

# Model-Assisted Inference for Conditional Average Treatment Effects with High-Dimensional Controls\*

Adam Baybutt and Manu Navjeevan<sup>†</sup>

University of California, Los Angeles

Revised July 31, 2022

## Abstract

Plausible identification of conditional average treatment effects (CATE's) often relies on controlling for a large number of variables to account for confounding factors. In these high-dimensional settings, estimation of the CATE requires first estimating first-stage models. The resulting second-stage CATE estimator is said to be doubly robust if it remains consistent when either one of the first-stage models is consistently estimated. Inference procedures for this estimator are termed “model-assisted” if they are asymptotically valid when either one of the first stage models is consistently estimated. These properties are especially relevant in high-dimensional settings where consistency of the first stage estimators often relies on correctly specifying their parametric forms.

While doubly robust estimators of the CATE have been studied, inference based on these estimators is typically not model-assisted. Using the popular augmented inverse propensity weighting signal, we propose an estimator for the CATE that is doubly robust and whose resulting Wald-type confidence intervals are model-assisted. We assume a logistic model for the propensity score and a linear model for the outcome regression, and estimate the parameters of these models using an  $\ell_1$  (Lasso) penalty to address the high dimensional covariates. Our proposed estimator remains consistent at the nonparametric rate and our proposed pointwise and uniform confidence intervals remain asymptotically valid even if one of the logistic propensity score or linear outcome regression models are misspecified. These results are obtained under similar conditions to existing analyses in the high-dimensional and nonparametric literatures.

---

\*We are grateful to Denis Chetverikov, Andres Santos, Zhipeng Liao, Jinyong Hahn, Rosa Matzkin and participants in UCLA's Econometrics Proseminar for helpful comments.

<sup>†</sup>Corresponding author: [mnavjeevan@g.ucla.edu](mailto:mnavjeevan@g.ucla.edu).

## 1 Introduction

Consider a potential outcomes framework (Rubin, 1974, 1978) where an observed outcome  $Y \in \mathbb{R}$  and treatment  $D \in \{0, 1\}$  are related to two latent potential outcomes  $Y_1, Y_0 \in \mathbb{R}$  via  $Y = DY_1 + (1 - D)Y_0$ . To account for unobserved confounding factors a common strategy is to assume the researcher has access to a vector of covariates  $Z = (Z_1; X) \in \mathbb{R}^{d_z \times d_x}$  such that the potential outcomes are independent of the treatment decision after conditioning on the observed covariates,  $(Y_1, Y_0) \perp D | Z$ . In this setting, we are interested in estimation of and inference on the conditional average treatment effect (CATE):

$$E[Y_1 - Y_0 | X = x]. \quad (1.1)$$

Estimation of the CATE generally requires first fitting propensity score and/or outcome regression models. When the number of control variables  $Z$  is large ( $d_z \gg n$ ), as is often needed to justify conditional unconfoundedness, these first stage models are typically estimated using regularized methods which converge slowly and rely on the correctness of ad-hoc parametric specifications for consistency.<sup>1</sup>

Fortunately, so long as both models are correctly specified, one can obtain a nonparametric-rate consistent estimator and valid inference procedure for the CATE by using the popular augmented inverse propensity weighted (aIPW) signal (Semenova and Chernozhukov, 2021; Fan et al., 2022). This is because the aIPW signal obeys an orthogonality condition at the true nuisance model values that limits the first stage estimation error passed on to the second stage estimator. Moreover, estimators based on the aIPW signal are doubly robust; consistency of the resulting second-stage estimators requires correct specification of only one of the first stage propensity score or outcome regression models. However inference based on these estimators is not model-assisted. Under misspecification the aIPW signal orthogonality fails and the second stage CATE estimators converge slowly, rendering resulting testing procedures and confidence intervals invalid.

This paper proposes a model-assisted inference procedure for the conditional average treatment effect when the number of control variables  $d_z$  is potentially much larger than the sample size  $n$ . The dimensionality of the conditioning variable,  $d_x$ , will remain fixed in our analysis. Our approach is based on Tan (2020) wherein model-assisted inference is developed for the average treatment effect. Following Semanova and Chernozhukov (2021) we take a series approach to estimating the CATE, using a projection of aIPW signal onto a growing set of basis functions. By assuming a logistic form for the propensity score model and a linear form for the outcome regression model, we construct novel  $\ell_1$ -regularized first-stage estimating equations to recover a limited orthogonality of the aIPW signal at the limiting values of the first stage estimators. This restricted orthogonality is enough to achieve model-assisted pointwise and uniform inference; the second-stage estimator's pointwise and uniform confidence intervals are valid even if one of the logistic or linear functional forms is misspecified.

To achieve model-assisted inference at all points in the support of the conditioning variable, we must obtain this restricted orthogonality for each basis term in the series approximation. This is accomplished by employing distinct first-stage estimating equations for each basis term used in the second-stage series approximation. This results in the number of first-stage estimators growing with the number of basis terms. These estimators converge uniformly to limiting values under standard conditions in high-dimensional analysis. Improving on prior work in model-assisted inference, our  $\ell_1$  regularized first-stage estimation incorporates

<sup>1</sup>See Bauer and Kohler (2019); Schmidt-Hieber (2020) for some recent nonparametric results in high-dimensional settings using deep neural networks.

a data-dependent penalty parameter based on the work of [Chetverikov and Sørensen \(2021\)](#). This allows practical implementation of our proposed estimation procedure with minimal knowledge of the underlying data generating process.

The use of multiple pairs of nuisance parameter estimates limits our ability to straightforwardly apply existing nonparametric results for series estimators ([Newey, 1997](#); [Belloni et al., 2015](#)). Under modified conditions, we analyze the asymptotic properties of our second-stage series estimator to re-derive pointwise normality and strong gaussian process approximation results. These modified conditions are in general slightly stronger than those of [Belloni et al. \(2015\)](#), though in certain special cases collapse exactly to the conditions of [Belloni et al. \(2015\)](#).

**Prior Literature.** In the same setting as this paper, [Tan \(2020\)](#) considers estimation of the average treatment effect. After assuming a logistic form for the propensity score and a linear form for the outcome regression, [Tan \(2020\)](#) proposes  $\ell_1$ -regularized first-stage estimators that allow for partial control of the derivative of the aIPW signal away from true nuisance values and thus allow for model-assisted inference. [Smucler et al. \(2019\)](#) extends the analysis of [Tan \(2020\)](#) to consider model-assisted inference for a larger class of finite dimensional target parameters with bilinear influence functions.

Inference on low-dimensional target parameters after estimation of high-dimensional nuisance parameters has its origins in the work of [Belloni et al. \(2012, 2013\)](#), who consider inference on causal parameters after estimation of nuisance parameters using Lasso and Post-Lasso methods. [Chernozhukov et al. \(2018\)](#) analyze estimation of low-dimensional target parameters using score functions that are Neyman-orthogonal with respect to nuisance parameters. They show that the resulting target parameter estimate is  $\sqrt{P_n}$ -consistent and asymptotically normal so long as the nuisance parameters are consistent at rate  $n^{-1/4}$ , a condition satisfied by many machine learning-based estimators.

[Stone \(1985\)](#) provides a primary analysis of nonparametric series estimators. [Newey \(1997\)](#) derives convergence rates and pointwise normality results which are updated by [Belloni et al. \(2015\)](#). Recently, [Semenova and Chernozhukov \(2021\)](#) consider series estimation of functional target parameters after high-dimensional nuisance estimation using the aIPW signal. [Fan et al. \(2022\)](#) provides a similar analysis using a second stage kernel estimator. The resulting second stage estimators are doubly robust but do not necessarily yield model-assisted inference. [Wu et al. \(2021\)](#) provide model-assisted inference procedures for covariate-specific treatment effects with discrete conditioning variables; their results depend on exact representation assumptions that are unlikely to hold with continuous covariates. Moreover, no uniform inference procedures are described.

Regularization using an  $\ell_1$ -penalty in linear regression models (Lasso) was first introduced by [Tibshirani \(1996\)](#) for selecting relevant features out of a high-dimensional set of controls. Near oracle rates of convergence for these estimators under sparsity are given by [Bickel et al. \(2009\)](#); [Belloni and Chernozhukov \(2013\)](#). [Chetverikov et al. \(2021\)](#) provide interesting new results with cross validated penalties. These works can be extended to consider rates of convergence for general  $\ell_1$ -penalized M-estimators as seen in [van der Greer \(2016\)](#); [Tan \(2017\)](#); [Chetverikov and Sørensen \(2021\)](#). In particular, [Chetverikov and Sørensen \(2021\)](#) propose a data-driven “bootstrap after cross-validation” approach to penalty parameter selection that is modified for and implemented in our setting.

**Paper Structure.** This paper proceeds as follows. Section 2 defines the problem and introduces methods for estimation and inference. Section 3 introduces, with exposition for intuition, the main result: valid inference for the second-stage series estimator if either the first-stage logistic propensity score model or linear outcome regression model is correctly specified. Section 4 shows uniform convergence of the first-stage estimators and describes the required conditions. Section 5 presents a full set of second-stage asymptotic results. Section 6 ties up a technical detail. Section 7 provides evidence from a simulation study while Section 8 applies our proposed estimator to examine the effect of maternal smoking on infant birth weight. Section 9 concludes. Proofs of main results are deferred to Appendix A.

