with  $t = 0$  generates the model  $P_n$ ; by varying  $t$  within  $\delta$ -ball, we generate models  $\mathcal{P}_n^\delta$ , containing the least favorable deviations. By [18],

the efficient score for the model given above is  $S_i$ , so we cannot have a better regular estimator than the estimator whose influence function is  $J^{-1}S_i$ , where  $J = E[S_i^2]$ . Since our model  $\mathcal{P}_n$  contains  $\mathcal{P}_n^\delta$ , all the formal conclusions about (local minimax) optimality of our estimators hold from theorems cited above (using subsequence arguments to handle models changing with  $n$ ). Our estimators are regular, since under  $\mathcal{P}_n^t$  with  $t = (O(1/\sqrt{n}), o(1))$ , their first order asymptotics do not change, as a consequence of Theorems in Section 2. (Though our theorems actually prove more than this.)

### APPENDIX C. INSTRUMENTAL LAD REGRESSION WITH ESTIMATED INPUTS

Throughout this section  $(x_i)_{i=1}^n$  are non-stochastic, and let  $\bar{E}[f(y, d, x)] := \mathbb{E}_n[E[f(y, d, x) | x]]$ ,

$$\psi_{\alpha, \beta, \theta}(y_i, d_i, x_i) = (1/2 - 1\{y_i \leq x_i' \beta + d_i \alpha\})(d_i - x_i' \theta) = (1/2 - 1\{y_i \leq x_i' \beta + d_i \alpha\})\{v_i - x_i'(\theta - \theta_0)\}$$

$$\text{and } \mathbb{G}_n(f) = n^{-1/2} \sum_{i=1}^n \{f(y_i, d_i, x_i) - E[f(y_i, d_i, x_i)]\}.$$

For fixed  $\alpha \in \mathbb{R}$  and  $\beta, \theta \in \mathbb{R}^p$ , define the function

$$\Gamma(\alpha, \beta, \theta) := \bar{E}[\psi_{\tilde{\alpha}, \tilde{\beta}, \tilde{\theta}}(y, d, x)] \Big|_{\tilde{\alpha}=\alpha, \tilde{\beta}=\beta, \tilde{\theta}=\theta}$$

where the expectation is conditional on  $(x_i)_{i=1}^n$ . For the notational convenience, let  $h = (\beta', \theta)'$ ,  $h_0 = (\beta_0', \theta_0)'$  and  $\hat{h} = (\hat{\beta}', \hat{\theta})'$ . The partial derivative of  $\Gamma(\alpha, \beta, \theta)$  with respect to  $\alpha$  is denoted by  $\Gamma_1(\alpha, \beta, \theta)$  and the partial derivative of  $\Gamma(\alpha, \beta, \theta)$  with respect to  $h = (\beta', \theta)'$  is denoted by  $\Gamma_2(\alpha, \beta, \theta)$ . Consider the following high-level condition. Here  $(\hat{\beta}', \hat{\theta})'$  is a generic estimator of  $(\beta_0', \theta_0)'$  (and not necessarily  $\ell_1$ -LAD and Lasso estimators, reps.), and  $\tilde{\alpha}$  is defined by  $\tilde{\alpha} \in \arg \min_{\alpha \in \mathcal{A}} L_n(\alpha)$  with this  $(\hat{\beta}', \hat{\theta})'$ , where  $\mathcal{A}$  here is also a generic (possibly random) compact interval. We assume that  $(\hat{\beta}', \hat{\theta})'$ ,  $\mathcal{A}$  and  $\tilde{\alpha}$  satisfy the following conditions.

**Condition ILAD.** Let  $\{(y_i, d_i)'\}_{i=1}^n$  be a sequence of independent random vectors generated according to models (1.1) and (1.4). Suppose (i)  $f_\epsilon(t) \vee |f'_\epsilon(t)| \leq C$  for all  $t \in \mathbb{R}$ ,  $\bar{E}[v^2] \geq c > 0$ , and  $\bar{E}[v^4] \vee \bar{E}[d^4] \leq C$ . Moreover, for some sequences  $\delta_n \searrow 0$  and  $\Delta_n \searrow 0$ , with probability at least  $1 - \Delta_n$ ,  
 (ii)  $\{\alpha : |\alpha - \alpha_0| \leq n^{-1/2}/\delta_n\} \subset \mathcal{A}$ , where  $\mathcal{A}$  is a (possibly random) compact interval;  
 (iii) the estimated parameters  $(\hat{\beta}', \hat{\theta})'$  satisfy

$$\left\{1 \vee \max_{1 \leq i \leq n} (E[|v_i| | x_i] \vee |x_i'(\hat{\theta} - \theta_0)|)\right\}^{1/2} \|x'(\hat{\beta} - \beta_0)\|_{2,n} \leq \delta_n n^{-1/4}, \quad \|x'(\hat{\theta} - \theta_0)\|_{2,n} \leq \delta_n n^{-1/4}, \quad (\text{C.34})$$

$$\sup_{\alpha \in \mathcal{A}} |\mathbb{G}_n(\psi_{\alpha, \hat{\beta}, \hat{\theta}} - \psi_{\alpha, \beta_0, \theta_0})| \leq \delta_n; \quad (\text{C.35})$$

(iv) the estimator  $\tilde{\alpha}$  satisfies  $|\tilde{\alpha} - \alpha_0| \leq \delta_n$ .

**Comment C.1.** Condition ILAD suffices to make the impact of the estimation of instruments negligible on the first order asymptotics of the estimator  $\tilde{\alpha}$ . We note that Condition ILAD covers several different estimators including both estimators proposed in Algorithms 1 and 2.

The following lemma summarizes the main inferential result based on the high level Condition ILAD.

**Lemma 1.** *Under Condition ILAD we have, for  $\sigma_n^2 = 1/(4f_\epsilon^2\bar{\mathbb{E}}[v^2])$ ,*

$$\sigma_n^{-1}\sqrt{n}(\tilde{\alpha} - \alpha_0) \rightsquigarrow N(0, 1) \text{ and } nL_n(\alpha_0) \rightsquigarrow \chi^2(1),$$

*Proof of Lemma 1.* We shall separate the proof into two parts.

Part 1. (Proof for the first assertion). Observe that

$$\begin{aligned} \mathbb{E}_n[\psi_{\tilde{\alpha}, \hat{\beta}, \hat{\theta}}(y, d, x)] &= \mathbb{E}_n[\psi_{\alpha_0, \beta_0, \theta_0}(y, d, x)] + \mathbb{E}_n[\psi_{\tilde{\alpha}, \hat{\beta}, \hat{\theta}}(y, d, x) - \psi_{\alpha_0, \beta_0, \theta_0}(y, d, x)] \\ &= \mathbb{E}_n[\psi_{\alpha_0, \beta_0, \theta_0}(y, d, x)] + \Gamma(\tilde{\alpha}, \hat{\beta}, \hat{\theta}) \\ &\quad + n^{-1/2}\mathbb{G}_n(\psi_{\tilde{\alpha}, \hat{\beta}, \hat{\theta}} - \psi_{\alpha_0, \beta_0, \theta_0}) + n^{-1/2}\mathbb{G}_n(\psi_{\alpha_0, \beta_0, \theta_0} - \psi_{\alpha_0, \beta_0, \theta_0}) \\ &= I + II + III + IV. \end{aligned}$$

By Condition ILAD(iii) (C.35) we have with probability at least  $1 - \Delta_n$  that  $|III| \leq \delta_n n^{-1/2}$ . We wish to show that

$$|II + (f_\epsilon \bar{\mathbb{E}}[v^2])(\tilde{\alpha} - \alpha_0)| \lesssim_P \delta_n n^{-1/2} + \delta_n |\tilde{\alpha} - \alpha_0|. \quad (\text{C.36})$$

Observe that

$$\begin{aligned} \Gamma(\alpha, \hat{\beta}, \hat{\theta}) &= \Gamma(\alpha, \beta_0, \theta_0) + \Gamma(\alpha, \hat{\beta}, \hat{\theta}) - \Gamma(\alpha, \beta_0, \theta_0) \\ &= \Gamma(\alpha, \beta_0, \theta_0) + \{\Gamma(\alpha, \hat{\beta}, \hat{\theta}) - \Gamma(\alpha, \beta_0, \theta_0) - \Gamma_2(\alpha, \beta_0, \theta_0)'(\hat{h} - h_0)\} + \Gamma_2(\alpha, \beta_0, \theta_0)'(\hat{h} - h_0). \end{aligned}$$

Since  $\Gamma(\alpha_0, \beta_0, \theta_0) = 0$ , by Taylor's theorem, there exists some point  $\tilde{\alpha}$  between  $\alpha_0$  and  $\alpha$  such that  $\Gamma(\alpha, \beta_0, \theta_0) = \Gamma_1(\tilde{\alpha}, \beta_0, \theta_0)(\alpha - \alpha_0)$ . By its definition, we have

$$\Gamma_1(\alpha, \beta, \theta) = -\bar{\mathbb{E}}[f_\epsilon(x'(\beta - \beta_0) + d(\alpha - \alpha_0))d(d - x'\theta)] = -\bar{\mathbb{E}}[f_\epsilon(x'(\beta - \beta_0) + d(\alpha - \alpha_0))d\{v - x'(\theta - \theta_0)\}].$$

Since  $f_\epsilon = f_\epsilon(0)$  and  $d_i = x_i'\theta_0 + v_i$  with  $\mathbb{E}[v_i | x_i] = 0$ , we have  $\Gamma_1(\alpha_0, \beta_0, \theta_0) = -f_\epsilon \bar{\mathbb{E}}[dv] = -f_\epsilon \bar{\mathbb{E}}[v^2]$ . Also

$$|\Gamma_1(\alpha, \beta_0, \theta_0) - \Gamma_1(\alpha_0, \beta_0, \theta_0)| \leq |\bar{\mathbb{E}}[\{f_\epsilon(0) - f_\epsilon(d(\alpha - \alpha_0))\}dv]| \leq C|\alpha - \alpha_0|\bar{\mathbb{E}}[d^2v].$$

Hence  $\Gamma_1(\tilde{\alpha}, \beta_0, \theta_0) = -f_\epsilon \bar{\mathbb{E}}[v^2] + O(1)|\tilde{\alpha} - \alpha_0|$ .

Observe that

$$\Gamma_2(\alpha, \beta, \theta) = \begin{pmatrix} -\bar{\mathbb{E}}[f_\epsilon(x'(\beta - \beta_0) + d(\alpha - \alpha_0))(d - x'\theta)x] \\ -\bar{\mathbb{E}}[(1/2 - 1\{y \leq x'\beta + d\alpha\})x] \end{pmatrix}.$$

Note that since  $\bar{\mathbb{E}}[f_\epsilon(0)(d - x'\theta_0)x] = f_\epsilon \bar{\mathbb{E}}[vx] = 0$  and  $\bar{\mathbb{E}}[(1/2 - 1\{y \leq x'\beta_0 + d\alpha_0\})x] = \bar{\mathbb{E}}[(1/2 - 1\{\epsilon \leq 0\})x] = 0$ , we have  $\Gamma_2(\alpha_0, \beta_0, \theta_0) = 0$ . Moreover,

$$\begin{aligned} |\Gamma_2(\alpha, \beta_0, \theta_0)'(\hat{h} - h_0)| &= |\{\Gamma_2(\alpha, \beta_0, \theta_0) - \Gamma_2(\alpha_0, \beta_0, \theta_0)\}'(\hat{h} - h_0)| \\ &\leq |\bar{\mathbb{E}}[\{f_\epsilon(d(\alpha - \alpha_0)) - f_\epsilon(0)\}vx'](\hat{\beta} - \beta_0)| \\ &\quad + |\bar{\mathbb{E}}[\{F(d(\alpha - \alpha_0)) - F(0)\}x'](\hat{\theta} - \theta_0)| \\ &\leq O(1)\{\|x'(\hat{\beta} - \beta_0)\|_{2,n} + \|x'(\hat{\theta} - \theta_0)\|_{2,n}\}|\alpha - \alpha_0| \\ &= O_P(\delta_n)|\alpha - \alpha_0|. \end{aligned}$$

Hence  $|\Gamma_2(\tilde{\alpha}, \beta_0, \theta_0)'(\hat{h} - h_0)| \lesssim_P \delta_n |\tilde{\alpha} - \alpha_0|$ .

Denote by  $\Gamma_{22}(\alpha, \beta, \theta)$  the Hessian matrix of  $\Gamma(\alpha, \beta, \theta)$  with respect to  $h = (\beta', \theta)'$ . Then

$$\Gamma_{22}(\alpha, \beta, \theta) = \begin{pmatrix} -\bar{\mathbb{E}}[f'_\epsilon(x'(\beta - \beta_0) + d(\alpha - \alpha_0))(d - x'\theta)xx'] & \bar{\mathbb{E}}[f_\epsilon(x'(\beta - \beta_0) + d(\alpha - \alpha_0))xx'] \\ \bar{\mathbb{E}}[f_\epsilon(x'(\beta - \beta_0) + d(\alpha - \alpha_0))xx'] & 0 \end{pmatrix},$$

so that

$$\begin{aligned} (\hat{h} - h_0)' \Gamma_{22}(\alpha, \beta, \theta) (\hat{h} - h_0) &\leq |(\hat{\beta} - \beta_0)' \bar{\mathbb{E}}[f'_\epsilon(x'(\beta - \beta_0) + d(\alpha - \alpha_0))(d - x'\theta)xx'] (\hat{\beta} - \beta_0)| \\ &\quad + 2|(\hat{\beta} - \beta_0)' \bar{\mathbb{E}}[f_\epsilon(x'(\beta - \beta_0) + d(\alpha - \alpha_0))xx'] (\hat{\theta} - \theta_0)| \\ &\leq C \{ \max_{1 \leq i \leq n} \mathbb{E}[|d_i - x'_i \theta| | x_i|] \|x'(\hat{\beta} - \beta_0)\|_{2,n}^2 + 2\|x'(\hat{\beta} - \beta_0)\|_{2,n} \cdot \|x'(\hat{\theta} - \theta_0)\|_{2,n} \}. \end{aligned}$$

Here  $|d_i - x'_i \theta| = |v_i - x'_i(\theta - \theta_0)| \leq |v_i| + |x'_i(\theta - \theta_0)|$ . Hence by Taylor's theorem together with ILAD(iii), we conclude that

$$|\Gamma(\check{\alpha}, \hat{\beta}, \hat{\theta}) - \Gamma(\check{\alpha}, \beta_0, \theta_0) - \Gamma_2(\check{\alpha}, \beta_0, \theta_0)'(\hat{h} - h_0)| \lesssim_P \delta_n n^{-1/2}.$$

This leads to the expansion in (C.36).