**Notation.** For any measure  $F$  and any function  $f$ , define the  $L^2$  norm,  $\|f\|_{F,2} = (\mathbb{E}_F[f^2])^{1/2}$  and the  $L^1$  norm  $\|f\|_{F,1} = \mathbb{E}_F|f|$ . For any vector in  $\mathbb{R}^p$  let  $\|k\|_p$  for  $p \geq [1; \infty]$  denote the  $\ell_p$  norm,  $\|k\|_p = (\sum_{i=1}^p |a_i|^p)^{1/p}$  and  $\|k\|_1 = \sum_{i=1}^p |a_i|$ . If the subscript is unspecified, we are using the  $\ell_2$  norm. For two vectors  $a, b \in \mathbb{R}^p$ , let  $a \cdot b = (a_i b_i)_{i=1}^p$  denote the Hadamard (element-wise) product. We adopt the convention that for  $a \in \mathbb{R}^p$  and  $c \in \mathbb{R}$ ,  $a + c = (a_i + c)_{i=1}^p$ . For a matrix  $A \in \mathbb{R}^{m \times n}$  let  $\|A\|_F$  denote the operator norm and  $\|A\|_1 = \sum_{j=1}^n \sum_{i=1}^m |A_{ij}|$ . For any real valued function  $f$  let  $\mathbb{E}_n[f] = \frac{1}{n} \sum_{i=1}^n f(X_i)$  denote the empirical expectation and  $G_n[f] = \frac{1}{\sqrt{n}} \sum_{i=1}^n (f(X_i) - \mathbb{E}[f])$  denote the empirical process. For two sequences of random variables  $\{a_n\}_{n \in \mathbb{N}}$  and  $\{b_n\}_{n \in \mathbb{N}}$ , we say  $a_n = o_p(b_n)$  or  $a_n = O_p(b_n)$  if  $a_n/b_n \rightarrow 0$  in probability and say  $a_n = o_p(b_n)$  if  $a_n/b_n \rightarrow 0$ .

## 2 Setup and Methods

We will assume that the researcher observes i.i.d data and that, as discussed, conditioning on  $Z$  is sufficient to control for all confounding factors affecting both the treatment decision  $D$  and the potential outcomes,  $Y_1$  and  $Y_0$ .

**Assumption 2.1** (Identification).

- (i)  $(Y_1, D, Z, g_{i=1}^n)$  are independent and identically distributed.
- (ii)  $(Y_1, Y_0) \perp D \mid Z$ .
- (iii) There exists a value  $\alpha \in (0, 1)$  such that  $\alpha < \mathbb{E}[D \mid Z = z] < 1 - \alpha$  almost surely in  $Z$ .

Our analysis allows the dimensionality of  $Z$  to grow much faster than sample size ( $d_Z \ll n$ ), while assuming the dimensionality of  $X$  remains fixed ( $d_X \ll n$ ). The parameter of interest is the conditional average treatment effect:  $\mathbb{E}[Y_1 - Y_0 \mid X = x]$ . However, for this paper we largely focus on estimation and inference for the conditional average counterfactual outcome:

$$g_0(x) := \mathbb{E}[Y_1 \mid X = x]; \tag{2.1}$$

A model-assisted inference procedure for the other conditional counterfactual outcome,  $\mathbb{E}[Y_0 \mid X = x]$ , follows a similar procedure and is described in Section 6. The procedures can be combined for model-assisted estimation and inference of the CATE.

We utilize series estimation targeting  $g_0(x)$ . The augmented inverse propensity weighted (aIPW) signal we employ,

$$Y(\cdot; m) = \frac{DY}{(Z)} - \frac{D}{(Z)} + m(Z); \tag{2.2}$$

is a function of a fitted propensity score model,  $\hat{p}(Z)$ ; and a fitted outcome regression model,  $m(Z)$ . After estimating these two first-stage models, whose true values are given  $p^*(Z) := \Pr(D = 1 | Z)$  and  $m^*(Z) := E[Y | D = 1; Z]$ , the target parameter  $\theta_0(x)$  is estimated by taking a quasi-projection of the aIPW signal (2.2) onto a growing set of  $k$ , weakly positive, basis terms

$$p^k(x) := \rho_1(x); \dots; \rho_k(x) \in \mathbb{R}_+^k. \quad (2.3)$$

The basis terms are required to be weakly positive as they will be used as weights within the first-stage estimators (2.5)-(2.6). Negative weights break the convexity of the objective function. Examples of weakly positive basis functions are B-splines or Fourier series basis terms<sup>1</sup>. To ensure that the basis terms are well behaved, we will make assumptions on  $k_{1,1} := \sup_{x \in \mathcal{X}} k p^k(x) k_1$ ,  $k_{1,2} := \sup_{x \in \mathcal{X}} k p^k(x) k_2$ , and the eigenvalues of the design matrix  $Q := E[\rho^k(x) \rho^k(x)^\top]$ .

**First-Stage Setup and Estimation.** We estimate a separate propensity score model and outcome regression model for each basis term. Following Tan (2020) we assume a logistic regression form for the propensity score model and a linear form for the outcome regression model:

$$\begin{aligned} \hat{p}(Z; \beta) &= 1 + \exp(-\beta^\top Z)^{-1}; \\ m(Z; \gamma) &= \beta^\top Z. \end{aligned} \quad (2.4)$$

To deal with the high-dimensional  $Z$ , we implement  $\ell_1$ -penalized optimization to estimate the parameters of the models. Importantly, these first stage estimating equations are designed as in Tan (2020) so that their first order conditions directly limit the bias passed on to the second-stage series estimator, as will be described in Section 3. For each  $j = 1; \dots; k$ ; the parameters of (2.4),  $\beta_j \in \mathbb{R}^{d_Z}$ ; are estimated by

$$\hat{\beta}_j := \arg \min E_n[\rho_j(X) f D e^{-\beta_j^\top Z} + (1 - D) \beta_j^\top Z g] + \lambda_j k_1; \quad (2.5)$$

$$\hat{\beta}_j := \arg \min E_n[\rho_j(Z) D e^{-\beta_j^\top Z} (Y - \beta_j^\top Z)^2] + \lambda_j k_1; \quad (2.6)$$

The penalty parameters  $\lambda_j$  and  $\lambda_j$ , for each  $j = 1; \dots; k$ , will be chosen via a data dependent technique described in Section 4.1. Under standard assumptions the nuisance estimators  $\hat{\beta}_j; \hat{\beta}_j$  will converge uniformly over  $j = 1; \dots; k$  to the population minimizers

$$\bar{\beta}_j := \arg \min E[\rho_j(X) f D e^{-\beta_j^\top Z} + (1 - D) \beta_j^\top Z g]; \quad (2.7)$$

$$\bar{\beta}_j := \arg \min E[\rho_j(Z) D e^{-\beta_j^\top Z} (Y - \beta_j^\top Z)^2]. \quad (2.8)$$

We will assume that the population minimizers are sufficiently sparse. To control the sparsity, define  $S_{\lambda_j} := \{j : \bar{\beta}_j \neq 0\}$ ,  $S_{\lambda_j} := \{j : \bar{\beta}_j \neq 0\}$  and  $s_k := \max_{1 \leq j \leq k} |S_{\lambda_j} \cap S_{\lambda_j}|$ . Uniform convergence of the parameter estimators (2.5) and (2.6) to their population minimizers (2.7) and (2.8) will ensure uniform convergence of the estimated propensity score and outcome regression models to limiting values

$$\hat{\beta}_j(Z) := (Z; \hat{\beta}_j) \rightarrow_p \bar{\beta}_j(Z) := (Z; \bar{\beta}_j); \quad \hat{m}_j(Z) := m(Z; \hat{\beta}_j) \rightarrow_p \bar{m}_j(Z) := m(Z; \bar{\beta}_j);$$

<sup>1</sup>B-splines are linear transformations of regression splines with lower multicollinearity. See De Boor (2001) for more information on B-splines.

If the logistic propensity score model (linear outcome regression model) is correctly specified, the limiting propensity score model (outcome regression model) will coincide with the true value, i.e. for all  $j = 1; \dots; k$   $\bar{m}_j(Z) = m_j(Z)$  ( $\bar{m}_j(Z) = m_j(Z)$ ). In general however they may differ. As discussed, our inference results will remain valid so long as one of either the logistic propensity score model or the linear outcome regression model is correctly specified.

**Second Stage Setup and Estimation.** Under Assumption 2.1, the aIPW signal  $Y(\cdot; \cdot)$  is doubly robust, that is for square integrable  $\bar{m}_j \in \mathcal{M}$  and  $m \in \mathcal{M}$ ,

$$\begin{aligned} E[Y_1 | X = x] &= E[Y(\cdot; m^?) | X = x] \\ &= E[Y(\cdot; m^?) | X = x] \\ &= E[Y(\cdot; m) | X = x]: \end{aligned} \quad (2.9)$$

This double robustness of the aIPW signal allows for identification of the CATE even if only one of the nuisance models is correctly specified. If either  $\bar{m}_j = m_j^?$  or  $\bar{m}_j = m_j^?$ , we can write for all  $j = 1; \dots; k$ :

$$\begin{aligned} Y(\cdot; \bar{m}_j) &= g_0(x) + r_j; \quad E[Y(\cdot; \bar{m}_j) | X = x] = 0 \\ &= g_k(x) + r_k(x) + r_j \end{aligned} \quad (2.10)$$

where  $g_0(x)$  is the conditional counterfactual outcome (2.1),  $g_k(x) := p^k(x)^{\theta^k}$  is the projection of  $g_0(x)$  onto the first  $k$  basis terms, and  $r_k(x) := g_0(x) - g_k(x)$  denotes the approximation error from this projection. The least squares parameter  $\theta^k$  governing the projection in  $g_k(x)$  can be identified by the projection of the aIPW signal onto the basis terms  $p^k(x)$ :<sup>2</sup>

$$\begin{aligned} \theta^k &:= Q^{-1} E[p^k(X) Y_1] \\ &= Q^{-1} E[p^k(X) Y(\cdot; m^?)] \\ &= Q^{-1} E[p^k(X) Y(\cdot; \bar{m}_j)]; \quad \theta^j = 1; \dots; k: \end{aligned} \quad (2.11)$$

To estimate  $\theta^k$  we combine the  $k$  estimates according to:

$$\hat{\theta}^k = \hat{\mathcal{Q}}^{-1} E_n \begin{bmatrix} p_1(X) Y(b_1; \bar{m}_1) \\ \vdots \\ p_k(X) Y(b_k; \bar{m}_k) \end{bmatrix}; \quad (2.12)$$

where  $\hat{\mathcal{Q}} := E_n[p^k(x)p^k(x)^{\theta^k}]$ . Our final estimator for  $g_0(x)$  is  $\hat{g}(x) := p^k(x)^{\hat{\theta}^k}$ . Note the separate error terms for each  $j = 1; \dots; k$  in (2.10), which are collected together in the vector  $\delta^k := (\delta_1; \dots; \delta_k)$ .

**Inference Procedure.** Having defined our estimator  $\hat{g}(x)$  we now describe a typical inference procedure. The variance of our second stage estimator is governed by the following variance matrices:

$$\begin{aligned} \tilde{\Omega} &:= Q^{-1} E[\hat{r} p^k(x) (\theta^k + r_k) g^{\theta^k} p^k(x) (\theta^k + r_k) g^{\theta^k}] Q^{-1} \\ \Omega_0 &:= Q^{-1} E[\hat{r} p^k(x) \theta^k g^{\theta^k} p^k(x) \theta^k g^{\theta^k}] Q^{-1} \end{aligned} \quad (2.13)$$

<sup>2</sup>The second and third equalities of (2.11) hold under Assumption 2.1 for any  $\bar{m}_j \in \mathcal{M}$  and  $m \in \mathcal{M}$  by double robustness of the aIPW signal (2.9).

where  $\odot$  represents the Hadamard (element-wise) product and, abusing notation somewhat, for a vector  $a \in \mathbb{R}^k$  and scalar  $c \in \mathbb{R}$  we let  $a + c = (a_i + c)_{i=1}^k$ . We will estimate these matrices using the plug-in empirical analog

$$\hat{\mathbf{b}} := \hat{\mathbf{b}}^{-1} E_n[\hat{f} \hat{p}^k(x) \hat{b}^k g \hat{f} \hat{p}^k(x) \hat{b}^k g^T] \hat{\mathbf{b}}^{-1} \quad (2.14)$$

for  $\hat{b}^k := (\hat{b}_1; \dots; \hat{b}_k)$  and  $\hat{b}_j := Y(\hat{b}_j; \hat{m}_j) - \hat{g}(x)$  for  $j = 1; \dots; k$ . Define  $\hat{b}(x) := k \hat{\mathbf{b}}^{-1=2} \hat{p}^k(x) k = \hat{P} \hat{\mathbf{b}}^{-1}$  and let the test statistic  $t(x)$  be given:

$$t(x) := \frac{\hat{g}(x) - g(x)}{\hat{b}(x)} \quad (2.15)$$

Denote the bootstrapped test statistic as

$$\hat{t}^b(x) := \frac{\hat{p}^k(x) \hat{\mathbf{b}}^{-1=2}}{k \hat{\mathbf{b}}^{-1=2} \hat{p}^k(x) k} N_k^b \quad (2.16)$$

where  $N_k^b$  is a bootstrap draw from  $N(0; I_k)$ . Define the confidence bands

$$\hat{l}(x); \hat{r}(x) := \hat{g}(x) - c^2(1 - \alpha) \hat{b}(x); \hat{g}(x) + c^2(1 - \alpha) \hat{b}(x) \quad (2.17)$$

The critical value  $c^2(1 - \alpha)$  is either  $c_p^2(1 - \alpha)$ , the  $(1 - \alpha)$  quantile of a  $N(0; 1)$  distribution, for pointwise inference or  $c_u^2(1 - \alpha)$ , the  $(1 - \alpha)$ -quantile of  $\sup_{x \in \mathcal{X}} \hat{t}_n^b(x)$ , for uniform inference. Below, we will show that under our construction of  $\hat{g}(x)$ , this standard inference procedure will remain valid even under misspecification of either nuisance model.

### 3 Theory Overview

We now present the main results of the paper and give some intuition into how they are obtained.

#### 3.1 Main Results

**Proposition 3.1** (Pointwise Normality). *Suppose that the conditions of Theorems 4.3 and 5.4 hold. If additionally the first stage sparsity bounds of (4.10) and (4.13) hold, then so long as either the logistic propensity score model or linear outcome regression model is correctly specified we have for any  $\alpha \in S^{k-1}$ ,*

$$\hat{P} \frac{\hat{b}^k - \alpha \hat{b}^k}{k \hat{\mathbf{b}}^{-1=2} k} \xrightarrow{d} N(0; 1) \quad (3.1)$$

Moreover, for any  $x \in \mathcal{X}$  and  $\hat{b}(x) := \hat{\mathbf{b}}^{-1=2} \hat{p}^k(x)$ ,

$$\hat{P} \frac{\hat{p}(x) \hat{b}^k - \alpha \hat{p}(x) \hat{b}^k}{k \hat{b}(x) k} \xrightarrow{d} N(0; 1) \quad (3.2)$$

If the approximation error is negligible relative to the estimation error, namely  $\hat{P} \hat{r}_k(x) = o(k \hat{\mathbf{b}}^{-1=2} \hat{p}^k(x) k)$ , then

$$\hat{P} \frac{\hat{g}(x) - g(x)}{k \hat{b}(x) k} \xrightarrow{d} N(0; 1) \quad (3.3)$$

**Proposition 3.2** (Uniform Inference). *Suppose that the conditions of Theorems 4.3 and 5.5 hold. If additionally the first stage sparsity bounds of (4.10) and (4.13) hold, then so long as either the logistic propensity*

score model or linear regression model is correctly specified the uniform confidence intervals stated in (2.17) satisfy

$$\Pr \left( g(x) \in [\underline{l}(x); \bar{l}(x)] : \delta x \in X \right) = 1 - o_p(1)$$

Proposition 3.1 states that our method yields model-assisted pointwise inference, that is an inference procedure based on the test statistic  $t(x)$  in (2.15) remains valid when using the pointwise critical value  $c_p^2(1 - \alpha)$  even if either the logistic propensity score model or linear outcome regression model is misspecified, but not both. Proposition 3.2 gives conditions for model-assisted uniform inference and shows that the uniform confidence intervals described in (2.17) remain valid under misspecification when using the uniform critical value  $c_U^2(1 - \alpha)$ .

In Section 5 we provide further results for the nonparametric estimator  $\hat{g}(x)$ , including pointwise and uniform convergence rates. These results will rely on first-stage convergence results established in Section 4. For now, we give some intuition for how model-assisted confidence bands are obtained for the nonparametric estimator.

### 3.2 Managing First-Stage Bias

We focus on controlling the size of the vector  $\mathbf{B}^k$ , defined in (3.4), which measures the bias passed on from first-stage estimation to the second-stage estimate  $\hat{b}^k$ . Limiting the size of  $\mathbf{B}^k$  is crucial in showing convergence of  $\hat{b}^k$  to the true parameter  $b^k$  and thus consistency of the nonparametric estimator  $\hat{g}(x)$ .

$$\mathbf{B}^k := E_n \begin{pmatrix} \rho_1(X) & Y(\hat{\tau}_1; \hat{m}_1) & Y(\tau_1; m_1) \\ \vdots & \vdots & \vdots \\ \rho_k(X) & Y(\hat{\tau}_k; \hat{m}_k) & Y(\tau_k; m_k) \end{pmatrix} \quad (3.4)$$

For exposition, we will consider a single term of (3.4),  $\mathbf{B}_j^k$ , which roughly measures the first stage estimation bias taken on from adding the  $j^{\text{th}}$  basis term to our series approximation of  $g_0(x)$ . The discussion that follows is a bit informal, instead of considering the derivatives with respect to the true parameters below our proof strategy will directly use the Kuhn-Tucker conditions of the optimization routines in (2.5)-(2.6). However, the general intuition is the same as is used in the proofs.

In addition to the doubly robust identification property (2.9), the aIPW signal is typically useful in the high-dimensional setting because it obeys an orthogonality condition:<sup>1</sup>

$$E[\tau(Y(\tau; m^?) - jZ)] = E[\tau_m(Y(\tau; m^?) - jZ)] = 0 \quad (3.5)$$

When both the propensity score model and outcome regression model are correctly specified we can (loosely speaking) examine the bias  $\mathbf{B}_j^k$  by replacing  $\tau_j = \tau^?$  and  $m_j = m^?$  and considering the following first order

<sup>1</sup>Robustness and orthogonality are indeed closely related, see Theorem 6.2 in Newey and McFadden (1994) for a discussion. The orthogonality condition is sometimes called *Neyman Orthogonality* and is discussed in depth in Chernozhukov et al. (2018); Semenova and Chernozhukov (2021).