We now proceed to bound the fourth term. By Condition ILAD(iii) we have with probability at least  $1 - \Delta_n$  that  $|\check{\alpha} - \alpha_0| \leq \delta_n$ . Observe that

$$\begin{aligned} (\psi_{\alpha, \beta_0, \theta_0} - \psi_{\check{\alpha}, \hat{\beta}, \hat{\theta}})(y_i, d_i, x_i) &= (1\{y_i \leq x'_i \beta_0 + d_i \alpha_0\} - 1\{y_i \leq x'_i \hat{\beta} + d_i \check{\alpha}\})v_i \\ &= (1\{\epsilon_i \leq 0\} - 1\{\epsilon_i \leq d_i(\alpha - \alpha_0)\})v_i, \end{aligned}$$

so that  $|(\psi_{\alpha, \beta_0, \theta_0} - \psi_{\check{\alpha}, \hat{\beta}, \hat{\theta}})(y_i, d_i, x_i)| \leq 1\{|\epsilon_i| \leq \delta_n |d_i|\} |v_i|$  whenever  $|\alpha - \alpha_0| \leq \delta_n$ . Since the class of functions  $\{(y, d, x) \mapsto (\psi_{\alpha, \beta_0, \theta_0} - \psi_{\check{\alpha}, \hat{\beta}, \hat{\theta}})(y, d, x) : |\alpha - \alpha_0| \leq \delta_n\}$  is a VC subgraph class with VC index bounded by some constant independent of  $n$ , using (a version of) Theorem 2.14.1 in [30], we have

$$\sup_{|\alpha - \alpha_0| \leq \delta_n} |\mathbb{G}_n(\psi_{\alpha, \beta_0, \theta_0} - \psi_{\check{\alpha}, \hat{\beta}, \hat{\theta}})| \lesssim_P (\bar{\mathbb{E}}[1\{|\epsilon_i| \leq \delta_n |d_i|\} v_i^2])^{1/2} \lesssim_P \delta_n^{1/2}.$$

This implies that  $|IV| \lesssim_P \delta_n^{1/2} n^{-1/2}$ .

Combining these bounds on II, III and IV, we have the following stochastic expansion

$$\mathbb{E}_n[\psi_{\check{\alpha}, \hat{\beta}, \hat{\theta}}(y, d, x)] = -(f_\epsilon \bar{\mathbb{E}}[v^2])(\check{\alpha} - \alpha_0) + \mathbb{E}_n[\psi_{\alpha_0, \beta_0, \theta_0}(y, d, x)] + O_P(\delta_n^{1/2} n^{-1/2}) + O_P(\delta_n) |\check{\alpha} - \alpha_0|.$$

Let  $\alpha^* = \alpha_0 + (f_\epsilon \bar{\mathbb{E}}[v^2])^{-1} \mathbb{E}_n[\psi_{\alpha_0, \beta_0, \theta_0}(y, d, x)]$ . Then  $\alpha^* \in \mathcal{A}$  with probability  $1 - o(1)$  since  $|\alpha^* - \alpha_0| \lesssim_P n^{-1/2}$ . It is not difficult to see that the above stochastic expansion holds with  $\check{\alpha}$  replaced by  $\alpha^*$ , so that

$$\mathbb{E}_n[\psi_{\alpha^*, \hat{\beta}, \hat{\theta}}(y, d, x)] = -(f_\epsilon \bar{\mathbb{E}}[v^2])(\alpha^* - \alpha_0) + \mathbb{E}_n[\psi_{\alpha_0, \beta_0, \theta_0}(y, d, x)] + O_P(\delta_n^{1/2} n^{-1/2}) = O_P(\delta_n^{1/2} n^{-1/2}).$$

Therefore,  $|\mathbb{E}_n[\psi_{\check{\alpha}, \hat{\beta}, \hat{\theta}}(y, d, x)]| \leq |\mathbb{E}_n[\psi_{\alpha^*, \hat{\beta}, \hat{\theta}}(y, d, x)]| = O_P(\delta_n^{1/2} n^{-1/2})$ , so that

$$(f_\epsilon \bar{\mathbb{E}}[v^2])(\check{\alpha} - \alpha_0) = \mathbb{E}_n[\psi_{\alpha_0, \beta_0, \theta_0}(y, d, x)] + O_P(\delta_n^{1/2} n^{-1/2}),$$

which immediately implies that  $\sigma_n^{-1} \sqrt{n}(\check{\alpha} - \alpha_0) \rightsquigarrow N(0, 1)$  since by the Lyapunov CLT,

$$(\bar{\mathbb{E}}[v^2]/4)^{-1/2} \sqrt{n} \mathbb{E}_n[\psi_{\alpha_0, \beta_0, \theta_0}(y, d, x)] \rightsquigarrow N(0, 1).$$

Part 2. (Proof for the second assertion). First consider the denominator of  $L_n(\alpha_0)$ . We have that

$$\begin{aligned} |\mathbb{E}_n[\widehat{v}^2] - \mathbb{E}_n[v^2]| &= |\mathbb{E}_n[(\widehat{v} - v)(\widehat{v} + v)]| \leq \|\widehat{v} - v\|_{2,n} \|\widehat{v} + v\|_{2,n} \\ &\leq \|x'(\widehat{\theta} - \theta_0)\|_{2,n} (2\|v\|_{2,n} + \|x'(\widehat{\theta} - \theta_0)\|_{2,n}) \lesssim_P \delta_n, \end{aligned}$$

where we have used the fact that  $\|v\|_{2,n} \lesssim_P (\bar{\mathbb{E}}[v^2])^{1/2} = O(1)$  (which is guaranteed by ILAD(i)).

Next consider the numerator of  $L_n(\alpha_0)$ . Since  $\bar{\mathbb{E}}[\psi_{\alpha_0, \beta_0, \theta_0}(y_i, d_i, x_i)] = 0$  we have

$$\mathbb{E}_n[\psi_{\alpha_0, \widehat{\beta}, \widehat{\theta}}(y, d, x)] = n^{-1/2} \mathbb{G}_n(\psi_{\alpha_0, \widehat{\beta}, \widehat{\theta}} - \psi_{\alpha_0, \beta_0, \theta_0}) + \Gamma(\alpha_0, \widehat{\beta}, \widehat{\theta}) + \mathbb{E}_n[\psi_{\alpha_0, \beta_0, \theta_0}(y, d, x)].$$

By Condition ILAD(iii) and the previous calculation, we have

$$|\mathbb{G}_n(\psi_{\alpha_0, \widehat{\beta}, \widehat{\theta}} - \psi_{\alpha_0, \beta_0, \theta_0})| \lesssim_P \delta_n \text{ and } |\Gamma(\alpha_0, \widehat{\beta}, \widehat{\theta})| \lesssim_P \delta_n n^{-1/2}.$$

Therefore, using the simple identity that  $nA_n^2 = nB_n^2 + n(A_n - B_n)^2 + 2nB_n(A_n - B_n)$  with

$$A_n = \mathbb{E}_n[\psi_{\alpha_0, \widehat{\beta}, \widehat{\theta}}(y, d, x)] \text{ and } B_n = \mathbb{E}_n[\psi_{\alpha_0, \beta_0, \theta_0}(y, d, x)] \lesssim_P (\bar{\mathbb{E}}[v^2])n^{-1/2},$$

we have

$$nL_n(\alpha_0) = \frac{4n|\mathbb{E}_n[\psi_{\alpha_0, \widehat{\beta}, \widehat{\theta}}(y, d, x)]|^2}{\mathbb{E}_n[\widehat{v}^2]} = \frac{4n|\mathbb{E}_n[\psi_{\alpha_0, \beta_0, \theta_0}(y, d, x)]|^2}{\bar{\mathbb{E}}[v^2]} + O_P(\delta_n)$$

since  $\bar{\mathbb{E}}[v^2] \geq c$  is bounded away from zero. The result then follows since

$$(\bar{\mathbb{E}}[v^2]/4)^{-1/2} \sqrt{n} \mathbb{E}_n[\psi_{\alpha_0, \beta_0, \theta_0}(y, d, x)] \rightsquigarrow N(0, 1).$$

□

**Comment C.2** (On 1-step procedure). An inspection of the proof leads to the following stochastic expansion:

$$\begin{aligned} \mathbb{E}_n[\psi_{\widehat{\alpha}, \widehat{\beta}, \widehat{\theta}}(y, d, x)] &= -(f_\epsilon \bar{\mathbb{E}}[v^2])(\widehat{\alpha} - \alpha_0) + \mathbb{E}_n[\psi_{\alpha_0, \beta_0, \theta_0}(y, d, x)] \\ &\quad + O_P(\delta_n^{1/2} n^{-1/2} + \delta_n n^{-1/4} |\widehat{\alpha} - \alpha_0| + |\widehat{\alpha} - \alpha_0|^2), \end{aligned}$$

where  $\widehat{\alpha}$  is any consistent estimator of  $\alpha_0$ . Hence provided that  $|\widehat{\alpha} - \alpha_0| = o_P(n^{-1/4})$ , the remainder term in the above expansion is  $o_P(n^{-1/2})$ , and the 1-step estimator  $\check{\alpha}$  defined by

$$\check{\alpha} = \widehat{\alpha} + (\mathbb{E}_n[f_\epsilon \widehat{v}^2])^{-1} \mathbb{E}_n[\psi_{\widehat{\alpha}, \widehat{\beta}, \widehat{\theta}}(y, d, x)]$$

has the following stochastic expansion:

$$\begin{aligned} \check{\alpha} &= \widehat{\alpha} + \{f_\epsilon \bar{\mathbb{E}}[v^2] + o_P(n^{-1/4})\}^{-1} \{- (f_\epsilon \bar{\mathbb{E}}[v^2])(\widehat{\alpha} - \alpha_0) + \mathbb{E}_n[\psi_{\alpha_0, \beta_0, \theta_0}(y, d, x)] + o_P(n^{-1/2})\} \\ &= \alpha_0 + (f_\epsilon \bar{\mathbb{E}}[v^2])^{-1} \mathbb{E}_n[\psi_{\alpha_0, \beta_0, \theta_0}(y, d, x)] + o_P(n^{-1/2}), \end{aligned}$$

so that  $\sigma_n^{-1} \sqrt{n}(\check{\alpha} - \alpha_0) \rightsquigarrow N(0, 1)$ .

## APPENDIX D. PROOFS FOR SECTION 3

**D.1. Proof of Theorem 1.** The proof of Theorem 1 uses the properties of Post- $\ell_1$ -LAD and Post-Lasso. We will collect these properties together with required regularity conditions in Appendix F. In what follows we denote  $\phi_{\min}(m) := \phi_{\min}(m, \mathbb{E}_n[\tilde{x}\tilde{x}'])$ ,  $\phi_{\max}(m) := \phi_{\max}(m, \mathbb{E}_n[\tilde{x}\tilde{x}'])$ ,  $\bar{\phi}_{\min}(m) := \phi_{\min}(m, \mathbb{E}[\tilde{x}\tilde{x}'])$ ,  $\bar{\phi}_{\max}(m) := \phi_{\max}(m, \mathbb{E}[\tilde{x}\tilde{x}'])$ . The proof focuses on Algorithm 1, while the proof for Algorithm 2 is deferred to the Supplementary Appendix, since it is basically the same proof.

We will verify Condition ILAD and the desired result then follows from Lemma 1 and noting that  $|\mathbb{E}_n[\mathbb{E}[v^2 | x]] - \mathbb{E}[v^2]| \lesssim_P \delta_n$  and  $\mathbb{E}[v^2]$  is bounded away from zero under Condition 1.

The assumptions on the error density  $f_e(\cdot)$  in Condition ILAD(i) are assumed in Condition 1(iv). The moment conditions on  $d_i$  and  $v_i$  in Condition ILAD(i) are assumed in Condition 1(ii).

Because Condition 1(v) and (vi), by Lemma 4 we have for some  $\tilde{\ell}_n \rightarrow \infty$  we have  $\kappa' \leq \phi_{\min}(\tilde{\ell}_n) \leq \phi_{\max}(\tilde{\ell}_n) \leq \kappa''$  with probability  $1 - \Delta_n$ . In turn,  $\kappa_c$  is bounded away from zero with probability  $1 - \Delta_n$  for  $n$  sufficiently large, see [9].

Step 1 relies on Post- $\ell_1$ -LAD. By assumption with probability  $1 - \Delta_n$  we have  $\hat{s} = \|\tilde{\beta}\|_0 \leq Cs$ . Thus,  $\phi_{\min}(\hat{s} + s)$  is bounded away from zero since  $\hat{s} + s \leq \ell_n s$  for large enough  $n$  with probability  $1 - \Delta_n$ . Moreover, Condition PLAD in Appendix F is implied by Condition 1. The required side condition of Lemma 6 is satisfied by relations (H.52) and (H.53). By Lemma 6 we have  $|\hat{\alpha} - \alpha_0| \lesssim_P \sqrt{s \log(p \vee n)/n} \leq o(1) \log^{-1} n$  under  $s^3 \log^3(p \vee n) \leq \delta_n n$ . Note that this implies  $\{\alpha : |\alpha - \alpha_0| \leq n^{-1/2} \log n\} \subset \mathcal{A}$  (with probability  $1 - o(1)$ ) which is required in ILAD(ii) and the (shrinking) definition of  $\mathcal{A}$  establishes the initial rate of ILAD(iv). By Lemma 7 in Appendix F we have  $\|x'(\tilde{\beta} - \beta_0)\|_{2,n} \lesssim_P \sqrt{s \log(n \vee p)/n}$  since the required side condition holds. Indeed, for  $\tilde{x}_i = (d_i, x'_i)'$  and  $\delta = (\delta_d, \delta'_x)'$ , because  $\phi_{\min}(\hat{s} + s)$  is bounded away from zero,  $\phi_{\max}(\hat{s} + s)$  is bounded from above, and the fact that  $\mathbb{E}_n[|d|^3] \lesssim_P \mathbb{E}[|d|^3] = O(1)$ ,

$$\begin{aligned} \inf_{\|\delta\|_0 \leq s + Cs} \frac{\|\tilde{x}'\delta\|_{2,n}^3}{\mathbb{E}_n[|\tilde{x}'\delta|^3]} &\geq \inf_{\|\delta\|_0 \leq s + Cs} \frac{\{\phi_{\min}(s + Cs)\}^{3/2} \|\delta\|^3}{4\mathbb{E}_n[|x'\delta_x|^3] + 4|\delta_d|^3 \mathbb{E}_n[|d|^3]} \\ &\geq \inf_{\|\delta\|_0 \leq s + Cs} \frac{\{\phi_{\min}(s + Cs)\}^{3/2} \|\delta\|^3}{4K_x \|\delta_x\|_1 \phi_{\max}(s + Cs) \|\delta_x\|^2 + 4\|\delta\|^3 \mathbb{E}_n[|d|^3]} \\ &\geq \frac{\{\phi_{\min}(s + Cs)\}^{3/2}}{4K_x \sqrt{s + Cs} \phi_{\max}(s + Cs) + 4\mathbb{E}_n[|d|^3]} \gtrsim_P \frac{1}{K_x \sqrt{s}}. \end{aligned}$$

Therefore, since  $K_x^2 s^2 \log^2(p \vee n) \leq \delta_n n$  and  $\lambda \lesssim \sqrt{n \log(p \vee n)}$  we have

$$\frac{n \sqrt{\phi_{\min}(s + Cs)}}{\lambda \sqrt{s + \sqrt{sn \log(p \vee n)}}} \inf_{\|\delta\|_0 \leq s + Cs} \frac{\|\tilde{x}'\delta\|_{2,n}^3}{\mathbb{E}_n[|\tilde{x}'\delta|^3]} \gtrsim_P \frac{\sqrt{n}}{K_x s \log(p \vee n)} \rightarrow \infty.$$

Step 2 relies on Post-Lasso. Condition HL in Appendix F is implied by Condition 1. Indeed,  $\mathbb{E}_n[x_j^2]$  is bounded away from zero and from above with probability  $1 - o(1)$  by Conditions 1 (v) and (vi) and Lemma 4. Next note that by  $c \leq \mathbb{E}[v^2 | x]$  and  $\mathbb{E}[|v|^3 | x] \leq C$  so that  $\max_{1 \leq j \leq p} \{\mathbb{E}_n[|x_j|^3 \mathbb{E}[|v|^3 | x]]\}^{1/3} / \{\mathbb{E}_n[x_j^2 \mathbb{E}[v^2 | x]]\}^{1/2} \leq K_x^{1/3} C'$  with probability  $1 - \Delta_n$ . Thus, Condition HL(ii) holds under  $K_x^{1/3} \sqrt{\log(p \vee n)} = o(n^{1/6})$ . Condition HL(iii) follows by Lemma 2 applied twice with  $\zeta_i = v_i$  and  $\tilde{\zeta}_i = d_i$  under the condition that  $K_x^4 \log p \leq \delta_n n$ . By Lemma 9 in Appendix F we have  $\|x'(\tilde{\theta} - \theta_0)\|_{2,n} \lesssim_P \sqrt{s \log(n \vee p)/n}$  and  $\|\tilde{\theta}\|_0 \lesssim s$  with probability  $1 - o(1)$ .