Taylor expansion:

$$\begin{aligned} \mathbf{B}_j^k &= E_n[\rho_j(X)Y(\hat{\tau}_j; \hat{m}_j)] - E_n[\rho_j(X)Y(\tau_j^*; m_j^*)] \\ &= E_n[\rho_j(X) \underbrace{r_{\tau_j^*, m_j^*} Y(\tau_j^*; m_j^*)}_{O_p(n^{-1/2}) \text{ by (3.5)}}] \begin{matrix} \hat{\tau}_j \\ \hat{m}_j \end{matrix} \begin{matrix} \tau_j^* \\ m_j^* \end{matrix} + o_p(n^{-1/2}). \end{aligned} \quad (3.6)$$

By orthogonality of the aIPW signal the gradient term is close to zero, which guarantees that the bias is  $\sqrt{n}$ -asymptotically negligible, i.e.  $o_p(n^{-1/2})$ , even if the nuisance parameters converge slowly to the true values,  $\tau_j^*$  and  $m_j^*$ .<sup>2</sup> Since the bias disappears faster than the rate of convergence in Proposition 3.1, the researcher can ignore first stage nuisance parameter estimation error and treat  $\tau_j^*$  and  $m_j^*$  as known when analyzing the asymptotic properties of the nonparametric estimator.

Indeed since  $X$  is a subvector of  $Z$ , under mild conditions aIPW signal orthogonality yields that the gradient  $E_n[\rho_j(X) \underbrace{r_{\tau_j^*, m_j^*} Y(\tau_j^*; m_j^*)}_{O_p(n^{-1/2}) \text{ by (3.5)}}]$  is asymptotically bounded uniformly for all  $j = 1, \dots, k$ . In this case, we would not need separate nuisance estimators for each basis function. The same consistent estimators  $\hat{\tau}$  and  $\hat{m}$  could be used at each element of  $\mathbf{b}^k$  and the analysis of (3.6) would limit the size of the entire vector  $\mathbf{B}^k$ . This is the approach followed by [Semenova and Chernozhukov \(2021\)](#).

However, if the functional form of either the propensity score model or outcome regression model is misspecified, this approach does not work. Consider the first order Taylor expansion around the probability limits  $\tau_j$  and  $m_j$ .

$$\begin{aligned} \mathbf{B}_j^k &= E_n[\rho_j(X)Y(\hat{\tau}_j; \hat{m}_j)] - E_n[\rho_j(X)Y(\tau_j; m_j)] \\ &= E_n[\rho_j(X) \underbrace{r_{\tau_j, m_j} Y(\tau_j; m_j)}_{O_p(n^{-1/2})}] \begin{matrix} \hat{\tau}_j \\ \hat{m}_j \end{matrix} \begin{matrix} \tau_j \\ m_j \end{matrix} + o_p(n^{-1/2}) \end{aligned} \quad (3.7)$$

So long as either one of  $\tau_j = \tau_j^*$  or  $m_j = m_j^*$ , the double robustness of the aIPW signal (2.9) still allows for identification of  $E[\rho_j(X)Y_1]$ .<sup>3</sup> However, the aIPW orthogonality tells us nothing about the expectation of the gradient away from the true parameters,  $\tau_j^*; m_j^*$ ; if either  $\tau_j \neq \tau_j^*$  or  $m_j \neq m_j^*$  there is no reason to believe that the gradient on the right hand side of (3.7) is mean zero. In general, the bias  $\mathbf{B}_j^k$  will then diminish at the rate of convergence of our nuisance parameters. Because we have high dimensional controls, this rate of convergence will generally be much slower than the parametric rate needed to establish Proposition 3.1.

To get around this, we follow a similar approach as [Tan \(2020\)](#) by designing the first-stage objective functions (2.5)-(2.6) such that the resulting first-order conditions control the bias passed on to the second stage. Since we have assumed specific parametric forms for our propensity score and outcome regression models we will consider expansions around the parameters themselves.

$$\begin{aligned} \mathbf{B}_j^k &= E_n[\rho_j(X)Y(\hat{\tau}_j; \hat{m}_j)] - E_n[\rho_j(X)Y(\tau_j; m_j)] \\ &= E_n[\rho_j(X) \underbrace{r_{\tau_j, m_j} Y(\tau_j; m_j)}_{O_p(n^{-1/2})}] \begin{matrix} \hat{\tau}_j \\ \hat{m}_j \end{matrix} \begin{matrix} \tau_j \\ m_j \end{matrix} + o_p(n^{-1/2}) \end{aligned} \quad (3.8)$$

<sup>2</sup>Typically all that is required is that  $k \hat{\tau}_j - \tau_j^k = o_p(n^{-1/4})$  and  $k \hat{m}_j - m_j^k = o_p(n^{-1/4})$  in order to make the second order remainder term  $\sqrt{n}$ -negligible

<sup>3</sup>That is  $E_n[\rho_j(X)Y(\tau_j; m_j)] = E_n[\rho_j(X)Y_1]$

After substituting the forms of  $\bar{y}_j(z) = (z; \bar{y}_j)$  and  $\bar{m}_j(z) = m(z; \bar{y}_j)$  described in (2.4) we obtain

$$Y(\bar{y}_j; \bar{m}_j) = DY(1 + e^{\beta_j Z}) - D(1 + e^{\beta_j Z})^{-\theta} Z +^{-\theta} Z:$$

Differentiating with respect to the parameters  $\bar{y}_j$  and  $\bar{y}_j$  gives

$$E[\rho_j(X) r_{j; j} Y(\bar{y}_j; \bar{m}_j)] = E \left[ \frac{\rho_j(X) D e^{\beta_j Z} (Y -^{-\theta} Z) Z^{\#}}{\rho_j(X) f D(1 + e^{\beta_j Z}) Z + Z g} \right] \quad (3.9)$$

However, by definition  $\bar{y}_j$  and  $\bar{y}_j$  solve the minimization problems defined in (2.7)-(2.8), the population analogs of our finite sample estimating equations. The first order conditions of these minimization problems yield that

$$E \left[ \frac{\rho_j(X) D e^{\beta_j Z} (DY -^{-\theta} Z) Z}{\rho_j(X) f D(1 + e^{\beta_j Z}) Z + Z g} \right] = 0 \quad (3.10)$$

$\left. \begin{array}{l} \text{First order condition of } \bar{y}_j \\ \hline \rho_j(X) f D(1 + e^{\beta_j Z}) Z + Z g \\ \hline \rho_j(X) D e^{\beta_j Z} (DY -^{-\theta} Z) Z \\ \hline \text{First order condition of } \bar{y}_j \end{array} \right\}$

Examining the first order conditions in (3.10), we see that they exactly give us control over the gradient (3.9). This gives us that the bias examined in expansion (3.8) is  $\rho_{\bar{n}}$ -asymptotically negligible even under misspecification of the propensity score or outcome regression models.<sup>4</sup>

Notice that the orthogonality condition obtained in (3.10) is weaker than that of (3.5). Orthogonality is obtained with respect to the parameters of the logistic propensity score model and linear outcome model rather than the nuisance functions themselves. The bias is controlled only when approaching the limiting values along specific parametric paths; hence the term model-assisted inference. Because we are assuming constant parametric models, this is exactly what is needed to recover valid inference.

Control of this gradient under misspecification is not provided using other estimating equations, such as maximum likelihood for the logistic propensity score model or ordinary least squares for the linear outcome regression model. Moreover, control over the gradient of  $\mathbf{B}_j^k$  from (3.4) is not provided by the first-order conditions for  $\bar{y}_l$  and  $\bar{y}_l$  for  $l \neq j$ :

$$E[\rho_j(X) r_{j; j} Y(\bar{y}_j; \bar{m}_j)] = E \left[ \frac{\rho_j(X) D e^{\beta_j Z} (Y -^{-\theta} Z) Z^3}{\rho_j(X) f D(1 + e^{\beta_j Z}) Z + Z g} \right] \quad (3.11)$$

$\left. \begin{array}{l} \text{First order condition of } \bar{y}_l \\ \hline \rho_l(X) f D(1 + e^{\beta_l Z}) Z + Z g \\ \hline \rho_l(X) D e^{\beta_l Z} (Y -^{-\theta} Z) Z \\ \hline \text{First order condition of } \bar{y}_l \end{array} \right\}$

Obtaining valid inference at all points  $x \in X$  under misspecification as in Proposition 3.1 requires showing  $\rho_{\bar{n}}$ -negligible first stage estimation bias for any linear combination of the vector (3.4). As outlined above, this requires using  $k$  separate pairs of nuisance parameter estimator to obtain  $k$  separate pairs of first order conditions, one for each term of the vector.

<sup>4</sup>Correct specification is only necessary for identification.

Of course it is not enough to analyze each term of  $\mathbf{B}^k$  individually. Showing Proposition 3.1 and the rest of our nonparametric results we will require that these terms, and thus the nuisance parameter estimators, converge uniformly. Section 4 provides conditions under which this uniform convergence is obtained.

## 4 First-Stage Convergence

In order to show nonparameteric rate convergence of the estimator described in Section 2, we seek to establish the following condition, which limits the bias passed from first-stage estimation onto the second-stage estimator.

**Condition 1** (No Effect of First-Stage Bias).

$$\max_{1 \leq j \leq k} E_n[\rho_j(X)Y(b_j; \hat{m}_j)] - E_n[\rho_j(X)Y(\cdot; \bar{m}_j)] = o_p(n^{-1/2}k^{-1/2}). \quad (4.1)$$

In particular, Condition 1 implies that the bias vector  $\mathbf{B}^k$  from (3.4) satisfies  $k\mathbf{B}^k k_2 = o_p(n^{-1/2})$ . Establishing this condition largely follows the logic laid out in Section 3 after establishing the uniform consistency of the first-stage estimators,  $b_j$  and  $\hat{b}_j$ , described in (2.5)-(2.6), to their corresponding probability limits. This section will establish constraints under which this uniform first-stage convergence is obtained. These provisions are similar to those existing in the  $\lambda_1$ -penalized estimation literature (Bickel et al., 2009; Belloni and Chernozhukov, 2013; Chetverikov and Sørensen, 2021; Bühlmann and van de Geer, 2011).

### 4.1 Penalty Parameter Selection

We improve on existing work on model-assisted inference by incorporating a data dependent two step penalty parameter selection procedure based on the work of Chetverikov and Sørensen (2021). Before presenting uniform convergence results it is useful to describe this procedure.

For each  $j = 0; 1; \dots; k$ ; we start with a pilot penalty parameters given by

$$\hat{c}_j^{\text{pilot}} = c_j \frac{\ln^3(d_Z)}{n} \quad \text{and} \quad \bar{c}_j^{\text{pilot}} = \bar{c}_j \frac{\ln^3(d_Z)}{n} \quad (4.2)$$

for some constants  $c_j; \bar{c}_j$  selected from the interval  $[\underline{c}_n; \bar{c}_n]$  with  $\underline{c}_n > 0$ . Using  $\hat{c}_j^{\text{pilot}}$  and  $\bar{c}_j^{\text{pilot}}$  in lieu of  $c_j$  and  $\bar{c}_j$  in (2.5)-(2.6) we generate pilot estimators  $\hat{b}_j^{\text{pilot}}$  and  $\bar{b}_j^{\text{pilot}}$ . These pilot estimators are used to generate plug in estimators  $\hat{\theta}_j$  and  $\bar{\theta}_j$  of the residuals

$$\begin{aligned} \hat{\theta}_j &:= \rho_j(X) f D e^{\hat{b}_j^{\text{pilot}} Z} + (1 - D) g \\ \bar{\theta}_j &:= \rho_j(X) D e^{\bar{b}_j^{\text{pilot}} Z} (Y - \bar{b}_j^{\text{pilot}} Z); \end{aligned} \quad (4.3)$$

We then use a multiplier bootstrap procedure to select our final penalty parameters  $\hat{c}_j$  and  $\bar{c}_j$ .

$$\begin{aligned} \hat{c}_j &= c_0 \quad (1 - \alpha) \text{-quantile of } \max_{1 \leq l \leq d_Z} \int E_n[e_i \hat{\theta}_j Z_l | j] \text{ given } f Y_i; D_i; Z_i g_{i=1}^n; \\ \bar{c}_j &= \bar{c}_0 \quad (1 - \alpha) \text{-quantile of } \max_{1 \leq l \leq d_Z} \int E_n[e_i \bar{\theta}_j Z_l | j] \text{ given } f Y_i; D_i; Z_i g_{i=1}^n \end{aligned} \quad (4.4)$$

where  $e_1; \dots; e_n$  are independent standard normal random variables generated independently of the data

$f(Y_i; D_i; X_i)g_{i=1}^n$  and  $c_0 > 1$  is a fixed constant<sup>1</sup>. In line with other work we take  $c_0 = 1.1$  in simulations. So long as our residual estimates converge in empirical mean square to limiting values,

$$\begin{aligned} U_{\cdot j} &:= p_j(X) f D e^{-j^{\text{pilot}} Z} + (1 - D) g \\ U_{\cdot j} &:= p_j(X) D e^{-j^{\text{pilot}} Z} (Y - e^{-j^{\text{pilot}} Z}); \end{aligned} \quad (4.5)$$

the choice of penalty parameter in (4.4) will ensure that the penalty parameter dominates the noise with high probability. This allows for consistent variable selection and coefficient estimation.

**Remark 4.1** (Pilot Penalty Selection). In practice, the researcher has a fair bit of flexibility in choosing the constants  $c_{\cdot j}$  and  $c_{\cdot j}$  in (4.2). The optimal choice of these constants may depend on the underlying data degenerating process. We recommend using cross validation to pick these constants from a fixed-cardinality set of possible values. In line with Assumption 4.1(v), the values in the set should be chosen to be on the order of the maximum value of  $k p^k(X_i) k_{\cdot 1}$  observed in the data. So long as this set of possible values is bounded and of cardinality on an order less than  $n$ , the results of this section hold for a cross-validated choice of these constants.

## 4.2 Uniform First-Stage Convergence

With the penalty parameter selection procedure established, we now demonstrate uniform convergence of our first-stage estimators. To do so, we will rely on the following first-stage assumptions. The first assumption will allow us to generate a uniform bound on the convergence rate of the logistic propensity score model parameters.

**Assumption 4.1** (Logistic Model Convergence).

- (i) The regressors  $Z$  are bounded,  $\max_{1 \leq i \leq d_z} |Z_{ij}| \leq C_0$  almost surely.
- (ii) There is a constant  $B_0$  such that  $\|Z_{\cdot j}\| \leq B_0$  almost surely for all  $j$ .
- (iii) There exist fixed constants  $\epsilon_0 > 1$  and  $1 > \epsilon_0 > 0$  such that for each  $j = 1, \dots, l$  the following empirical compatibility condition holds for the empirical hessian matrix  $\tilde{\Sigma}_{\cdot j} := E_n[De^{-jZ} Z Z^{\top}]$ . For any  $b \geq R^{d_z}$ 

$$\prod_{1 \leq j \leq d_z} |b_{ij}| \leq \epsilon_0 \prod_{1 \leq j \leq d_z} |b_{ij}| \Rightarrow \frac{1}{\epsilon_0} \prod_{1 \leq j \leq d_z} |b_{ij}|^2 \leq |S_{\cdot j}| b^{\epsilon_0 \tilde{\Sigma}_{\cdot j}} b;$$
- (iv) There exist fixed constants  $c_U$  and  $C_U > 0$  such that for all  $j = 1, \dots, k$ ,  $E[U_{\cdot j}^4] \leq (c_U k_{\cdot 1} C_U)^4$  and  $\min_{1 \leq i \leq d_z} E[U_{\cdot j}^2 Z_{ij}^2] \geq c_U$ .
- (v) The following sparsity bounds hold with  $s_{k, \cdot} := \max_{1 \leq j \leq k} |S_{\cdot j}|$ .

$$\frac{k_{\cdot 1}}{c_n \ln(d_z n)} \leq 0; \quad \frac{k_{\cdot 1} s_{k, \cdot} \tilde{c}_n^2 \ln^5(d_z n)}{n} \leq 0; \quad \text{and} \quad \frac{4}{k_{\cdot 1}} \frac{\ln^7(d_z k n)}{n} \leq 0;$$

The first part of Assumption 4.1 assumes that the regressors are bounded. It can be relaxed with sufficient assumptions on the tail behavior of the controls. We should note that compactness of  $X$  is generally required

<sup>1</sup>The constant  $c_0$  can actually be different for the propensity score model and for the outcome regression model, and can also vary for each  $j = 1, \dots, k$ . All that matters is that each constant satisfies the requirements on  $c_0$  in Theorem 4.3. This complicates notation, however, so the constant is assumed to be the same across estimation procedures.

by nonparametric estimators. The second assumption bounds all limiting propensity scores  $\pi_j(Z)$  away from zero uniformly. The third assumption is an empirical compatibility condition on the weighted first-stage design matrix. It is slightly weaker than the restricted eigenvalue conditions often assumed in the literature (Bickel et al., 2009; Belloni et al., 2012). The penultimate condition is an identifiability constraint that limits the moments of the noise and bounds it away from zero uniformly over all estimation procedures. The last condition requires the validity of the bootstrap penalty parameter selection procedure and is comparable to the requirements needed for the bootstrap after cross validation technique described by Chetverikov and Sørensen (2021). The main difference is the additional assumption on the growth rate of the basis functions,  $k_{:,1}$  which is to ensure uniform stability of the estimation procedures (2.5)-(2.6).

The next set of assumptions enable us to generate a uniform bound on the convergence rate of the linear outcome regression model parameters. The assumptions needed for this convergence are a bit stronger than those needed for the propensity score model convergence in order to manage the reliance of outcome regression estimates on the propensity model parameter estimates in (2.6).

**Assumption 4.2** (Linear Model Convergence).

- (i) The errors  $Y_1 - m_j(Z)$  are uniformly subgaussian conditional on  $Z$  in the following sense. There exist fixed positive constants  $D_0$  and  $D_1$  such that for any  $j$ :

$$D_0 \mathbb{E} \exp \left( \frac{1}{D_0} \sum_{i=1}^k \beta_i Y_1 - m_j(Z) \right) \leq D_1 \exp \left( \frac{1}{D_1} \sum_{i=1}^k \beta_i^2 \right)$$

almost surely.