The rates established above for  $\tilde{\theta}$  and  $\tilde{\beta}$  imply (C.34) in ILAD(iii) since by Condition 1(ii)  $\mathbb{E}[|v_i|] \leq (\mathbb{E}[v_i^2])^{1/2} = O(1)$  and  $\max_{1 \leq i \leq n} |x'_i(\tilde{\theta} - \theta_0)| \lesssim_P K_x \sqrt{s^2 \log(p \vee n)/n} = o(1)$ .

We now verify the last requirement in Condition ILAD(iii). Consider the following class of functions

$$\mathcal{F}_s = \{(y, d, x) \mapsto 1\{y \leq x'\beta + d\alpha\} : \alpha \in \mathbb{R}, \|\beta\|_0 \leq Cs\},$$

which is the union of  $\binom{p}{Cs}$  VC-subgraph classes of functions with VC indices bounded by  $C's$ . Hence

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

Likewise, consider the following class of functions  $\mathcal{G}_{s,r} = \{(y, d, x) \mapsto x'\theta : \|\theta\|_0 \leq Cs, \|x'\theta\|_{2,n} \leq r\}$ . Then

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

where  $G_{s,r}(y, d, x) = \max_{\|\theta\|_0 \leq Cs, \|x'\theta\|_{2,n} \leq r} |x'\theta|$ .

Note that

$$\sup_{\alpha \in \mathcal{A}} |\mathbb{G}_n(\psi_{\alpha, \tilde{\beta}, \tilde{\theta}} - \psi_{\alpha, \beta_0, \theta_0})| \leq \sup_{\alpha \in \mathcal{A}} |\mathbb{G}_n(\psi_{\alpha, \tilde{\beta}, \tilde{\theta}} - \psi_{\alpha, \tilde{\beta}, \theta_0})| \quad (\text{D.37})$$

$$+ \sup_{\alpha \in \mathcal{A}} |\mathbb{G}_n(\psi_{\alpha, \tilde{\beta}, \theta_0} - \psi_{\alpha, \beta_0, \theta_0})|. \quad (\text{D.38})$$

Consider to bound (D.37). Observe that

$$\psi_{\alpha, \beta, \theta}(y_i, d_i, x_i) - \psi_{\alpha, \beta_0, \theta_0}(y_i, d_i, x_i) = -(1/2 - 1\{y_i \leq x'_i \beta + d_i \alpha\}) x'_i (\theta - \theta_0),$$

and consider the class of functions  $\mathcal{H}_{s,r}^1 = \{(y, d, x) \mapsto (1/2 - 1\{y \leq x'\beta + d\alpha\})x'(\theta - \theta_0) : \alpha \in \mathbb{R}, \|\beta\|_0 \leq Cs, \|\theta\|_0 \leq Cs, \|x'(\theta - \theta_0)\|_{2,n} \leq r\}$  with  $r \lesssim \sqrt{s \log(p \vee n)/n}$ . Then by Lemma 11 together with the above entropy calculations (and some straightforward algebras), we have

$$\sup_{g \in \mathcal{H}_{s,r}^1} |\mathbb{G}_n(g)| \lesssim_P \sqrt{s \log(p \vee n)} \sqrt{s \log(p \vee n)/n} = o_P(1),$$

where  $s^2 \log^2(p \vee n) \leq \delta_n n$  is used. Since  $\|x'(\tilde{\theta} - \theta_0)\|_{2,n} \lesssim_P \sqrt{s \log(n \vee p)/n}$  and  $\|\tilde{\beta}\|_0 \vee \|\tilde{\theta}\|_0 \lesssim s$  with probability  $1 - o(1)$ , we conclude that (D.37) =  $o_P(1)$ .

Lastly consider to bound (D.38). Observe that

$$\psi_{\alpha, \beta, \theta_0}(y_i, d_i, x_i) - \psi_{\alpha, \beta_0, \theta_0}(y_i, d_i, x_i) = -(1\{y_i \leq x'_i \beta + d_i \alpha\} - 1\{y_i \leq x'_i \beta_0 + d_i \alpha\}) v_i,$$

where  $v_i = d_i - x'_i \theta_0$ , and consider the class of functions  $\mathcal{H}_{s,r}^2 = \{(y, d, x) \mapsto (1\{y \leq x'\beta + d\alpha\} - 1\{y \leq x'\beta_0 + d\alpha\})(d - x'\theta_0) : \alpha \in \mathbb{R}, \|\beta\|_0 \leq Cs, \|x'(\beta - \beta_0)\|_{2,n} \leq r\}$  with  $r \lesssim \sqrt{s \log(p \vee n)/n}$ . Then by Lemma 11 together with the above entropy calculations (and some straightforward algebras), we have

$$\sup_{g \in \mathcal{H}_{s,r}^2} |\mathbb{G}_n(g)| \lesssim_P \sqrt{s \log(p \vee n)} \sup_{g \in \mathcal{H}_{s,r}^2} \sqrt{\mathbb{E}_n[g(y, d, x)^2] \vee \mathbb{E}[g(y, d, x)^2]}.$$

Here we have

$$\mathbb{E}[g(y, d, x)^2] \leq C \|x'(\beta - \beta_0)\|_{2,n} (\mathbb{E}[v^4])^{1/2} \lesssim \sqrt{s \log(p \vee n)/n}.$$

On the other hand,

$$\sup_{g \in \mathcal{H}_{s,r}^2} \mathbb{E}_n[g(y, d, x)^2] \leq n^{-1/2} \sup_{g \in \mathcal{H}_{s,r}^2} \mathbb{G}_n(g^2) + \sup_{g \in \mathcal{H}_{s,r}^2} \mathbb{E}[g(y, d, x)^2], \quad (\text{D.39})$$

and apply Lemma 11 to the first term on the right side of (D.39). Then we have

$$\begin{aligned} \sup_{g \in \mathcal{H}_{s,r}^2} \mathbb{G}_n(g^2) &\lesssim_P \sqrt{s \log(p \vee n)} \sup_{g \in \mathcal{H}_{s,r}^2} \sqrt{\mathbb{E}_n[g(y, d, x)^4] \vee \mathbb{E}[g(y, d, x)^4]} \\ &\lesssim \sqrt{s \log(p \vee n)} \sqrt{\mathbb{E}_n[v^4] \vee \mathbb{E}[v^4]} \lesssim_P \sqrt{s \log(p \vee n)} \sqrt{\mathbb{E}[v^4]}. \end{aligned}$$

Since  $\|x'(\tilde{\beta} - \beta_0)\|_{2,n} \lesssim_P \sqrt{s \log(n \vee p)/n}$  and  $\|\tilde{\beta}\|_0 \leq Cs$  with probability  $1 - \Delta_n$ , we conclude that

$$(\text{D.38}) \lesssim_P \sqrt{s \log(p \vee n)} (s \log(p \vee n)/n)^{1/4} = o(1),$$

where  $s^3 \log^3(p \vee n) \leq \delta_n n$  is used.

**D.2. Auxiliary Technical Results for Proofs of Section 3.** In this section we collect two auxiliary technical results. Their proofs are given in the supplementary appendix.

**Lemma 2.** *Let  $x_1, \dots, x_n$  be non-stochastic vectors in  $\mathbb{R}^p$  with  $\max_{1 \leq i \leq n} \|x_i\|_\infty \leq K_x$ . Let  $\zeta_1, \dots, \zeta_n$  be independent random variables such that  $\mathbb{E}[|\zeta_i|^q] < \infty$  for some  $q \geq 4$ . Then with probability at least  $1 - 8\tau$ ,*

$$\max_{1 \leq j \leq p} |(\mathbb{E}_n - \mathbb{E})[x_j^2 \zeta^2]| \leq 4 \sqrt{\frac{\log(2p/\tau)}{n}} K_x^2 (\mathbb{E}[|\zeta|^q]/\tau)^{4/q}.$$

**Lemma 3.** *Let  $T = \text{supp}(\beta_0)$ ,  $|T| = \|\beta_0\|_0 \leq s$  and  $\|\hat{\beta}_{T^c}\|_1 \leq \mathbf{c} \|\hat{\beta}_T - \beta_0\|_1$ . Moreover, let  $\hat{\beta}^{(2m)}$  denote the vector formed by the largest  $2m$  components of  $\hat{\beta}$  in absolute value and zero in the remaining components. Then for  $m \geq s$  we have that  $\hat{\beta}^{(2m)}$  satisfies*

$$\|x'(\hat{\beta}^{(2m)} - \beta_0)\|_{2,n} \leq \|x'(\hat{\beta} - \beta_0)\|_{2,n} + \sqrt{\phi_{\max}(m)/m} \mathbf{c} \|\hat{\beta}_T - \beta_0\|_1,$$

where  $\phi_{\max}(m)/m \leq 2\phi_{\max}(s)/s$  and  $\|\hat{\beta}_T - \beta_0\|_1 \leq \sqrt{s} \|x'(\hat{\beta} - \beta_0)\|_{2,n}/\kappa_{\mathbf{c}}$ .

**Lemma 4.** *Under Condition 1, for  $\tilde{x}_i = (d_i, x_i)'$ , there is  $\tilde{\ell}_n \rightarrow \infty$  such that with probability  $1 - o(1)$  we have*

$$\sup_{\|\delta\|_0 \leq \tilde{\ell}_n s} \left| \frac{\|\tilde{x}'\delta\|_{2,n}}{\|\tilde{x}'\delta\|_{P,2}} - 1 \right| = o(1).$$

## APPENDIX E. PROOFS FOR SECTION 4

**E.1. A Maximal Inequality.** For a class of measurable functions  $\mathcal{F}$  equipped with the envelope  $F = \sup_{f \in \mathcal{F}} |f|$ , let  $N(\epsilon, \mathcal{F}, \|\cdot\|_{Q,2})$  denote the  $\epsilon$ -covering number of the class of functions  $\mathcal{F}$  with respect to the  $L^2(Q)$  seminorm  $\|\cdot\|_{Q,2}$ , where  $Q$  is finitely discrete, and let  $\text{ent}(\epsilon, \mathcal{F}) = \log \sup_Q N(\epsilon \|F\|_{Q,2}, \mathcal{F}, \|\cdot\|_{Q,2})$  denote the uniform covering entropy.

**Lemma 5** ([10]). *Let  $\mathcal{F}$  be a suitably measurable class of functions. Suppose that  $F = \sup_{f \in \mathcal{F}} |f|$  with  $\|F\|_{Q,q} < \infty$  for some  $q \geq 2$ . Let  $M_{P,q} = \{\mathbb{E}|\max_{i \leq n} F(w_i)|^q\}^{1/q}$ . Suppose that there exist constants*



$a \geq e$  and  $s \geq 1$  such that  $\text{ent}(\varepsilon, \mathcal{F}) \leq s\{\log a + \log(1/\varepsilon)\}$ ,  $0 < \varepsilon \leq 1$ . Then for  $\sigma > 0$  denoting a positive constant such that  $\sup_{f \in \mathcal{F}} \|f\|_{P,2} \leq \sigma \leq \|F\|_{P,2}$ :

$$\mathbb{E} \sup_{f \in \mathcal{F}} |\mathbb{G}_n(f)| \lesssim \sqrt{s\sigma^2 \log(a\|F\|_{P,2}/\sigma)} + n^{-1/2} s M_{P,2} \log(a\|F\|_{P,2}/\sigma).$$

Moreover, for every  $t \geq 1$  and all  $\mu > 0$ , with probability not less than  $1 - t^{-q/2}$ ,

$$\sup_{f \in \mathcal{F}} |\mathbb{G}_n(f)| \leq (1 + \mu) \mathbb{E} \sup_{f \in \mathcal{F}} |\mathbb{G}_n(f)| + K(q) \left\{ (\sigma + n^{-1/2} M_{P,q}) \sqrt{t} + \mu^{-1} n^{-1/2} M_{P,2} t \right\},$$

where  $K(q) > 0$  is a constant depending only on  $q$ .

**E.2. Proof of Theorem 2.** It suffices to establish the result under any sequence  $P = P_n \in \mathcal{P}_n$ . We shall suppress the dependency of  $P$  on  $n$  in the proof. We use  $C$  as a generic constant that may differ in each appearance, but that does not depend on the sequence  $P \in \mathcal{P}_n$ . Let

$$\mathbb{B}(w) = \max_{k \leq K} \sup_{\nu_\ell \in \mathcal{A}_\ell \times T_\ell, \ell \in \mathcal{L}} |\partial_{\nu_k} \mathbb{E}\{\psi_\ell(w, \nu) \mid z_\ell\}|, \quad \rho_n = n^{-1/2} [\{s \log(a_n)\}^{1/2} + n^{-1/2} s n^{\frac{1}{q}} \log(a_n)].$$

Step 1. (Preliminary Rate). In this step we claim that w.p.  $1 - o(1)$ ,  $\max_{\ell \in \mathcal{L}} |\widehat{\alpha}_\ell - \alpha_\ell| \leq C \rho_n$  for some constant  $C$ , independent of  $n$ . By definition  $|\mathbb{E}_n \psi_\ell\{w_\ell, \widehat{\alpha}_\ell, \widehat{h}_\ell(z_\ell)\}| \leq \inf_{\alpha \in \mathcal{A}_\ell} |\mathbb{E}_n \psi_\ell\{w_\ell, \alpha, \widehat{h}_\ell(z_\ell)\}| + \epsilon_n$  for each  $\ell \in \mathcal{L}$ , which implies via triangle inequality and Step 2 that w.p.  $1 - o(1)$  uniformly in  $\ell \in \mathcal{L}$ ,

$$|\mathbb{E}[\psi_\ell\{w_\ell, \alpha, h_\ell(z_\ell)\}]| \Big|_{\alpha = \widehat{\alpha}_\ell} \leq \epsilon_n + 2I_1 + 2I_2 \lesssim \rho_n, \quad (\text{E.40})$$

where we define  $I_1$  and  $I_2$  in Step 2. The second inequality in (E.40) is by Step 2 and by the assumption  $\epsilon_n = o(b_n^{-1} n^{-1/2})$ . Since by Condition 2  $2^{-1}\{|\Gamma_\ell(\widehat{\alpha}_\ell - \alpha_\ell)| \vee c_0\}$  is weakly smaller than the left side of (E.40) and  $\inf_{n \geq 1, \ell \in \mathcal{L}} |\Gamma_\ell| > c$ , conclude  $\max_{\ell \in \mathcal{L}} |\widehat{\alpha}_\ell - \alpha_\ell| \lesssim (\rho_n/c) \lesssim \rho_n$  w. p.  $1 - o(1)$ .