- (ii) For any  $j$  the following empirical compatibility condition holds with fixed constants  $\rho_0 > 1$  and  $\rho_1 > 0$  for the empirical hessian matrix  $\tilde{\Sigma}_{:,j}$ :

$$\lambda_{\min}(\tilde{\Sigma}_{:,j}) \geq \rho_0 \lambda_{\min}(\tilde{\Sigma}_{:,j}) \Rightarrow \frac{1}{\rho_0} \sum_{i=1}^k \beta_i^2 \leq \sum_{i=1}^k \beta_i \tilde{\Sigma}_{:,j} \beta_i$$

- (iii) For all  $j \leq k$ ,  $\mathbb{E}[U_{:,j}^4] \leq (k_{:,1} C_U)^4$  and  $\min_{1 \leq i \leq d_z} \mathbb{E}[U_{:,j}^2 Z_i^2] \geq c_U$ .

- (iv) The following sparsity bound holds with  $s_{k,:} = \max_{1 \leq j \leq k} \sum_{i=1}^k \beta_i$ :

$$\frac{2}{k_{:,1}} s_{k,:}^2 \leq \frac{c_n^2 \ln^5(d_z n)}{n} \leq 0;$$

The first condition of Assumption 4.2 assumes that tail behavior of the outcome regression errors are uniformly thin. The second and third conditions are the same empirical compatibility condition and identification conditions as in Assumption 4.1 but applied to the parameters and residuals for the outcome regression. The final part of Assumption 4.2 is a sparsity bound needed for validity of the penalty parameter selection procedure and mirrors those of Chetverikov and Sørensen (2021).

**Remark 4.2** (Constants). Many of the constants in Assumptions 4.1 and 4.2 are assumed to be fixed across all  $j$ . This is mainly to simplify the exposition of the results below. In practice, all constants can be allowed to grow slowly as the number of basis terms increases. However, the growth rate of these terms affects the first-stage sparsity bounds needed to show Condition 1.

We can now establish uniform convergence of the propensity score model and the outcome regression model parameters. The convergence will be presented in terms of the  $\ell_1$  norm of the coefficient vector as well as the symmetrized Bregman divergences,  $D^z_{\cdot j}(b_j; \bar{\cdot}_j)$  and  $D^z_{\cdot j}(b_j; \bar{\cdot}_j; \bar{g}_j)$ , defined by

$$\begin{aligned} D^z_{\cdot j}(b_j; \bar{\cdot}_j) &:= E_n \left[ \rho_j(X) Df e^{b_j^0 Z} - e^{b_j^0 Z} g^0 \right] Z^{\top} b_j^0 Z^{\top} g^0 ; \\ D^z_{\cdot j}(b_j; \bar{\cdot}_j; \bar{g}) &:= E_n \left[ \rho_j(X) Df e^{b_j^0 Z} - e^{b_j^0 Z} (g^0 - \bar{g}) \right] Z^{\top} b_j^0 Z^{\top} (g^0 - \bar{g})^2 ; \end{aligned} \quad (4.6)$$

We refer readers to discussion in [Tan \(2017\)](#) for details and motivation. For now it suffices to note that the Bregman divergence is the error resulting from approximating the non-penalized loss function at the estimated value with a first order Taylor expansion of the non-penalized loss function at the true values. Because our loss functions are convex, these errors will always be positive. Bounds on the Bregman divergence help directly control second order terms in the remainder of (3.8).

Theorems 4.1 and 4.2 establish uniform convergence of the propensity score and outcome regression parameters at nearly the near-oracle rate from [Belloni and Chernozhukov \(2013\)](#). Theorem 4.2 relies on a condition on the ratio between the penalty parameters that is discussed in Remark 4.3.

**Theorem 4.1** (Propensity Score Model). *Suppose that Assumption 4.1 holds, that  $\lambda_0 > (c_0 + 1)\lambda = (c_0 - 1)\lambda$ ,  $k = n \lambda \rightarrow 0$ , and that  $k \rightarrow 0$  as  $n \rightarrow \infty$ . Then the logistic propensity score model parameters are uniformly bounded:*

$$\sup_{1 \leq j \leq k} D^z_{\cdot j}(b_j; \bar{\cdot}_j) \leq P_{k,1}^2 \frac{S_k \ln(d_z)}{n} \quad \text{and} \quad \sup_{1 \leq j \leq k} \|b_j - \bar{g}_j\|_1 \leq P_{k,1} \sqrt{\frac{S_k^2 \ln(d_z)}{n}} \quad (4.7)$$

**Theorem 4.2** (Outcome Regression Model). *Suppose that Assumptions 4.1 and 4.2 hold. In addition assume that  $\lambda_0 > (c_0 + 1)\lambda = (c_0 - 1)\lambda$ ,  $k = n \lambda \rightarrow 0$ ,  $k \rightarrow 0$ , and that there is a fixed constant  $c > 0$  such that for all  $j$ ,  $\|g_j - \bar{g}_j\|_1 \leq c$ . Then the linear outcome regression model parameters are uniformly bounded:*

$$\sup_{1 \leq j \leq k} D^z_{\cdot j}(b_j; \bar{\cdot}_j; \bar{g}_j) \leq P_{k,1}^2 \frac{S_k \ln(d_z)}{n} \quad \text{and} \quad \sup_{1 \leq j \leq k} \|b_j - \bar{g}_j\|_1 \leq P_{k,1} \sqrt{\frac{S_k^2 \ln(d_z)}{n}} \quad (4.8)$$

**Remark 4.3** (Penalty Parameter Ratio). The existence of a fixed constant  $c > 0$  such  $\|g_j - \bar{g}_j\|_1 \leq c$  may seem a bit unnatural, but it can be enforced in practice by setting the penalty to instead be, say,  $\lambda_j^{\text{ratio}} := \max\{\lambda_j, c\}$ . The use of this alternative penalty parameter instead of  $\lambda_j$  does not effect any of the convergence results established in this document. In simulations, we find this constraint is rarely binding.

With the uniform convergence of the logistic regression parameters established, we can apply the logic of Section 3 to show convergence of the entire bias vector  $\mathbf{B}^k$  (3.4). Again, it should be emphasized that the result of Theorem 4.3 could not be obtained using standard maximum likelihood estimating equations for the logistic propensity score parameters or on ordinary least squares loss function for the outcome regression parameters. Moreover, the rate of convergence of Theorem 4.3 would not be obtained by using the estimating equations of [Tan \(2020\)](#). It is the particular form of the estimating equations in (2.5)-(2.6) that allows us to show a suitable rate of convergence.

**Theorem 4.3** (First-Stage Convergence). *Suppose that Assumptions 4.1 and 4.2. In addition assume that  $\lambda_0 > (c_0 + 1)\lambda = (c_0 - 1)\lambda$ ,  $k = n \lambda \rightarrow 0$ ,  $k \rightarrow 0$ , and there is a fixed constant  $c > 0$  such that for all  $j$ ,  $\|g_j - \bar{g}_j\|_1 \leq c$ .*

Then the elements of the vector (3.4) converge uniformly:

$$\max_{1 \leq j \leq k} E_n[\rho_j(X)Y(b_j; \mathfrak{m}_j)] - E_n[\rho_j(X)Y(\bar{b}_j; \bar{\mathfrak{m}}_j)] \leq P \frac{s_{k,1}^2 \ln(d_z)}{n} \quad (4.9)$$

In particular, Condition 1 holds under the following additional sparsity bound

$$\frac{s_{k,1}^2 k^{1=2} \ln(d_z)}{P n} \leq 0: \quad (4.10)$$

For implementable inference using  $\hat{\Omega}$  to estimate the second-stage variance as in Proposition 3.1, we additionally wish to establish the following condition uniformly bounding the rate of convergence of the second moments for elements  $\mathbf{B}_j^k$  of the bias vector  $\mathbf{B}^k$ . Under Condition 2 and additional second-stage assumptions, Theorem 5.4 shows that the plug-in estimator for the matrices described in (2.13) is consistent.

**Condition 2** (Variance Estimation). Let  $m > 2$  be as in Assumption 5.3. Then,

$$s_{k,1} \max_{1 \leq j \leq k} E_n[\rho_j(X)^2 (Y(b_j; \mathfrak{m}_j) - Y(\bar{b}_j; \bar{\mathfrak{m}}_j))^2] = o_p(k^{-2} n^{1-m}) \quad (4.11)$$

**Theorem 4.4** (Second Moment Convergence). Suppose that the conditions of Theorem 4.3 hold. Then the elements of the vector (3.4) converge uniformly in empirical mean square at the following rate:

$$\max_{1 \leq j \leq k} E_n[\rho_j^2(X) (Y(b_j; \mathfrak{m}_j) - Y(\bar{b}_j; \bar{\mathfrak{m}}_j))^2] \leq P \frac{s_{k,1}^4 \ln(d_z)}{n} \quad (4.12)$$

In particular, Condition 2 holds if additionally

$$\frac{s_{k,1}^5 s_k^2 k^2 \ln(d_z)}{n^{(m-1)=m}} \leq 0: \quad (4.13)$$

**Remark 4.4.** The sparsity bound and conditions of Theorem 5.4 can be significantly weakened if the researcher is willing to use a cross fitting procedure, using one sample to estimate the nuisance parameters and another to evaluate the aIPW signal. This is because we could more directly follow Semenova and Chernozhukov (2021) and use that with high probability the quantities:

$$\max_{1 \leq i \leq n} \max_{1 \leq j \leq k} |Y_i(b_j; \mathfrak{m}_j) - Y_i(\bar{b}_j; \bar{\mathfrak{m}}_j)| \text{ and } \max_{1 \leq i \leq n} \max_{1 \leq j \leq k} |Y_i(b_j; \mathfrak{m}_j) - Y_i(\bar{b}_j; \bar{\mathfrak{m}}_j)|^2$$

are controlled by bounds that converge to zero. This could then be used in the proof of the matrix estimation result to show that the first-stage estimation error passed on to the variance estimator is negligible. Unfortunately, because we are not using cross-fitting, we cannot easily say anything about the behavior of the random variable  $Y$  on the event that  $k b_j - \bar{b}_j k$  and  $k \mathfrak{m}_j - \bar{\mathfrak{m}}_j k$  are small. Because the difference  $Y_i(b_j; \mathfrak{m}_j) - Y_i(\bar{b}_j; \bar{\mathfrak{m}}_j)$  depends on  $Y$ , not knowing its expectation on the event that  $b_j$  is close to  $\bar{b}_j$  and  $\mathfrak{m}_j$  is close to  $\bar{\mathfrak{m}}_j$  is problematic.