Step 2. (Define and bound  $I_1$  and  $I_2$ ) We claim that w. p.  $1 - o(1)$ :

$$\begin{aligned} I_1 &= \sup_{\alpha \in \mathcal{A}_\ell, \ell \in \mathcal{L}} |\mathbb{E}_n \psi_\ell\{w_\ell, \alpha, \widehat{h}_\ell(z_\ell)\} - \mathbb{E}_n \psi_\ell\{w_\ell, \alpha, h_\ell(z_\ell)\}| \lesssim \rho_n, \\ I_2 &= \sup_{\alpha \in \mathcal{A}_\ell, \ell \in \mathcal{L}} |\mathbb{E}_n \psi_\ell\{w_\ell, \alpha, h_\ell(z_\ell)\} - \mathbb{E} \psi_\ell\{w_\ell, \alpha, h_\ell(z_\ell)\}| \lesssim \rho_n. \end{aligned}$$

To show this, we can bound  $I_1 \leq I_{1a} + I_{1b}$  and  $I_2 \leq I_{2a}$ , where w. p.  $1 - o(1)$ ,

$$\begin{aligned} I_{1a} &= \sup_{\alpha \in \mathcal{A}_\ell, \ell \in \mathcal{L}, h \in \mathcal{H}_\ell \cup \{h_\ell\}} |\mathbb{E}_n \psi_\ell\{w_\ell, \alpha, h(z_\ell)\} - \mathbb{E} \psi_\ell\{w_\ell, \alpha, h(z_\ell)\}| \lesssim \rho_n, \\ I_{1b} &= \sup_{\alpha \in \mathcal{A}_\ell, \ell \in \mathcal{L}, h \in \mathcal{H}_\ell \cup \{h_\ell\}} |\mathbb{E} \psi_\ell\{w_\ell, \alpha, h(z_\ell)\} - \mathbb{E} \psi_\ell\{w_\ell, \alpha, h_\ell(z_\ell)\}| \lesssim \rho_n. \end{aligned}$$

The latter bounds hold by the following arguments.

By Taylor's expansion and triangle inequality, for  $\bar{h}_{\ell\alpha}(z_\ell)$  on a line connecting  $h(z_\ell)$  and  $h_\ell(z_\ell)$ ,

$$\begin{aligned} I_{1b} &\leq \sum_{m=1}^M \sup_{\alpha \in \mathcal{A}_\ell, \ell \in \mathcal{L}, h \in \mathcal{H}_\ell} |\mathbb{E} [\partial_{t_m} \mathbb{E} [\psi_\ell\{w_\ell, \alpha, \bar{h}_{\ell\alpha}(z_\ell)\} \mid z_\ell] \{h_m(z_\ell) - h_{\ell m}(z_\ell)\}]| \\ &\leq M \max_{m \leq M} \|\mathbb{B}\|_{P,2} \|h_m - h_{\ell m}\|_{P,2} \lesssim \rho_n, \end{aligned}$$

where the penultimate inequality holds by definition of  $\mathbf{B}$  and Holder's inequality, and the final inequality holds since by Condition 2  $\|\mathbf{B}\|_{P,2} \leq C$  and by Condition 3  $M$  is fixed and  $\|h_m - h_{\ell m}\|_{P,2} \lesssim \rho_n$  uniformly for all  $h \in \mathcal{H}_\ell$  and all  $m \leq M$ .

To bound  $I_{1a}$  we apply Lemma 5 to the empirical process indexed by  $\mathcal{F}' = \cup_{\ell \in \mathcal{L}} \mathcal{F}_\ell$  with the envelope  $F' \leq F$ , where  $\mathcal{F}_\ell = [\psi_\ell\{w, \alpha, h(z_\ell)\}, \alpha \in \mathcal{A}, h \in \mathcal{H}_\ell \cup \{h_\ell\}]$ . We note that  $\text{ent}(\varepsilon, \mathcal{F}') \lesssim \log(p_1) + \max_{\ell \in \mathcal{L}} \text{ent}(\varepsilon, \mathcal{F}_\ell)$ . Since  $\text{ent}(\varepsilon, \mathcal{H}_{\ell m}) \lesssim s \log(p/\varepsilon)$  by  $\mathcal{H}_{\ell m}$  consisting of  $p$  choose at most  $Cs$  VC subgraph classes, we concluded by Condition 3 that  $\text{ent}(\varepsilon, \mathcal{F}_\ell) \lesssim \{\log(e/\varepsilon) + Ms \log(a_n/\varepsilon)\}$ . Thus, recalling that  $a_n = \max(p_1, p, n, e)$ ,  $\text{ent}(\varepsilon, \mathcal{F}') \lesssim s \log(a_n/\varepsilon)$  and  $1 \lesssim \sigma^2 \leq \|F\|_{P,2}$ . By Lemma 5 with  $t = \log n$ , we conclude that there exists a constant  $C > 0$  such that with probability  $1 - o(1)$ ,

$$I_{1a} \leq \sup_{f \in \mathcal{F}'} |\mathbb{G}_n(f)| \leq Cn^{-1/2} \left( \|F\|_{P,2} \sqrt{s \log a_n} + n^{-1/2} s n^{\frac{1}{q}} \|F\|_{P,q} \log a_n \right) \lesssim \rho_n.$$

To arrive at the conclusion, we used the assumption that  $\|F\|_{P,q} \leq C$ ,  $a_n \geq n$ , and  $s \geq 1$ , and the elementary inequality  $M_{P,2} \leq n^{\frac{1}{q}} \|F\|_{P,q}$ .

Step 3. (Linearization) By definition

$$\sqrt{n} |\mathbb{E}_n \psi_\ell\{w_\ell, \hat{\alpha}_\ell, \hat{h}_\ell(z_\ell)\}| \leq \inf_{\alpha \in \mathcal{A}_\ell} \sqrt{n} |\mathbb{E}_n \psi_\ell\{w_\ell, \alpha, \hat{h}_\ell(z_\ell)\}| + \epsilon_n n^{1/2}.$$

Application of Taylor's theorem and the triangle inequality gives that with probability  $1 - o(1)$

$$\begin{aligned} & \max_{\ell \in \mathcal{L}} \left| \sqrt{n} \mathbb{E}_n \psi_\ell\{w, \alpha_\ell, h_\ell(z_\ell)\} + \Gamma_\ell \sqrt{n} (\hat{\alpha}_\ell - \alpha_\ell) + \Delta_\ell (\hat{h}_\ell - h_\ell) \right| \\ & \leq \epsilon_n \sqrt{n} + \max_{\ell \in \mathcal{L}} \left[ \inf_{\alpha \in \mathcal{A}_\ell} \sqrt{n} |\mathbb{E}_n \psi_\ell\{w_\ell, \alpha, \hat{h}_\ell(z_\ell)\}| + |II_1(\ell)| + |II_2(\ell)| \right] = o(b_n^{-1}), \end{aligned}$$

where  $II_1$  and  $II_2$  are defined in Step 4; the  $o(b_n^{-1})$  bound follows from Step 4, the assumption  $\epsilon_n \sqrt{n} = o(b_n^{-1})$ , and Step 5; and

$$\Delta_\ell (\hat{h}_\ell - h_\ell) = \sum_{m=1}^M \sqrt{n} \mathbb{E} \partial_{t_m} \mathbb{E} [\psi_\ell\{w_\ell, \alpha_\ell, h_\ell(z_\ell)\} \mid z_\ell] \{h_m(z_\ell) - h_{\ell m}(z_\ell)\} \Big|_{h_m = \hat{h}_m} = 0,$$

where the last equality occurs because of the orthogonality condition (3.18). Conclude using Condition 2 that with probability  $1 - o(1)$ :

$$\max_{\ell \in \mathcal{L}} |\Gamma_\ell^{-1} \sqrt{n} \mathbb{E}_n \psi_\ell\{w, \alpha_\ell, h_\ell(z_\ell)\} + \sqrt{n} (\hat{\alpha}_\ell - \alpha_\ell)| \leq o(b_n^{-1}) \max_{\ell \in \mathcal{L}} \{\text{mineg}(\Gamma_\ell)\}^{-1} = o(b_n^{-1}),$$

which implies the main claim of the theorem, since  $1 \lesssim \sigma_\ell \lesssim 1$  for all  $\ell \in \mathcal{L}$  by Condition 2.

Application of Lemma 5 to the empirical process indexed by  $[\Gamma_\ell^{-1} \psi_\ell\{w, \alpha_\ell, h_\ell(z_\ell)\}, \ell \in \mathcal{L}]$  and that  $\|F\|_{P,q} \leq C$  gives that with probability  $1 - o(1)$  for some constant  $C$ :

$$\max_{\ell \in \mathcal{L}} |\Gamma_\ell^{-1} \sqrt{n} \mathbb{E}_n \psi_\ell\{w, \alpha_\ell, h_\ell(z_\ell)\}| \leq C \sqrt{\log(p_1 \vee n)}.$$

Step 4. (Define and Bound  $II_1$  and  $II_2$ ). Define  $K = 1 + M$ ,  $\mu_\ell(z_\ell) = \{\mu_{\ell k}(z_\ell)\}_{k=1}^K = \{\alpha, \tilde{h}_\ell(z_\ell)\}'$ , where  $z_\ell \mapsto \tilde{h}_\ell(z_\ell)$  is a generic measurable function  $\mathcal{Z}_\ell \rightarrow T_\ell$ ,  $\nu_\ell(z_\ell) = \{\alpha_\ell, h_\ell(z_\ell)\}'$ ,  $\{\nu_{\ell k}(z_\ell)\}_{k=1}^K =$

$\{\alpha_\ell, h_\ell(z_\ell)\}'$ ;  $\widehat{\nu}_\ell(z_\ell) = \{\widehat{\nu}_{\ell k}(z_\ell)\}_{k=1}^K = \{\widehat{\alpha}_\ell, \widehat{h}_\ell(z_\ell)\}'$ , and let  $\bar{\nu}_\ell(z_\ell)$  be a vector on the line connecting  $\nu_\ell(w)$  and  $\mu_\ell(w)$ . We define

$$\begin{aligned} II_1(\ell) &= \sum_{r,k=1}^K \sqrt{n} \mathbb{E}(\partial_{\nu_k} \partial_{\nu_r} \mathbb{E}[\psi_\ell\{w_\ell, \bar{\nu}_\ell(z_\ell)\} | z_\ell] \{\mu_{\ell r}(z_\ell) - \nu_{\ell r}(z_\ell)\} \{\mu_{\ell k}(z_\ell) - \nu_{\ell k}(z_\ell)\}), \\ II_2(\ell) &= \mathbb{G}_n \left[ \psi_\ell\{w_\ell, \alpha, \tilde{h}_\ell(z_\ell)\} - \psi_\ell\{w_\ell, \alpha_\ell, h_\ell(z_\ell)\} \right], \end{aligned}$$

with expressions above are evaluated at  $\mu_{\ell k}(\cdot) = \widehat{\nu}_{\ell k}(\cdot)$ ,  $\mu_{\ell r}(\cdot) = \widehat{\nu}_{\ell r}(\cdot)$ ,  $\alpha = \widehat{\alpha}_\ell$ ,  $\tilde{h}_\ell(\cdot) = \widehat{h}_\ell(\cdot)$  after computing expectations. Using Condition 3 and the claim of Step 1, we conclude that with probability  $1 - o(1)$  uniformly in  $\ell \in \mathcal{L}$ :

$$\begin{aligned} |II_1(\ell)| &\leq \sum_{r,k=1}^K \sqrt{n} \mathbb{E} \{ C |\mu_{\ell r}(z_\ell) - \nu_{\ell r}(z_\ell)| |\mu_{\ell k}(z_\ell) - \nu_{\ell k}(z_\ell)| \} \\ &\leq C \sqrt{n} K^2 \max_{k \leq K} \|\mu_{\ell k} - \nu_{\ell k}\|_{P,2}^2 \lesssim \sqrt{n} \rho_n^2 = o(b_n^{-1}), \end{aligned}$$

where expressions above are evaluated at  $\mu_{\ell k}(\cdot) = \widehat{\nu}_{\ell k}(\cdot)$ ,  $\mu_{\ell r}(\cdot) = \widehat{\nu}_{\ell r}(\cdot)$  after computing expectations.

With probability  $1 - o(1)$  we have  $\max_{\ell \in \mathcal{L}} |II_1(\ell)| \lesssim \sup_{f \in \mathcal{F}''} |\mathbb{G}_n(f)|$  where

$$\mathcal{F}'' = [\psi_\ell\{w, \alpha, h(z_\ell)\} - \psi_\ell\{w, \alpha_\ell, h_\ell(z_\ell)\} : \ell \in \mathcal{L}, h \in \mathcal{H}_\ell, \alpha \in \mathcal{A}_{\ell n}]$$

and  $\mathcal{A}_{\ell n} := \{\alpha \in \mathcal{A}_\ell : |\alpha - \alpha_\ell| \leq C \rho_n\}$ . We note that similarly to Step 2,  $\text{ent}(\varepsilon, \mathcal{F}'') \lesssim s \log(p_1 a_n / \varepsilon) \lesssim s \log(a_n / \varepsilon)$ . We wish to apply Lemma 5. We can choose  $\sigma$  in Lemma 5 so that  $\sup_{f \in \mathcal{F}''} \|f\|_{P,2} \leq \sigma \lesssim \tau_n^{5/2}$ . Indeed, this follows from the following calculation:

$$\begin{aligned} \sup_{f \in \mathcal{F}''} \|f\|_{P,2}^2 &\leq \sup_{\ell \in \mathcal{L}, \mu_\ell \in \mathcal{A}_{\ell n} \times \mathcal{H}_\ell} \mathbb{E} \left\{ \mathbb{E} \left( [\psi_\ell\{w, \mu_\ell(z_\ell)\} - \psi_\ell\{w, \nu_\ell(z_\ell)\}]^2 \mid z_\ell \right) \right\}, \\ &\leq \sup_{\ell \in \mathcal{L}, \mu_\ell \in \mathcal{A}_{\ell n} \times \mathcal{H}_\ell} \mathbb{E} \{ C \|\mu_\ell(z_\ell) - \nu_\ell(z_\ell)\|^s \}, \\ &= \sup_{\ell \in \mathcal{L}, \mu_\ell \in \mathcal{A}_{\ell n} \times \mathcal{H}_\ell} C \|\mu_\ell - \nu_\ell\|_{P,\varsigma}^s \leq \sup_{\ell \in \mathcal{L}, \mu_\ell \in \mathcal{A}_{\ell n} \times \mathcal{H}_\ell} C \|\mu_\ell - \nu_\ell\|_{P,2}^s \lesssim \rho_n^s, \end{aligned}$$

where  $\mu_\ell$  is defined as before. Here the first inequality follows by the law of iterated expectations; the second inequality follows by Condition 3; and the last inequality follows from  $\varsigma \in [1, 2]$  by Condition 3 and the monotonicity of the norm  $\|\cdot\|_{P,q}$  in  $q \in [1, \infty]$ . Application of Lemma 5 and of the inequality  $M_{P,2} \leq n^{1/q} \|F\|_{P,q}$  gives that with probability  $1 - o(1)$ :

$$\max_{\ell \in \mathcal{L}} |II_2(\ell)| \leq \sup_{f \in \mathcal{F}''} |\mathbb{G}_n(f)| \leq C \left( \rho_n^{s/2} \sqrt{s \log a_n} + n^{-1/2} s n^{1/q} \|F_1\|_{P,q} \log a_n + \rho_n \right) = o(b_n^{-1}),$$

for some constant  $C$ , where the last equality follows from the growth conditions of Condition 3.