**Remark 4.5** (Alternative Penalty Parameter Selection). The additional sparsity assumptions needed by Assumptions 4.1 and 4.2 are stronger than what is required by Tan (2020) who only requires  $S \sqrt{\ln(d_z)/n} \leq 0$ , where  $S$  denotes the maximum sparsity index. This is due to the use of the data dependent penalty parameter

selection procedure described in Section 4.1 as opposed to the theoretical penalty parameter of Tan (2020). Under a similar theoretical penalty parameter we could establish Condition 1 under only the sparsity bound in (4.10). However, in line with other work, we find that in practice the data dependent method performs better than a theoretical penalty and only requires the researcher have minimal knowledge of the underlying data generating process.

For various reasons, the researcher may not want to implement the bootstrap penalty parameter procedure described in Section 4.1. If this is the case, we note that the pilot penalty parameters of (4.2) can be used directly. The constants  $c_{\cdot j}$  and  $\bar{c}_{\cdot j}$  can be selected via cross validation from a growing set  $\Lambda_n$   $[\underline{c}_n; \bar{c}_n]$ . If the penalty parameters are selected in this fashion, Conditions 1 and 2 can be obtained under Assumption 4.1(i)-(iv), Assumption 4.2(i)-(iii) and the following modified rate conditions:

$$\frac{k_j \Lambda_n j}{n^2} \not\rightarrow 0; \frac{k_{\cdot 1}}{\underline{c}_n \ln(d_z n)} \not\rightarrow 0; \frac{\bar{c}_n^2 s_k k^{1=2} \ln^3(d_z n)}{\rho \bar{n}} \not\rightarrow 0; \text{ and } \frac{k_{\cdot 1}^2 \bar{c}_n^2 s_k k^2 \ln^3(d_z n)}{n^{(m-1)=m}} \not\rightarrow 0 \quad (4.14)$$

These conditions are comparable to the rates in Assumptions 4.1 and 4.2 and equations (4.10)-(4.13). However, the rate of convergence on the bias established in Theorem 4.3 is somewhat slower. Appendix D provides details for implementation as well as formally shows the bounds in (4.14).

## 5 Additional Second-Stage Results

In this section, we present a set consistency and distributional results for our second-stage estimator  $\hat{g}(x)$  described in Section 2. Additional results, including pointwise and uniform linearization lemmas, can be found in Appendix C. These results will be presented under the assumption that Conditions 1-2 from Section 4 hold and for the most part will follow existing results in the nonparametric literature (Newey, 1997; Belloni et al., 2015). The main remaining difficulty is managing the  $k$  different error terms presented in (2.10). This requires some modification to the second-stage results.

In the standard setting, consistency of nonparametric estimator relies on certain conditions on the error terms. In our setting, we will require that these assumptions hold uniformly over  $k$  the error terms. We note though that there is a non-trivial dependence structure between that limiting aIPW signals. This strong dependence gives plausibility to our uniform conditions. For example, if the logistic propensity score model is correctly specified and the limiting outcome regression models are uniformly bounded conditional on  $Z$ , our conditions reduce exactly to the conditions of Belloni et al. (2015). In general, however, the uniform conditions suggest that a degree of undersmoothing is optimal when implementing our estimation procedure.

### 5.1 Pointwise Inference

Pointwise inference will rely on the following assumptions in tandem with Condition 1. Assumption 5.1 is all that is needed for the estimator to achieve the nonparametric rate of convergence even under misspecification, while Assumption 5.2 will ensure that the second stage estimator is pointwise asymptotically linear (see Appendix C).

**Assumption 5.1** (Second-Stage Consistency Assumptions).

- (i) Uniformly over all  $n$ , the eigenvalues of  $Q = E[\rho^k(x) \rho^k(x)^\eta]$  are bounded from above and away from zero.



- (ii) The conditional variance of the error terms is uniformly bounded in the following sense. There exist constants  $\bar{c}_k^2$  and  $\underline{c}_k^2$  such that for any  $j = 1, 2, \dots, k$  we have that  $\bar{c}_k^2 \text{Var}(r_k(x)) \leq \underline{c}_k^2 < 1$ ;
- (iii) For each  $n$  and  $k$  there are finite constants  $c_k$  and  $\bar{c}_k$  such that for each  $f \in \mathcal{F}$

$$\|r_k\|_{L^2} = (E[r_k(x)^2])^{1/2} \leq c_k \text{ and } \|r_k\|_{L^1} = \sup_{x \in X} |r_k(x)| \leq \bar{c}_k.$$

**Assumption 5.2** (Uniform Integrability). Let  $\bar{c}_k := \max_{1 \leq j \leq k} \bar{c}_{kj}$ . Assume that

- (i)  $\sup_{x \in X} E[\bar{c}_k^2 \mathbf{1}_{\bar{c}_k > \bar{c}_k} | \mathcal{F}_k(x)] \rightarrow 0$  as  $n \rightarrow \infty$  for any  $\bar{c}_k > 0$ .
- (ii)  $\sup_{x \in X} E[\bar{c}_k^2 \mathbf{1}_{\bar{c}_k > \bar{c}_k} | \mathcal{F}_k(x)] \rightarrow 0$  as  $n \rightarrow \infty$  for any  $\bar{c}_k > 0$ .

As mentioned, note that these are exactly the conditions required by [Belloni et al. \(2015\)](#), with the modification that the bounds on conditional variance and other moment conditions on the error term hold uniformly over  $j = 1, \dots, k$ . The assumptions on the series terms being used in the approximation can be shown to be satisfied by a number of commonly used functional bases, such as polynomial bases or splines, under adequate normalizations and smoothness of the underlying regression function. Readers should refer to [Newey \(1997\)](#), [Chen \(2007\)](#), or [Belloni et al. \(2015\)](#) for a more in depth discussion of these assumptions.<sup>1</sup>

**Theorem 5.1** (Pointwise Normality). *Suppose that Condition 1, Assumption 5.1, and Assumption 5.2 holds. In addition suppose that  $\frac{1}{k} \log k \rightarrow 0$ . Then so long as either the logistic propensity score model or linear outcome regression model is correctly specified, for any  $s \in \mathcal{S}^k$ :*

$$\frac{1}{\bar{n}} \frac{\partial \log p^k(x)}{\partial \Omega} \xrightarrow{d} N(0; 1) \quad (5.1)$$

where generally  $\Omega = \tilde{\Omega}$  but if  $\bar{c}_k \neq 0$  then we can set  $\Omega = \Omega_0$ . Moreover, for any  $x \in X$  and  $s(x) := \Omega^{1/2} p^k(x)$ ,

$$\frac{1}{\bar{n}} \frac{\partial \log p^k(x)}{\partial s(x)} \xrightarrow{d} N(0; 1) \quad (5.2)$$

and if the approximation error is negligible relative to the estimation error, namely  $\bar{n} r_k(x) = o(\bar{n} s(x))$ , then

$$\frac{1}{\bar{n}} \frac{\partial \log g(x)}{\partial s(x)} \xrightarrow{d} N(0; 1) \quad (5.3)$$

In conjunction with the verification of Condition 1 established in Section 4, Theorem 5.1 verifies the validity of the pointwise inference procedure defined in Section 2 based on the critical values  $\hat{c}_p^2(1 - \alpha)$  up to procurement of a consistent estimator for the variance matrix  $\Omega$ : Theorem 5.4 will provide conditions under which the plug in estimator  $\hat{\Omega}$  defined in (2.14) is consistent.

## 5.2 Uniform Convergence

Next, we turn to strengthening the pointwise results to hold uniformly over all points  $x \in X$ . This will require stronger conditions. We will make the following assumptions on the tail behavior of the error terms which strengthen Assumption 5.2.

<sup>1</sup>In practice, we recommend the use of B-splines in order to satisfy the first requirement that the basis functions are weakly positive and to reduce instability of the convex optimization programs described in (2.5)-(2.6).

**Assumption 5.3** (Uniform Limit Theory). Let  $\bar{k} = \sup_{x \in \mathcal{X}} \sum_{j=1}^k |j_j|$ ,  $\rho^k(x) = \rho^k(x)/k$ , and let

$$\frac{L}{k} := \sup_{\substack{x: x^0 \in \mathcal{X} \\ x \notin x^0}} \frac{k(x) - (x^0)k}{kx - x^0k}$$

Further for any integer  $s$  let  $\bar{k}^s = \sup_{x \in \mathcal{X}} E[j_k^s | \mathcal{X} = x]$ . For some  $m > 2$  assume

- (i) The regression errors satisfy  $\sup_{x \in \mathcal{X}} E[\max_{1 \leq i \leq n} |j_{ki}|^m | \mathcal{X} = x] \leq P n^{1-m}$
- (ii) The basis functions are such that (a)  $\frac{2^{m-(m-2)} \log k}{k} = n^{-1}$ , (b)  $(\frac{2}{k} - \frac{2^m}{k}) \log \frac{L}{k} \leq \log k$ , and (c)  $\log \frac{2^m}{k} \leq \log k$ .