Step 5. (Auxiliary). We show that with probability  $1 - o(1)$ ,  $\inf_{\alpha \in \mathcal{A}_\ell} \sqrt{n} |\mathbb{E}_n \psi_\ell\{w_\ell, \alpha, \widehat{h}_\ell(z_\ell)\}| = o(b_n^{-1})$ . We have that with probability  $1 - o(1)$ ,  $\inf_{\alpha \in \mathcal{A}_\ell} \sqrt{n} |\mathbb{E}_n \psi_\ell\{w_\ell, \alpha, \widehat{h}_\ell(z_\ell)\}| \leq \sqrt{n} |\mathbb{E}_n \psi_\ell\{w_\ell, \bar{\alpha}_\ell, \widehat{h}_\ell(z_\ell)\}|$ , where  $\bar{\alpha}_\ell = \alpha_\ell - \Gamma_\ell^{-1} \mathbb{E}_n \psi_\ell\{w_\ell, \alpha_\ell, h_\ell(z)\}$ , since  $\bar{\alpha}_\ell \in \mathcal{A}_\ell$  for all  $\ell \in \mathcal{L}$  with probability  $1 - o(1)$ , and in fact  $\max_{\ell \in \mathcal{L}} |\bar{\alpha}_\ell - \alpha_\ell| \lesssim \sqrt{\log(p_1 \vee n)/n} = o(1)$  by the last sentence of Step 3. Then, arguing similarly to Steps 3 and 4, we can conclude that with probability  $1 - o(1)$ : uniformly in  $\ell \in \mathcal{L}$ ,

$$\sqrt{n} |\mathbb{E}_n \psi_\ell\{w_\ell, \bar{\alpha}_\ell, \widehat{h}_\ell(z_\ell)\}| \leq \sqrt{n} |\mathbb{E}_n \psi_\ell\{w_\ell, \alpha_\ell, h_\ell(z_\ell)\}| + \Gamma_\ell (\bar{\alpha}_\ell - \alpha_\ell) + \Delta_\ell (\widehat{h}_\ell - h_\ell) + o(b_n^{-1}),$$

where the first term on the right side vanishes by definition of  $\bar{\alpha}_\ell$  and by  $\Delta_\ell(\hat{h}_\ell - h_\ell) = 0$ .  $\square$

## REFERENCES

- [1] R. Baraniuk, M. Davenport, R. DeVore, and M. Wakin. A simple proof of the restricted isometry property for random matrices. *Constructive Approximation*, 28:253–263, 2008.
- [2] A. Belloni, D. Chen, V. Chernozhukov, and C. Hansen. Sparse models and methods for optimal instruments with an application to eminent domain. *Econometrica*, 80(6):2369–2430, November 2012.
- [3] A. Belloni and V. Chernozhukov.  $\ell_1$ -penalized quantile regression for high dimensional sparse models. *Annals of Statistics*, 39(1):82–130, 2011.
- [4] A. Belloni and V. Chernozhukov. Least squares after model selection in high-dimensional sparse models. *Bernoulli*, 19(2):521–547, 2013.
- [5] A. Belloni, V. Chernozhukov, and I. Fernandez-Val. Conditional Quantile Processes based on Series or Many Regressors. *ArXiv e-prints*, May 2011.
- [6] A. Belloni, V. Chernozhukov, and C. Hansen. Inference on treatment effects after selection amongst high-dimensional controls. *ArXiv*, 2011.
- [7] A. Belloni, V. Chernozhukov, and C. Hansen. Inference for high-dimensional sparse econometric models. *Advances in Economics and Econometrics. 10th World Congress of Econometric Society, held in August 2010*, III:245–295, 2013.
- [8] A. Belloni, V. Chernozhukov, and K. Kato. Robust inference in high-dimensional sparse quantile regression models. *Working Paper*, 2013.
- [9] P. J. Bickel, Y. Ritov, and A. B. Tsybakov. Simultaneous analysis of Lasso and Dantzig selector. *Annals of Statistics*, 37(4):1705–1732, 2009.
- [10] V. Chernozhukov, D. Chetverikov, and K. Kato. Gaussian approximation of suprema of empirical processes. *arXiv preprint arXiv:1212.6885*, 2012.
- [11] V. Chernozhukov, D. Chetverikov, and K. Kato. Gaussian approximations and multiplier bootstrap for maxima of sums of high-dimensional random vectors. *Annals of Statistics*, 41(6):2786–2819, 2013.
- [12] Victor Chernozhukov and Christian Hansen. Instrumental variable quantile regression: A robust inference approach. *Journal of Econometrics*, 142:379–398, 2008.
- [13] Xuming He and Qi-Man Shao. On parameters of increasing dimensions. *J. Multivariate Anal.*, 73(1):120–135, 2000.
- [14] Bing-Yi Jing, Qi-Man Shao, and Qiyang Wang. Self-normalized Cramer-type large deviations for independent random variables. *Ann. Probab.*, 31(4):2167–2215, 2003.
- [15] K. Kato. Group lasso for high dimensional sparse quantile regression models. Preprint, ArXiv, 2011.
- [16] Roger Koenker. *Quantile regression*, volume 38 of *Econometric Society Monographs*. Cambridge University Press, Cambridge, 2005.
- [17] Michael R. Kosorok. *Introduction to Empirical Processes and Semiparametric Inference*. Series in Statistics. Springer, Berlin, 2008.
- [18] Sokbae Lee. Efficient semiparametric estimation of a partially linear quantile regression model. *Econometric theory*, 19:1–31, 2003.
- [19] Hannes Leeb and Benedikt M. Pötscher. Model selection and inference: facts and fiction. *Economic Theory*, 21:21–59, 2005.
- [20] Hannes Leeb and Benedikt M. Pötscher. Sparse estimators and the oracle property, or the return of Hodges’ estimator. *J. Econometrics*, 142(1):201–211, 2008.
- [21] J. Neyman. Optimal asymptotic tests of composite statistical hypotheses. In U. Grenander, editor, *Probability and Statistics, the Harold Cramer Volume*. New York: John Wiley and Sons, Inc., 1959.
- [22] J. Neyman.  $c(\alpha)$  tests and their use. *Sankhya*, 41:1–21, 1979.
- [23] J. L. Powell. Censored regression quantiles. *Journal of Econometrics*, 32:143–155, 1986.
- [24] Joseph P. Romano and Michael Wolf. Stepwise multiple testing as formalized data snooping. *Econometrica*, 73(4):1237–1282, July 2005.
- [25] M. Rudelson and S. Zhou. Reconstruction from anisotropic random measurements. *ArXiv:1106.1151*, 2011.

- [26] Mark Rudelson and Roman Vershynin. On sparse reconstruction from fourier and gaussian measurements. *Communications on Pure and Applied Mathematics*, 61:10251045, 2008.
- [27] R. Tibshirani. Regression shrinkage and selection via the lasso. *J. Roy. Statist. Soc. Ser. B*, 58:267–288, 1996.
- [28] S. A. van de Geer. High-dimensional generalized linear models and the lasso. *Annals of Statistics*, 36(2):614–645, 2008.
- [29] A. W. van der Vaart. *Asymptotic Statistics*. Cambridge Series in Statistical and Probabilistic Mathematics, 1998.
- [30] A. W. van der Vaart and J. A. Wellner. *Weak Convergence and Empirical Processes*. Springer Series in Statistics, 1996.
- [31] Lie Wang.  $L_1$  penalized lad estimator for high dimensional linear regression. *ArXiv*, 2012.
- [32] Cun-Hui Zhang and Stephanie S. Zhang. Confidence intervals for low-dimensional parameters with high-dimensional data. *ArXiv.org*, (arXiv:1110.2563v1), 2011.
- [33] S. Zhou. Restricted eigenvalue conditions on subgaussian matrices. *ArXiv:0904.4723v2*, 2009.

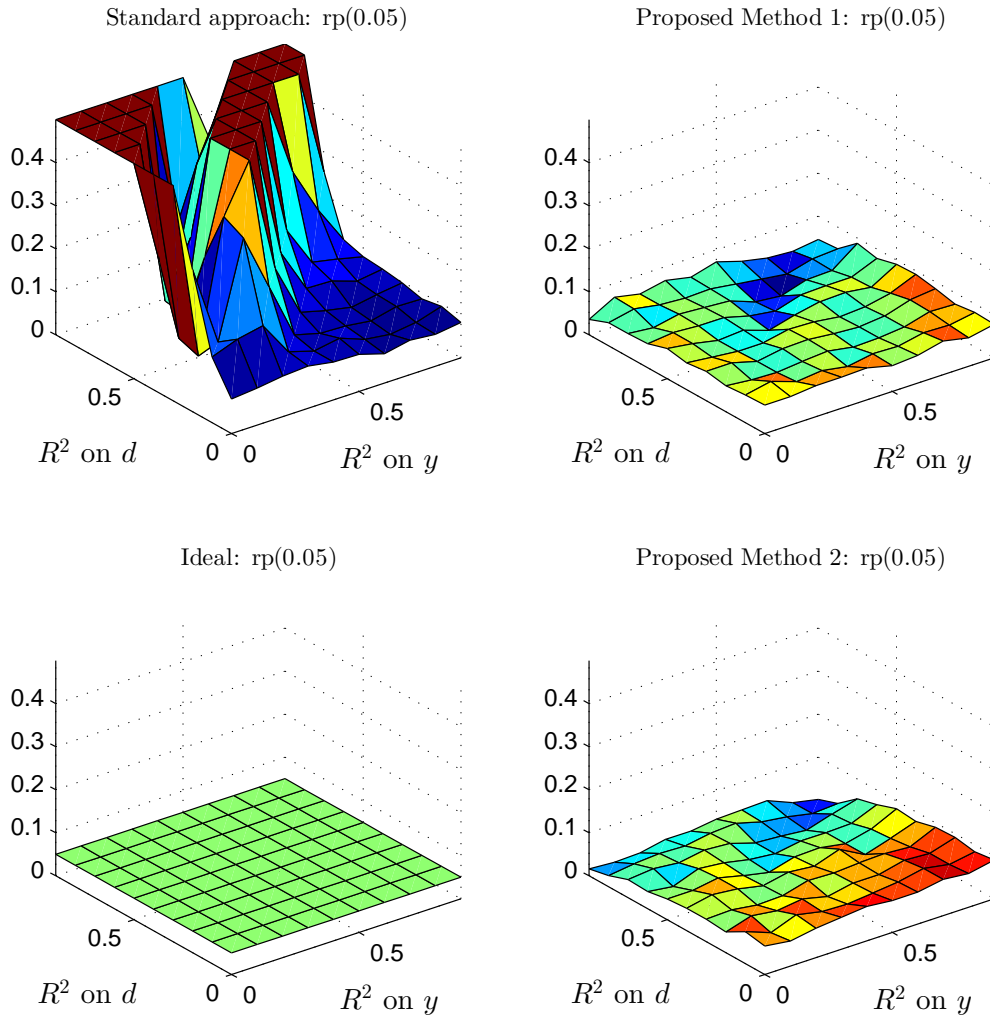


FIGURE 1. The figure displays the empirical rejection probabilities of the nominal 5% level tests of a true hypothesis based on different testing procedures: the top left plot is based on the standard post-model selection procedure based on  $\tilde{\alpha}$ , the top right plot is based on the proposed post-model selection procedure based on  $\tilde{\alpha}$ , and the bottom left plot is based on another proposed procedure based on the statistic  $L_n$ . The results are based on 500 replications for each of the 100 combinations of  $R^2$ 's in the primary and auxiliary equations in (4.23). Ideally we should observe the 5% rejection rate (of a true null) uniformly across the parameter space (as in bottom right plot).

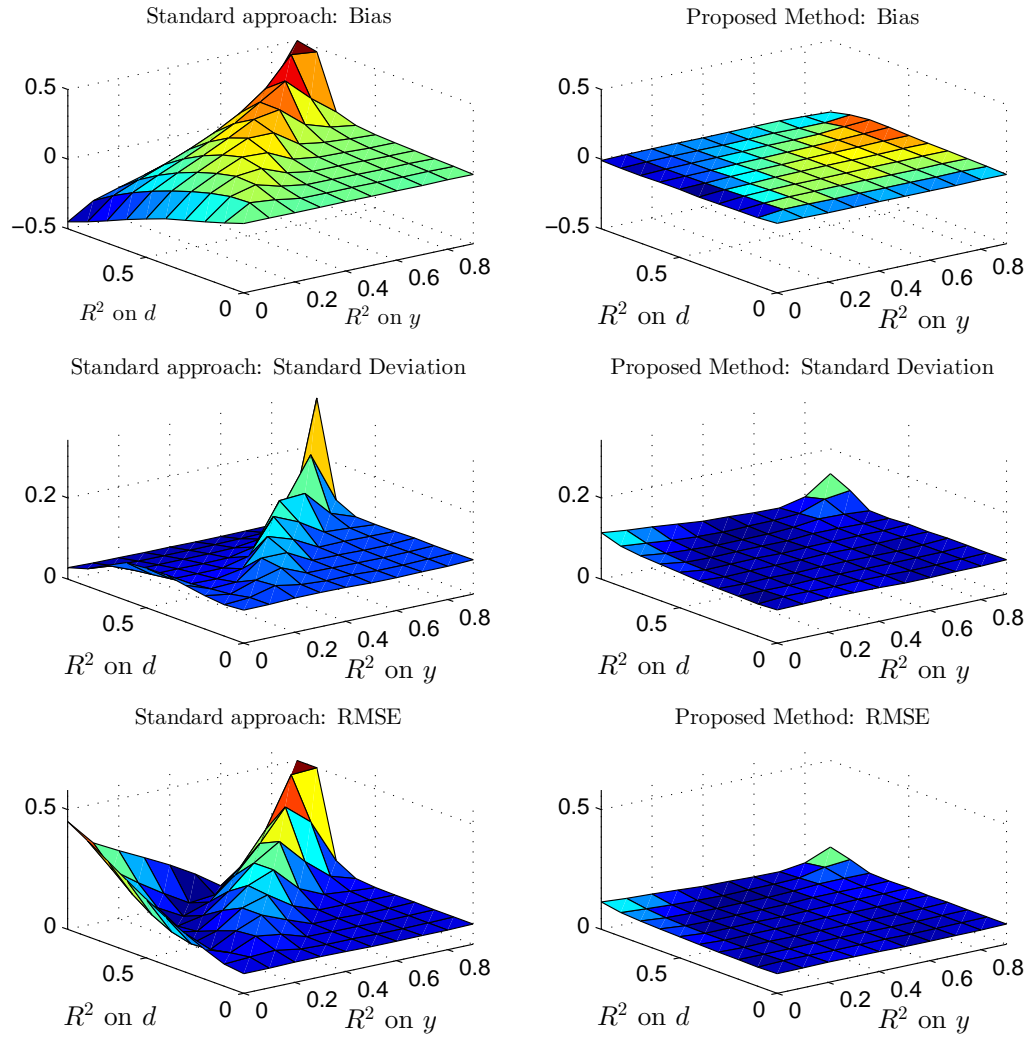


FIGURE 2. The figure displays mean bias (top row), standard deviation (middle row), and root mean square error (bottom row) for the the proposed post-model selection estimator  $\tilde{\alpha}$  (right column) and the standard post-model selection estimator  $\tilde{\alpha}$  (left column). The results are based on 500 replications for each of the 100 combinations of  $R^2$ 's in the primary and auxiliary equations in (4.23).

# Supplementary Appendix for “Uniform Post Selection Inference for LAD Regression Models”

## APPENDIX F. AUXILIARY RESULTS FOR $\ell_1$ -LAD AND HETEROSCEDASTIC LASSO

In this section we state relevant theoretical results on the performance of the estimators  $\ell_1$ -LAD, Post- $\ell_1$ -LAD, heteroscedastic Lasso, and heteroscedastic Post-Lasso. These results were developed in [3] and [2]. The main design condition relies on the restricted eigenvalue proposed in [9], namely for  $\tilde{x}_i = (d_i, x_i)'$

$$\kappa_{\mathbf{c}} = \inf_{\|\delta_{T^c}\|_1 \leq \mathbf{c}\|\delta_T\|_1} \|\tilde{x}'\delta\|_{2,n}/\|\delta_T\|, \quad (\text{F.41})$$

where  $\mathbf{c} = (c+1)/(c-1)$  for the slack constant  $c > 1$ , see [9]. It is well known that well behaved sparse eigenvalues imply that  $\kappa_{\mathbf{c}}$  is bounded away from zero if  $\mathbf{c}$  is bounded for any subset  $T \subset \{1, \dots, p\}$  with  $|T| \leq s$ .

**F.1.  $\ell_1$ -Penalized LAD.** For a data generating process such that  $\text{pr}(y_i \leq \tilde{x}_i'\eta_0 \mid \tilde{x}_i) = 1/2$ , independent across  $i$  ( $i = 1, \dots, n$ ) we consider the estimation of  $\eta_0$  via the  $\ell_1$ -penalized LAD regression estimate

$$\hat{\eta} \in \arg \min_{\eta} \mathbb{E}_n[|y - \tilde{x}'\eta|] + \frac{\lambda}{n} \|\Psi\eta\|_1$$

where  $\Psi^2 = \text{diag}(\mathbb{E}_n[\tilde{x}_1^2], \dots, \mathbb{E}_n[\tilde{x}_p^2])$  is a diagonal matrix of penalty loadings. As established in [3] and [31], under the event that

$$\frac{\lambda}{n} \geq 2c \|\Psi^{-1} \mathbb{E}_n[(1/2 - 1\{y \leq \tilde{x}'\eta_0\})\tilde{x}]\|_{\infty}, \quad (\text{F.42})$$

the estimator above achieves good theoretical guarantees under mild design conditions. Although  $\eta_0$  is unknown, we can set  $\lambda$  so that the event in (F.42) holds with high probability. In particular, the pivotal rule discussed in [3] proposes to set  $\lambda = c'n\Lambda(1 - \gamma \mid \tilde{x})$  for  $c' > c$  and  $\gamma \rightarrow 0$  where

$$\Lambda(1 - \gamma \mid \tilde{x}) := (1 - \gamma)\text{-quantile of } 2\|\Psi^{-1} \mathbb{E}_n[(1/2 - 1\{U \leq 1/2\})\tilde{x}]\|_{\infty}, \quad (\text{F.43})$$

and where  $U_i$  are independent uniform random variables on  $(0, 1)$ , independent of  $\tilde{x}_1, \dots, \tilde{x}_n$ . We suggest  $\gamma = 0.1/\log n$  and  $c' = 1.1c$ . This quantity can be easily approximated via simulations. Below we summarize required regularity conditions.

**Condition PLAD.** Assume that  $\|\eta_0\|_0 = s \geq 1$ ,  $c \leq \mathbb{E}_n[\tilde{x}_j^2] \leq C$  for all  $1 \leq j \leq p$ , the conditional density of  $y_i$  given  $d_i$ , denoted by  $f_i(\cdot)$ , and its derivative are bounded by  $\bar{f}$  and  $\bar{f}'$ , respectively, and  $f_i(\tilde{x}_i'\eta_0) \geq \underline{f} > 0$  is bounded away from zero uniformly in  $n$ .

Condition PLAD is implied by Condition 1. The assumption on the conditional density is standard in the quantile regression literature even with fixed  $p$  or  $p$  increasing slower than  $n$  (see respectively [16] and [5]). Next we present bounds on the prediction norm of the  $\ell_1$ -LAD estimator.

**Lemma 6** (Estimation Error of  $\ell_1$ -LAD). *Under Condition PLAD, and using  $\lambda = c'n\Lambda(1 - \gamma \mid \tilde{x})$ , we have with probability  $1 - 2\gamma - o(1)$  for  $n$  large enough*

$$\|\tilde{x}'(\hat{\eta} - \eta_0)\|_{2,n} \lesssim \frac{\lambda\sqrt{s}}{n\kappa_{\mathbf{c}}} + \frac{1}{\kappa_{\mathbf{c}}} \sqrt{\frac{s \log(p/\gamma)}{n}},$$



provided  $\left\{ \frac{n\kappa_{\mathbf{c}}}{\lambda\sqrt{s}} + \frac{n\kappa_{\mathbf{c}}}{\sqrt{sn \log([p \vee n]/\gamma)}} \right\} \frac{\bar{f}f'}{\underline{f}} \inf_{\delta \in \Delta_{\mathbf{c}}} \frac{\|\tilde{x}'\delta\|_{2,n}^3}{\mathbb{E}_n[|\tilde{x}'\delta|^3]} \rightarrow \infty$ .

Lemma 6 establishes the rate of convergence in the prediction norm for the  $\ell_1$ -LAD estimator in a parametric setting. The extra growth condition required for identification is mild. For instance we typically have  $\lambda \lesssim \sqrt{\log(n \vee p)/n}$  and for many designs of interest we have  $\inf_{\delta \in \Delta_{\mathbf{c}}} \|\tilde{x}'\delta\|_{2,n}^3 / \mathbb{E}_n[|\tilde{x}'\delta|^3]$  bounded away from zero (see [3]). For more general designs we have

$$\inf_{\delta \in \Delta_{\mathbf{c}}} \frac{\|\tilde{x}'\delta\|_{2,n}^3}{\mathbb{E}_n[|\tilde{x}'\delta|^3]} \geq \inf_{\delta \in \Delta_{\mathbf{c}}} \frac{\|\tilde{x}'\delta\|_{2,n}}{\|\delta\|_1 \max_{i \leq n} \|\tilde{x}_i\|_{\infty}} \geq \frac{\kappa_{\mathbf{c}}}{\sqrt{s}(1 + \mathbf{c}) \max_{i \leq n} \|\tilde{x}_i\|_{\infty}}$$

which implies the extra growth condition under  $K_x^2 s^2 \log(p \vee n) \leq \delta_n \kappa_{\mathbf{c}}^2 n$ .

In order to alleviate the bias introduced by the  $\ell_1$ -penalty, we can consider the associated post-model selection estimate associated with a selected support  $\hat{T}$

$$\tilde{\eta} \in \arg \min_{\eta} \left\{ \mathbb{E}_n[|y - \tilde{x}'\eta|] : \eta_j = 0 \text{ if } j \notin \hat{T} \right\}. \quad (\text{F.44})$$

The following result characterizes the performance of the estimator in (F.44), see [3] for the proof.

**Lemma 7** (Estimation Error of Post- $\ell_1$ -LAD). *Assume the conditions of Lemma 6 hold,  $\text{supp}(\hat{\eta}) \subseteq \hat{T}$ , and let  $\hat{s} = |\hat{T}|$ . Then we have for  $n$  large enough*

$$\|\tilde{x}'(\tilde{\eta} - \eta_0)\|_{2,n} \lesssim_P \sqrt{\frac{(\hat{s} + s) \log(n \vee p)}{n \phi_{\min}(\hat{s} + s)}} + \frac{\lambda\sqrt{s}}{n\kappa_{\mathbf{c}}} + \frac{1}{\kappa_{\mathbf{c}}} \sqrt{\frac{s \log(p/\gamma)}{n}},$$

provided  $\left\{ \frac{n\sqrt{\phi_{\min}(\hat{s}+s)}}{\lambda\sqrt{s}} + \frac{n\sqrt{\phi_{\min}(\hat{s}+s)}}{\sqrt{sn \log([p \vee n]/\gamma)}} \right\} \frac{\bar{f}f'}{\underline{f}} \inf_{\|\delta\|_0 \leq \hat{s}+s} \frac{\|\tilde{x}'\delta\|_{2,n}^3}{\mathbb{E}_n[|\tilde{x}'\delta|^3]} \rightarrow_P \infty$ .

Lemma 7 provides the rate of convergence in the prediction norm for the post model selection estimator despite of possible imperfect model selection. The rates rely on the overall quality of the selected model (which is at least as good as the model selected by  $\ell_1$ -LAD) and the overall number of components  $\hat{s}$ . Once again the extra growth condition required for identification is mild. For more general designs we have

$$\inf_{\|\delta\|_0 \leq \hat{s}+s} \frac{\|\tilde{x}'\delta\|_{2,n}^3}{\mathbb{E}_n[|\tilde{x}'\delta|^3]} \geq \inf_{\|\delta\|_0 \leq \hat{s}+s} \frac{\|\tilde{x}'\delta\|_{2,n}}{\|\delta\|_1 \max_{i \leq n} \|\tilde{x}_i\|_{\infty}} \geq \frac{\sqrt{\phi_{\min}(\hat{s} + s)}}{\sqrt{\hat{s} + s} \max_{i \leq n} \|\tilde{x}_i\|_{\infty}}.$$

**Comment F.1.** In Step 1 of Algorithm 2 we use  $\ell_1$ -LAD with  $\tilde{x}_i = (d_i, x_i)'$ ,  $\hat{\delta} := \hat{\eta} - \eta_0 = (\hat{\alpha} - \alpha_0, \hat{\beta}' - \beta_0)'$ , and we are interested on rates for  $\|x'(\hat{\beta} - \beta_0)\|_{2,n}$  instead of  $\|\tilde{x}'\hat{\delta}\|_{2,n}$ . However, it follows that

$$\|x'(\hat{\beta} - \beta_0)\|_{2,n} \leq \|\tilde{x}'\hat{\delta}\|_{2,n} + |\hat{\alpha} - \alpha_0| \cdot \|d\|_{2,n}.$$

Since  $s \geq 1$ , without loss of generality we can assume the component associated with the treatment  $d_i$  belongs to  $T$  (at the cost of increasing the cardinality of  $T$  by one which will not affect the rate of convergence). Therefore we have that

$$|\hat{\alpha} - \alpha_0| \leq \|\hat{\delta}_T\| \leq \|\tilde{x}'\hat{\delta}\|_{2,n} / \kappa_{\mathbf{c}}.$$

In most applications of interest  $\|d\|_{2,n}$  and  $1/\kappa_{\mathbf{c}}$  are bounded from above with high probability. Similarly, in Step 1 of Algorithm 1 we have that the Post- $\ell_1$ -LAD estimator satisfies

$$\|x'(\tilde{\beta} - \beta_0)\|_{2,n} \leq \|\tilde{x}'\tilde{\delta}\|_{2,n} \left( 1 + \|d\|_{2,n} / \sqrt{\phi_{\min}(\hat{s} + s)} \right).$$

**F.2. Heteroscedastic Lasso.** In this section we consider the equation (1.4) in the form

$$d_i = x_i' \theta_0 + v_i, \quad \mathbb{E}[v_i | x_i] = 0, \quad (\text{F.45})$$

where we observe  $\{(d_i, x_i')\}_{i=1}^n$ ,  $(x_i)_{i=1}^n$  are non-stochastic, and  $(v_i)_{i=1}^n$  are independent across  $i$  but not necessarily identically distributed. The unknown support of  $\theta_0$  is denoted by  $T_d$  and it satisfies  $|T_d| \leq s$ . To estimate  $\theta_0$  and consequently  $v_i$ , we compute

$$\hat{\theta} \in \arg \min_{\theta} \mathbb{E}_n[(d - x'\theta)^2] + \frac{\lambda}{n} \|\hat{\Gamma}\theta\|_1 \quad \text{and set } \hat{v}_i = d_i - x_i' \hat{\theta}, \quad i = 1, \dots, n, \quad (\text{F.46})$$

where  $\lambda$  and  $\hat{\Gamma}$  are the associated penalty level and loadings which are potentially data-driven. In this case the following regularization event plays an important role

$$\frac{\lambda}{n} \geq 2c \|\hat{\Gamma}^{-1} \mathbb{E}_n[x(d - x'\theta_0)]\|_{\infty}. \quad (\text{F.47})$$

As discussed in [9], [4] and [2], the event above implies that the estimator  $\hat{\theta}$  satisfies  $\|\hat{\theta}_{T_d^c}\|_1 \leq \mathbf{c} \|\hat{\theta}_{T_d} - \theta_0\|_1$  where  $\mathbf{c} = (c+1)/(c-1)$ . Thus rates of convergence for  $\hat{\theta}$  and  $\hat{v}_i$  defined on (F.46) can be established based on the restricted eigenvalue  $\kappa_c$  defined in (F.41) with  $\tilde{x}_i = x_i$  and  $T = T_d$ .

The following are sufficient high-level conditions where again the sequences  $\Delta_n$  and  $\delta_n$  go to zero and  $C$  is a positive constant independent of  $n$ , and let  $\bar{\mathbb{E}}[f(d, x)] := \frac{1}{n} \sum_{i=1}^n \mathbb{E}[f(d_i, x_i)]$ .

**Condition HL.** For the model (F.45), for  $s = s_n \geq 1$  we have  $\|\theta_0\|_0 \leq s$  and

- (i)  $c \leq \mathbb{E}_n[x_j^2] \leq C$  for all  $j = 1, \dots, p$ ;  $c \leq \mathbb{E}[v_i^2 | x_i] \leq C$  for all  $i = 1, \dots, n$ ,
- (ii)  $\max_{1 \leq j \leq p} \{(\bar{\mathbb{E}}[|x_j v|^3])^{1/3} / (\bar{\mathbb{E}}[|x_j v|^2])^{1/2}\} \Phi^{-1}(1 - \gamma/2p) \leq \delta_n n^{1/6}$ ,
- (iii)  $\max_{1 \leq j \leq p} |(\mathbb{E}_n - \bar{\mathbb{E}})[x_j^2 v^2]| + \max_{j \leq p} |(\mathbb{E}_n - \bar{\mathbb{E}})[x_j^2 d^2]| \leq \delta_n$ , with probability  $1 - \Delta_n$ .

Condition HL is implied by Condition 1. Several primitive moment conditions imply the various cross moments bounds. These conditions also allow us to invoke moderate deviation theorems for self-normalized sums from [14] to bound some important error components. Despite heteroscedastic non-Gaussian noise, Those results allows a sharp choice of penalty level and loadings was analyzed in [2] which is summarized by the following lemma.

Valid options for setting the penalty level and the loadings for  $j = 1, \dots, p$ , are

$$\begin{aligned} \text{initial} \quad \hat{\gamma}_j &= \sqrt{\mathbb{E}_n[x_j^2 (d - \bar{d})^2]}, & \lambda &= 2c\sqrt{n}\Phi^{-1}(1 - \gamma/(2p)), \\ \text{refined} \quad \hat{\gamma}_j &= \sqrt{\mathbb{E}_n[x_j^2 \hat{v}^2]}, & \lambda &= 2c\sqrt{n}\Phi^{-1}(1 - \gamma/(2p)), \end{aligned} \quad (\text{F.48})$$

where  $c > 1$  is a constant,  $\gamma \in (0, 1)$ ,  $\bar{d} := \mathbb{E}_n[d]$  and  $\hat{v}_i$  is an estimate of  $v_i$  based on Lasso with the initial option (or iterations). [2] established that using either of the choices in (F.48) implies that the regularization event (F.47) holds with high probability. Next we present results on the performance of the estimators generated by Lasso.

**Lemma 8.** *Under Condition HL and setting  $\lambda = 2c'\sqrt{n}\Phi^{-1}(1 - \gamma/2p)$  for  $c' > c > 1$ , and using penalty loadings as in (F.48), there is an uniformly bounded  $\mathbf{c}$  such that we have*

$$\|\hat{v} - v\|_{2,n} = \|x'(\hat{\theta} - \theta_0)\|_{2,n} \lesssim_P \frac{\lambda\sqrt{s}}{n\kappa_c} \quad \text{and} \quad \|\hat{v}_i - v_i\|_{\infty} \leq \|\hat{\theta} - \theta_0\|_1 \max_{i \leq n} \|x_i\|_{\infty}.$$

Associated with Lasso we can define the Post-Lasso estimator as

$$\tilde{\theta} \in \arg \min_{\theta} \left\{ \mathbb{E}_n[(d - x'\theta)^2] : \theta_j = 0 \text{ if } \hat{\theta}_j = 0 \right\} \text{ and set } \tilde{v}_i = d_i - x_i'\tilde{\theta}. \quad (\text{F.49})$$

That is, the Post-Lasso estimator is simply the least squares estimator applied to the covariates selected by Lasso in (F.46). Sparsity properties of the Lasso estimator  $\hat{\theta}$  under estimated weights follows similarly to the standard Lasso analysis derived in [2]. By combining such sparsity properties and the rates in the prediction norm we can establish rates for the post-model selection estimator under estimated weights. The following result summarizes the properties of the Post-Lasso estimator and relies on sparse eigenvalues of the empirical Gram matrix

$$\kappa' \leq \min_{\|\delta\|_0 \leq \ell_{n,s}} \frac{\|x'\delta\|_{2,n}^2}{\|\delta\|^2} \leq \max_{\|\delta\|_0 \leq \ell_{n,s}} \frac{\|x'\delta\|_{2,n}^2}{\|\delta\|^2} \leq \kappa'' \quad (\text{F.50})$$

**Lemma 9** (Model Selection Properties of Lasso and Properties of Post-Lasso). *Suppose that Condition HL and (F.50) hold. Consider the Lasso estimator with penalty level and loadings specified as in Lemma 8. Then the data-dependent model  $\hat{T}_d$  selected by the Lasso estimator  $\hat{\theta}$  satisfies with probability  $1 - \Delta_n$ :*

$$\|\tilde{\theta}\|_0 = |\hat{T}_d| \lesssim s. \quad (\text{F.51})$$

Moreover, the Post-Lasso estimator obeys

$$\|\tilde{v} - v\|_{2,n} = \|x'(\tilde{\theta} - \theta_0)\|_{2,n} \lesssim_P \sqrt{\frac{s \log(p \vee n)}{n}}.$$

## APPENDIX G. ALTERNATIVE IMPLEMENTATION VIA DOUBLE SELECTION

An alternative proposal for the method is reminiscent of the double selection method proposed in [6] for partial linear models. This version replaces Step 3 with a LAD regression of  $y$  on  $d$  and all covariates selected in Steps 1 and 2 (i.e. the union of the selected sets). The method is described as follows:

**Algorithm 3.** (A Double Selection Method)

**Step 1** Run Post- $\ell_1$ -LAD of  $y_i$  on  $d_i$  and  $x_i$ :

$$(\hat{\alpha}, \hat{\beta}) \in \arg \min_{\alpha, \beta} \mathbb{E}_n[|y - d\alpha - x'\beta|] + \frac{\lambda_1}{n} \|\Psi(\alpha, \beta)'\|_1.$$

**Step 2** Run Heteroscedastic Lasso of  $d_i$  on  $x_i$ :

$$\hat{\theta} \in \arg \min_{\theta} \mathbb{E}_n[(d - x'\theta)^2] + \frac{\lambda_2}{n} \|\hat{\Gamma}\theta\|_1.$$

**Step 3** Run LAD regression of  $y_i$  on  $d_i$  and the covariates selected in Step 1 and 2:

$$(\check{\alpha}, \check{\beta}) \in \arg \min_{\alpha, \beta} \{ \mathbb{E}_n[|y - d\alpha - x'\beta|] : \text{supp}(\beta) \subseteq \text{supp}(\hat{\beta}) \cup \text{supp}(\hat{\theta}) \}.$$

The double selection algorithm has three steps: (1) select covariates based on the standard  $\ell_1$ -LAD regression, (2) select covariates based on heteroscedastic Lasso of the treatment equation, and (3) run a LAD regression with the treatment and all selected covariates.

This approach can also be analyzed through Lemma 1 since it creates instruments implicitly. To see that let  $\widehat{T}^*$  denote the variables selected in Step 1 and 2:  $\widehat{T}^* = \text{supp}(\widehat{\beta}) \cup \text{supp}(\widehat{\theta})$ . By the first order conditions for  $(\check{\alpha}, \check{\beta})$  we have

$$\|\mathbb{E}_n[\varphi(y - d\check{\alpha} - x'\check{\beta})(d, x'_{\widehat{T}^*})']\| = O\left\{\left(\max_{1 \leq i \leq n} |d_i| + K_x |\widehat{T}^*|^{1/2}\right)(1 + |\widehat{T}^*|/n)\right\},$$

which creates an orthogonal relation to any linear combination of  $(d_i, x'_{i\widehat{T}^*})'$ . In particular, by taking the linear combination  $(d_i, x'_{i\widehat{T}^*})(1, -\widetilde{\theta}'_{\widehat{T}^*})' = d_i - x'_{i\widehat{T}^*} \widetilde{\theta}_{\widehat{T}^*} = d_i - x'_i \widetilde{\theta} = \widehat{z}_i$ , which is the instrument in Step 2 of Algorithm 1, we have

$$\mathbb{E}_n[\varphi(y - d\check{\alpha} - x'\check{\beta})\widehat{z}] = O\left\{\|(1, -\widetilde{\theta}')'\| \left(\max_{1 \leq i \leq n} |d_i| + K_x |\widehat{T}^*|^{1/2}\right)(1 + |\widehat{T}^*|/n)\right\}.$$

As soon as the right side is  $o_P(n^{-1/2})$ , the double selection estimator  $\check{\alpha}$  approximately minimizes

$$\widetilde{L}_n(\alpha) = \frac{|\mathbb{E}_n[\varphi(y - d\alpha - x'\check{\beta})\widehat{z}]|^2}{\mathbb{E}_n[\{\varphi(y - d\check{\alpha} - x'\check{\beta})\}^2 \widehat{z}^2]},$$

where  $\widehat{z}_i$  is the instrument created by Step 2 of Algorithm 1. Thus the double selection estimator can be seen as an iterated version of the method based on instruments where the Step 1 estimate  $\widetilde{\beta}$  is updated with  $\check{\beta}$ .

#### APPENDIX H. PROOF OF THEOREM 1: ALGORITHM 2

*Proof of Theorem 1, Algorithm 2.* We will verify Condition ILAD and the desired result then follows from Lemma 1 and noting that  $|\mathbb{E}_n[\mathbb{E}[v^2 | x]] - \mathbb{E}[v^2]| \lesssim_P \delta_n$  and  $\mathbb{E}[v^2]$  is bounded away from zero under Condition 1.

The assumptions on the error density  $f_\varepsilon(\cdot)$  in Condition ILAD(i) are assumed in Condition 1(iv). The moment conditions on  $d_i$  and  $v_i$  in Condition ILAD(i) are assumed in Condition 1(ii).

Because Condition 1(v) and (vi), by Lemma 4 we have for some  $\tilde{\ell}_n \rightarrow \infty$  we have  $\kappa' \leq \phi_{\min}(\tilde{\ell}_n) \leq \phi_{\max}(\tilde{\ell}_n) \leq \kappa''$  with probability  $1 - \Delta_n$ . In turn,  $\kappa_c$  is bounded away from zero with probability  $1 - \Delta_n$  for  $n$  sufficiently large, see [9].

Step 1 relies on  $\ell_1$ -LAD. Condition PLAD is implied by Condition 1. By Lemma 6 and Comment F.1 we have

$$\|x'(\widehat{\beta} - \beta_0)\|_{2,n} \lesssim_P \sqrt{s \log(n \vee p)/n} \quad \text{and} \quad |\widehat{\alpha} - \alpha_0| \lesssim_P \sqrt{s \log(p \vee n)/n} \lesssim o(1) \log^{-1} n$$

because  $s^3 \log^3(n \vee p) \leq \delta_n n$  and the required side condition holds. Indeed, without loss of generality assume that  $T$  contains the treatment so that for  $\tilde{x}_i = (d_i, x'_i)'$ ,  $\delta = (\delta_d, \delta'_x)'$ , because  $\kappa_c$  is bounded away from zero, and the fact that  $\mathbb{E}_n[|d|^3] \lesssim_P \mathbb{E}[|d|^3] = O(1)$ , we have

$$\begin{aligned} \inf_{\delta \in \Delta_c} \frac{\|\tilde{x}'\delta\|_{2,n}^3}{\mathbb{E}_n[|\tilde{x}'\delta|^3]} &\geq \inf_{\delta \in \Delta_c} \frac{\|\tilde{x}'\delta\|_{2,n}^2 \|\delta_T\| \kappa_c}{4\mathbb{E}_n[|x'\delta_x|^3] + 4\mathbb{E}_n[|d\delta_d|^3]} \geq \inf_{\delta \in \Delta_c} \frac{\|\tilde{x}'\delta\|_{2,n}^2 \|\delta_T\| \kappa_c}{4K_x \|\delta_x\|_1 \mathbb{E}_n[|x'\delta_x|^2] + 4|\delta_d|^3 \mathbb{E}_n[|d|^3]} \\ &\geq \inf_{\delta \in \Delta_c} \frac{\|\tilde{x}'\delta\|_{2,n}^2 \|\delta_T\| \kappa_c}{4K_x \|\delta_x\|_1 \{\|\tilde{x}'\delta\|_{2,n} + \|\delta_d d\|_{2,n}\}^2 + 4|\delta_d|^2 \mathbb{E}_n[|d|^3] \|\delta_T\|_1} \\ &\geq \inf_{\delta \in \Delta_c} \frac{\|\tilde{x}'\delta\|_{2,n}^2 \|\delta_T\|_1 \kappa_c / \sqrt{s}}{8K_x(1+c) \|\delta_T\|_1 \|\tilde{x}'\delta\|_{2,n}^2 + 8K_x(1+c) \|\delta_T\|_1 |\delta_d|^2 \{\|d\|_{2,n}^2 + \mathbb{E}_n[|d|^3]\}} \\ &\geq \frac{\kappa_c / \sqrt{s}}{8K_x(1+c) \{1 + \|d\|_{2,n}^2 / \kappa_c^2 + \mathbb{E}_n[|d|^3] / \kappa_c^2\}} \gtrsim_P \frac{1}{\sqrt{s} K_x}. \end{aligned} \tag{H.52}$$

Therefore, since  $\lambda \lesssim \sqrt{n \log(p \vee n)}$  we have

$$\frac{n\kappa_{\mathbf{c}}}{\lambda\sqrt{s} + \sqrt{sn \log(p \vee n)}} \inf_{\delta \in \Delta_{\mathbf{c}}} \frac{\|\tilde{x}'\delta\|_{2,n}^3}{\mathbb{E}_n[|\tilde{x}'\delta|^3]} \gtrsim_P \frac{\sqrt{n}}{K_x s \log(p \vee n)} \rightarrow_P \infty \quad (\text{H.53})$$

under  $K_x^2 s^2 \log^2(p \vee n) \leq \delta_n n$ . Note that the rate for  $\hat{\alpha}$  and the definition of  $\mathcal{A}$  implies  $\{\alpha : |\alpha - \alpha_0| \leq n^{-1/2} \log n\} \subset \mathcal{A}$  (with probability  $1 - o(1)$ ) which is required in ILAD(ii). Moreover, by the (shrinking) definition of  $\mathcal{A}$  we have the initial rate of ILAD(iv). Step 2 relies on Lasso. Condition HL is implied by Condition 1 and Lemma 2 applied twice with  $\zeta_i = v_i$  and  $\zeta_i = d_i$  under the condition that  $K_x^4 \log p \leq \delta_n n$ . By Lemma 8 we have  $\|x'(\hat{\theta} - \theta_0)\|_{2,n} \lesssim_P \sqrt{s \log(n \vee p)/n}$ . Moreover, by Lemma 9 we have  $\|\hat{\theta}\|_0 \lesssim s$  with probability  $1 - o(1)$ .

The rates established above for  $\hat{\theta}$  and  $\hat{\beta}$  imply (C.34) in ILAD(iii) since by Condition 1(ii)  $\mathbb{E}[|v_i| | x_i] \leq (\mathbb{E}[v_i^2 | x_i])^{1/2} = O(1)$  and  $\max_{1 \leq i \leq n} |x'_i(\hat{\theta} - \theta_0)| \lesssim_P K_x \sqrt{s^2 \log(p \vee n)/n} = o(1)$ .

To verify Condition ILAD(iii) (C.35), arguing as in the proof of Theorem 1, we can deduce that

$$\sup_{\alpha \in \mathcal{A}} |\mathbb{G}_n(\psi_{\alpha, \hat{\beta}, \hat{\theta}} - \psi_{\alpha, \beta_0, \theta_0})| = o_P(1).$$

This completes the proof. □

## APPENDIX I. PROOF OF AUXILIARY TECHNICAL RESULTS

*Proof of Lemma 2.* We shall use Lemma 10 ahead. Let  $Z_i = (x_i, \zeta_i)$  and define  $\mathcal{F} = \{f_j(x_i, \zeta_i) = x_{ij}^2 \zeta_i^2 : j = 1, \dots, p\}$ . Since  $\text{pr}(|X| > t) \leq \mathbb{E}[|X|^k]/t^k$ , for  $k = 2$  we have that  $\text{median}(|X|) \leq \sqrt{2\mathbb{E}[|X|^2]}$  and for  $k = q/4$  we have  $(1 - \tau)$ -quantile of  $|X|$  is bounded by  $(\mathbb{E}[|X|^{q/4}]/\tau)^{4/q}$ . Then we have

$$\max_{f \in \mathcal{F}} \text{median}(|\mathbb{G}_n(f(x_i, \zeta_i))|) \leq \sqrt{2\mathbb{E}[x_j^4 \zeta_i^4]} \leq K_x^2 \sqrt{2\mathbb{E}[\zeta^4]}$$

and

$$(1 - \tau)\text{-quantile of } \max_{j \leq p} \sqrt{\mathbb{E}_n[x_j^4 \zeta_i^4]} \leq (1 - \tau)\text{-quantile of } K_x^2 \sqrt{\mathbb{E}_n[\zeta^4]} \leq K_x^2 (\mathbb{E}[|\zeta|^q]/\tau)^{4/q}.$$

The conclusion follows from Lemma 10. □

*Proof of Lemma 3.* By the triangle inequality we have

$$\|x'(\hat{\beta}^{(2m)} - \beta_0)\|_{2,n} \leq \|x'(\hat{\beta} - \beta_0)\|_{2,n} + \|x'(\hat{\beta}^{(2m)} - \hat{\beta})\|_{2,n}.$$

Now let  $T^1$  denote the  $m$  largest components of  $\hat{\beta}$  and  $T^k$  corresponds to the  $m$  largest components of  $\hat{\beta}$  outside  $\cup_{d=1}^{k-1} T^d$ . It follows that  $\hat{\beta}^{(2m)} = \hat{\beta}_{T^1 \cup T^2}$ .

Next note that for  $k \geq 3$  we have  $\|\hat{\beta}_{T^{k+1}}\| \leq \|\hat{\beta}_{T^k}\|_1 / \sqrt{m}$ . Indeed, consider the problem  $\max\{\|v\|/\|u\|_1 : v, u \in \mathbb{R}^m, \max_i |v_i| \leq \min_i |u_i|\}$ . Given a  $v$  and  $u$  we can always increase the objective function by using  $\tilde{v} = \max_i |v_i|(1, \dots, 1)'$  and  $\tilde{u}' = \min_i |u_i|(1, \dots, 1)'$  instead. Thus, the maximum is achieved at  $v^* = u^* = (1, \dots, 1)'$ , yielding  $1/\sqrt{m}$ .

Thus, by  $\|\widehat{\beta}_{T^c}\|_1 \leq \mathbf{c}\|\delta_T\|_1$  and  $|T| = s$  we have

$$\begin{aligned} \|x'(\widehat{\beta}^{(2m)} - \widehat{\beta})\|_{2,n} &= \|x'_i \sum_{k=3}^K \widehat{\beta}_{T^k}\|_{2,n} \\ &\leq \sum_{k=3}^K \|x' \widehat{\beta}_{T^k}\|_{2,n} \leq \sqrt{\phi_{\max}(m)} \sum_{k=3}^K \|\widehat{\beta}_{T^k}\| \\ &\leq \sqrt{\phi_{\max}(m)} \sum_{k=2}^{K-1} \frac{\|\widehat{\beta}_{T^k}\|_1}{\sqrt{m}} \leq \sqrt{\phi_{\max}(m)} \frac{\|\widehat{\beta}_{(T^1)^c}\|_1}{\sqrt{m}} \\ &\leq \sqrt{\phi_{\max}(m)} \frac{\|\widehat{\beta}_{T^c}\|_1}{\sqrt{m}} \leq \sqrt{\phi_{\max}(m)} \mathbf{c} \frac{\|\delta_T\|_1}{\sqrt{m}}. \end{aligned}$$

□

*Proof of Lemma 4.* We will specify  $\ell_n$  below. Let  $\widetilde{K} = \max_{i \leq n} \|\tilde{x}_i\|_\infty$ . We have

$$\sup_{\|\delta\|_0 \leq \tilde{\ell}_n s} \left| \frac{\|\tilde{x}'\delta\|_{2,n}}{\|\tilde{x}'\delta\|_{P,2}} - 1 \right| \leq \{\bar{\phi}_{\min}(\tilde{\ell}_n s)\}^{-1} \sup_{\|\delta\|_0 \leq \tilde{\ell}_n s, \|\delta\|=1} |(\mathbb{E}_n - \mathbb{E})[(\tilde{x}'\delta)^2]| \delta_n.$$

By symmetrization for probabilities we have

$$\begin{aligned} &P(\sup_{\|\delta\|_0 \leq \tilde{\ell}_n s, \|\delta\|=1} |(\mathbb{E}_n - \mathbb{E})[(\tilde{x}'\delta)^2]| > t) \\ &\leq P(\sup_{\|\delta\|_0 \leq \tilde{\ell}_n s, \|\delta\|=1} |(\mathbb{E}_n - \mathbb{E})[(\tilde{x}'\delta)^2]| > t \mid \widetilde{K} \leq \bar{K}) + P(\widetilde{K} > \bar{K}) \\ &\leq 4P(\sup_{\|\delta\|_0 \leq \tilde{\ell}_n s, \|\delta\|=1} |\mathbb{E}_n[\varepsilon_i(\tilde{x}'\delta)^2]| > t/4 \mid \widetilde{K} \leq \bar{K}) + P(\widetilde{K} > \bar{K}) \end{aligned}$$

for any  $t \geq 2\bar{K}\sqrt{\tilde{\ell}_n s} \sqrt{\bar{\phi}_{\max}(\tilde{\ell}_n s)/n}$ . From Theorem 3.6 of [26], for

$$\delta_n = 2 \left( \bar{C}\bar{K}\sqrt{\tilde{\ell}_n s} \log(1 + \tilde{\ell}_n s) \sqrt{\log(p \vee n)} \sqrt{\log n} \right) / \sqrt{n},$$

where  $\bar{C}$  is the universal constant, we have

$$\mathbb{E} \left[ \sup_{\|\delta\|_0 \leq \tilde{\ell}_n s, \|\delta\|=1} |\mathbb{E}_n[\varepsilon_i(\tilde{x}'\delta)^2]| > t/4 \mid \widetilde{K} \leq \bar{K} \right] \leq \delta_n^2 + \delta_n \{\bar{\phi}_{\max}(\tilde{\ell}_n s)\}^{1/2}.$$

By Condition 1, we have  $\max_{1 \leq i \leq n} |d_i| \leq \tilde{\ell}_n n^{1/4}$  with probability  $1 - C/\tilde{\ell}_n$ . Setting  $\bar{K} := K_x \vee \tilde{\ell}_n n^{1/4}$  we have  $\widetilde{K} \leq \bar{K}$  with probability  $1 - \Delta_n - C/\tilde{\ell}_n$ . Next we show that  $\delta_n \rightarrow 0$ . Indeed, for  $\tilde{\ell}_n \rightarrow \infty$  slow enough, we have

$$\delta_n \lesssim \frac{\{K_x \vee \tilde{\ell}_n n^{1/4}\} \sqrt{\tilde{\ell}_n s} \log(1 + \tilde{\ell}_n s) \sqrt{\log(p \vee n)} \sqrt{\log n}}{\sqrt{n}} = \frac{K_x \vee \tilde{\ell}_n n^{1/4}}{n^{1/3}} \log^{3/2} n \sqrt{\frac{\tilde{\ell}_n s \log(p \vee n)}{n^{1/3}}} \rightarrow 0$$

since  $K_x^4 = o(n)$  and  $s \log(p \vee n) = o(n^{1/3})$ . □

## APPENDIX J. AUXILIARY PROBABILISTIC INEQUALITIES

Let  $Z_1, \dots, Z_n$  be independent random variables taking values in a measurable space  $(S, \mathcal{S})$ , and consider an empirical process  $\mathbb{G}_n(f) = n^{-1/2} \sum_{i=1}^n \{f(Z_i) - \mathbb{E}[f(Z_i)]\}$  indexed by a pointwise measurable class of functions  $\mathcal{F}$  on  $S$  (see [30], Chapter 2.3). Denote by  $\mathbb{P}_n$  the (random) empirical probability measure that assigns probability  $n^{-1}$  to each  $Z_i$ . Let  $N(\epsilon, \mathcal{F}, \|\cdot\|_{\mathbb{P}_n, 2})$  denote the  $\epsilon$ -covering number of  $\mathcal{F}$  with respect to the  $L^2(\mathbb{P}_n)$  seminorm  $\|\cdot\|_{\mathbb{P}_n, 2}$ .

The following maximal inequality is derived in [6].

**Lemma 10** (Maximal inequality for finite classes). *Suppose that the class  $\mathcal{F}$  is finite. Then for every  $\tau \in (0, 1/2)$  and  $\delta \in (0, 1)$ , with probability at least  $1 - 4\tau - 4\delta$ ,*

$$\max_{f \in \mathcal{F}} |\mathbb{G}_n(f)| \leq \left\{ 4\sqrt{2 \log(2|\mathcal{F}|/\delta)} Q(1 - \tau) \right\} \vee 2 \max_{f \in \mathcal{F}} \text{median}(|\mathbb{G}_n(f)|),$$

where  $Q(u) := u$ -quantile of  $\max_{f \in \mathcal{F}} \sqrt{\mathbb{E}_n[f(Z)^2]}$ .

The following maximal inequality is derived in [3].

**Lemma 11** (Maximal inequality for infinite classes). *Let  $F = \sup_{f \in \mathcal{F}} |f|$ , and suppose that there exist some constants  $\omega_n > 1$ ,  $\nu > 1$ ,  $m > 0$ , and  $h_n \geq h_0$  such that*

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

Set  $C := (1 + \sqrt{2\nu})/4$ . Then for every  $\delta \in (0, 1/6)$  and every constant  $K \geq \sqrt{2/\delta}$ , we have

$$\sup_{f \in \mathcal{F}} |\mathbb{G}_n(f)| \leq 4\sqrt{2}cKC \sqrt{m \log(n \vee h_n \vee \omega_n)} \max \left\{ \sup_{f \in \mathcal{F}} \sqrt{\mathbb{E}[f(Z)^2]}, \sup_{f \in \mathcal{F}} \sqrt{\mathbb{E}_n[f(Z)^2]} \right\},$$

with probability at least  $1 - \delta$ , provided that  $n \vee h_0 \geq 3$ ; the constant  $c < 30$  is universal.

## REFERENCES

- [1] R. Baraniuk, M. Davenport, R. DeVore, and M. Wakin. A simple proof of the restricted isometry property for random matrices. *Constructive Approximation*, 28:253–263, 2008.
- [2] A. Belloni, D. Chen, V. Chernozhukov, and C. Hansen. Sparse models and methods for optimal instruments with an application to eminent domain. *Econometrica*, 80(6):2369–2430, November 2012.
- [3] A. Belloni and V. Chernozhukov.  $\ell_1$ -penalized quantile regression for high dimensional sparse models. *Annals of Statistics*, 39(1):82–130, 2011.
- [4] A. Belloni and V. Chernozhukov. Least squares after model selection in high-dimensional sparse models. *Bernoulli*, 19(2):521–547, 2013.
- [5] A. Belloni, V. Chernozhukov, and I. Fernandez-Val. Conditional Quantile Processes based on Series or Many Regressors. *ArXiv e-prints*, May 2011.
- [6] A. Belloni, V. Chernozhukov, and C. Hansen. Inference on treatment effects after selection amongst high-dimensional controls. *ArXiv*, 2011.
- [7] A. Belloni, V. Chernozhukov, and C. Hansen. Inference for high-dimensional sparse econometric models. *Advances in Economics and Econometrics. 10th World Congress of Econometric Society, held in August 2010*, III:245–295, 2013.
- [8] A. Belloni, V. Chernozhukov, and K. Kato. Robust inference in high-dimensional sparse quantile regression models. *Working Paper*, 2013.
- [9] P. J. Bickel, Y. Ritov, and A. B. Tsybakov. Simultaneous analysis of Lasso and Dantzig selector. *Annals of Statistics*, 37(4):1705–1732, 2009.
- [10] V. Chernozhukov, D. Chetverikov, and K. Kato. Gaussian approximation of suprema of empirical processes. *arXiv preprint arXiv:1212.6885*, 2012.
- [11] V. Chernozhukov, D. Chetverikov, and K. Kato. Gaussian approximations and multiplier bootstrap for maxima of sums of high-dimensional random vectors. *Annals of Statistics*, 41(6):2786–2819, 2013.
- [12] Victor Chernozhukov and Christian Hansen. Instrumental variable quantile regression: A robust inference approach. *Journal of Econometrics*, 142:379–398, 2008.
- [13] Xuming He and Qi-Man Shao. On parameters of increasing dimensions. *J. Multivariate Anal.*, 73(1):120–135, 2000.
- [14] Bing-Yi Jing, Qi-Man Shao, and Qiyang Wang. Self-normalized Cramer-type large deviations for independent random variables. *Ann. Probab.*, 31(4):2167–2215, 2003.

- [15] K. Kato. Group lasso for high dimensional sparse quantile regression models. Preprint, ArXiv, 2011.
- [16] Roger Koenker. *Quantile regression*, volume 38 of *Econometric Society Monographs*. Cambridge University Press, Cambridge, 2005.
- [17] Michael R. Kosorok. *Introduction to Empirical Processes and Semiparametric Inference*. Series in Statistics. Springer, Berlin, 2008.
- [18] Sokbae Lee. Efficient semiparametric estimation of a partially linear quantile regression model. *Econometric theory*, 19:1–31, 2003.
- [19] Hannes Leeb and Benedikt M. Pötscher. Model selection and inference: facts and fiction. *Economic Theory*, 21:21–59, 2005.
- [20] Hannes Leeb and Benedikt M. Pötscher. Sparse estimators and the oracle property, or the return of Hodges’ estimator. *J. Econometrics*, 142(1):201–211, 2008.
- [21] J. Neyman. Optimal asymptotic tests of composite statistical hypotheses. In U. Grenander, editor, *Probability and Statistics, the Harold Cramer Volume*. New York: John Wiley and Sons, Inc., 1959.
- [22] J. Neyman.  $c(\alpha)$  tests and their use. *Sankhya*, 41:1–21, 1979.
- [23] J. L. Powell. Censored regression quantiles. *Journal of Econometrics*, 32:143–155, 1986.
- [24] Joseph P. Romano and Michael Wolf. Stepwise multiple testing as formalized data snooping. *Econometrica*, 73(4):1237–1282, July 2005.
- [25] M. Rudelson and S. Zhou. Reconstruction from anisotropic random measurements. *ArXiv:1106.1151*, 2011.
- [26] Mark Rudelson and Roman Vershynin. On sparse reconstruction from fourier and gaussian measurements. *Communications on Pure and Applied Mathematics*, 61:1025–1045, 2008.
- [27] R. Tibshirani. Regression shrinkage and selection via the lasso. *J. Roy. Statist. Soc. Ser. B*, 58:267–288, 1996.
- [28] S. A. van de Geer. High-dimensional generalized linear models and the lasso. *Annals of Statistics*, 36(2):614–645, 2008.
- [29] A. W. van der Vaart. *Asymptotic Statistics*. Cambridge Series in Statistical and Probabilistic Mathematics, 1998.
- [30] A. W. van der Vaart and J. A. Wellner. *Weak Convergence and Empirical Processes*. Springer Series in Statistics, 1996.
- [31] Lie Wang.  $L_1$  penalized lad estimator for high dimensional linear regression. *ArXiv*, 2012.
- [32] Cun-Hui Zhang and Stephanie S. Zhang. Confidence intervals for low-dimensional parameters with high-dimensional data. *ArXiv.org*, (arXiv:1110.2563v1), 2011.
- [33] S. Zhou. Restricted eigenvalue conditions on subgaussian matrices. *ArXiv:0904.4723v2*, 2009.