As before, Assumption 5.3 is very similar to its analogue in Belloni et al. (2015), with the modification that the conditions are required to hold for  $\bar{k}$  as opposed to  $k$ . Under this assumption, we derive doubly-robust uniform rates of convergence and model-assisted uniform inference procedures for the conditional counterfactual outcome  $g_0(x)$ .

**Theorem 5.2** (Uniform Rate of Convergence). *Suppose that Condition 1 and Assumptions 5.1-5.3 hold. Then so long as either the propensity score model or outcome regression model is correctly specified:*

$$\sup_{x \in \mathcal{X}} (x)^0 G_n[\rho^k(x) - k] \leq P \frac{1}{\log k} \tag{5.4}$$

Moreover, for

$$\begin{aligned} \bar{R}_{1n} &:= \frac{r}{k} \frac{2 \log k}{n} (n^{1-m} P \frac{1}{\log k} + \frac{P}{k} \bar{k} c_k) \\ \bar{R}_{2n} &:= P \frac{1}{\log k} \bar{k} c_k \end{aligned}$$

we have that

$$\sup_{x \in \mathcal{X}} \rho^k(x)^0 (b^k - k) \leq P \frac{k}{n} P \frac{1}{\log k} + \bar{R}_{1n} + \bar{R}_{2n} \tag{5.5}$$

and

$$\sup_{x \in \mathcal{X}} g(x) - g_0(x) \leq P \frac{k}{n} P \frac{1}{\log k} + \bar{R}_{1n} + \bar{R}_{2n} + \bar{k} c_k \tag{5.6}$$

Theorem 5.2 gives conditions under which our estimator converges in probability to the true conditional counterfactual outcome  $g_0(x)$ . In particular, this convergence happens uniformly at the rates defined in (5.5)-(5.6). If these two terms go to zero, the entire estimator will converge uniformly to the true conditional expectation of interest.

**Theorem 5.3** (Strong Approximation by a Gaussian Process). *Assume that Condition 1 holds and that Assumptions 5.1-5.3 hold with  $m = 3$ . In addition assume that (i)  $\bar{R}_{1n} = o_p(a_n^{-1})$  and (ii)  $a_n^6 k^4 (\frac{2}{k} - \frac{3}{k} + \frac{3}{k} c_k^2)^2 \log^2 n = n^{-1} \rightarrow 0$ . Then so long as either the propensity score model or outcome regression model is correctly specified, for some  $N_k \rightarrow N(0; I_k)$ :*

$$\frac{P}{n} \frac{(x)^0 (b - k)}{k (x)^0 \Omega^{1=2} k} = d \frac{(x)^0 \Omega^{1=2}}{k (x)^0 \Omega^{1=2} k} N_k + o_p(a_n^{-1}) \text{ in } \cdot^1(x) \tag{5.7}$$

so that for  $s(x) := \Omega^{1-2} p^k(x)$

$$\rho_{\frac{1}{n}} \frac{p^k(x)^{\theta} (b)}{k s(x)^k} =_d \frac{s(x)}{k s(x)^k} N_k + o_p(a_n^{-1}) \text{ in } \cdot^1(X) \quad (5.8)$$

and if  $\sup_{x \in \mathcal{X}} \rho_{\frac{1}{n}} \frac{p^k(x)^{\theta} (b)}{k s(x)^k} = o(a_n^{-1})$ , then

$$\rho_{\frac{1}{n}} \frac{g(x) - \bar{g}(x)}{k s(x)^k} =_d \frac{s(x)^{\theta}}{k s(x)^k} N_k + o_p(a_n^{-1}) \text{ in } \cdot^1(X) \quad (5.9)$$

where in general we take  $\Omega = \tilde{\Omega}$  but if  $\bar{R}_{2n} = o_p(a_n^{-1})$  then we can set  $\Omega = \Omega_0$  where  $\tilde{\Omega}$  and  $\Omega_0$  are as in (2.13).

Theorem 5.3 establishes conditions under which uniform confidence intervals can be created for the function  $g_0(x)$ . As noted by Belloni et al. (2015), this is distinctly different from a Donsker type convergence result for the estimator  $\hat{g}(x)$  as viewed as a random element of  $\cdot^1(X)$ . In particular, the covariance kernel is left completely unspecified and in general need not be well behaved.

### 5.3 Matrix Estimation and Uniform Inference

The final second-stage results presented establish the consistency of various matrix estimations and validity of the uniform confidence intervals described in Section 2. The matrix estimation results will allow for consistent estimation of the asymptotic variances in Theorem 5.1 and Theorem 5.3 and thus allow for practically implementable, asymptotically valid confidence intervals for  $g_0(x)$ , even under misspecification of one of the nuisance parameter models.

The matrices  $Q = E[p^k(x)p^k(x)^{\theta}]$  and  $\Sigma = E[\tilde{r}p^k(x) \quad (k+r_k)g \quad \tilde{r}p^k(x) \quad (k+r_k)g]$  are estimated using

$$\begin{aligned} \hat{Q} &= E_n[p^k(x)p^k(x)^{\theta}] \\ \hat{\Sigma} &= E_n[\tilde{r}p^k(x) \quad b^k g \tilde{r}p^k(x) \quad b^k g^{\theta}] \end{aligned} \quad (5.10)$$

where  $b^k := (b_1, \dots, b_k)$  and  $\tilde{r}_j = Y(b_j; \mathbf{m}_j) - \bar{g}(x)$  for each  $j = 1, \dots, k$ . In addition define  $v_n = (E[\max_{1 \leq j \leq k} \tilde{r}_j^2])^{1/2}$ .

**Theorem 5.4** (Matrix Estimation). *Suppose that Conditions 1 and 2 and Assumptions 5.1-5.3 hold. In addition, assume that  $\bar{R}_{1n} + \bar{R}_{2n} = o_p((\log k)^{1/2})$ . Then, so long as either the propensity score model or outcome regression model is correctly specified:*

$$\begin{aligned} \hat{Q} - Q &= o_p\left(\frac{\sqrt{\frac{2 \log k}{n}}}{n}\right) = o(1) \text{ and} \\ \hat{\Sigma} - \Sigma &= o_p\left(\frac{v_n \sqrt{\frac{2 \log k}{n}}}{n}\right) = o(1) \end{aligned}$$

Moreover, for  $\hat{\Omega} = \hat{Q}^{-1} \hat{\Sigma} \hat{Q}^{-1}$ :

$$\hat{\Omega} - \Omega = o_p\left(\frac{v_n \sqrt{\frac{2 \log k}{n}}}{n}\right) = o(1)$$

We now establish the validity of uniform inference based on the gaussian bootstrap critical values  $c_U^2(1 - \alpha)$

defined in Section 2.

**Theorem 5.5** (Validity of Uniform Confidence Bands). *Suppose Conditions 1 and 2 are satisfied and Assumptions 5.1{5.3 hold with  $m \geq 4$ . In addition suppose (i)  $R_{1n} + R_{2n} = o(\log^{1=2} n)$ , (ii)  $\|k\| \log^2 n = n^{1=2 - 1/m} = o(1)$ , (iii)  $\sup_{x \in \mathcal{X}} |r_k(x) - k p^k(x)| = o(\log^{-1=2} n)$ , and (iv)  $k^4 \frac{2}{k} (1 + l_k^3 r_k^3)^2 \log^5 n = n = o(1)$ . Then for  $t(x)$  defined in (2.15), so long as either the propensity score model or outcome regression model is satisfied*

$$\Pr \left( \sup_{x \in \mathcal{X}} |t(x) - c^2(1 - \alpha)| \leq c^2(1 - \alpha) + o(1) \right) = 1 - \alpha + o(1)$$

As a result, uniform confidence intervals formed in (2.17) satisfy

$$\Pr(g(x) \in [\underline{g}(x); \bar{g}(x)] : x \in \mathcal{X}) = 1 - \alpha + o(1)$$

The width of the confidence bands  $2c^2(1 - \alpha)b(x)$  satisfy

$$2c^2(1 - \alpha)b(x) = O_p \left( \frac{1}{\log n} \cdot \sqrt{\frac{2 \log n}{n}} \right)$$

In conjunction with theorems 4.3, 4.4 and 5.4, this shows Proposition 3.2. It also shows that the confidence intervals based on the uniform critical values  $c_{\alpha}^2(1 - \alpha)$  shrink at the uniform nonparametric rate derived in Theorem 5.2.

## 6 Estimation of the Conditional Average Treatment Effect

Up to now, we have mainly focused on doubly-robust estimation and model-assisted inference for the function

$$g_0(x) = E[Y_1 | X = x]$$

We conclude by noting that we can use a symmetric procedure to obtain model-assisted inference for the additional conditional counterfactual outcome

$$\tilde{g}_0(x) = E[Y_0 | X = x]$$

To do so, we use the alternate aIPW signal

$$Y_0(0; m_0) = \frac{(1 - D)Y}{1 - o(Z)} + \frac{1 - D}{1 - o(Z)} m_0(Z)$$

where as before the true value for  $\tilde{o}(z) = \Pr(D = 1 | Z = z)$  but now  $m_0^2(z) = E[Y | D = 0; Z = z]$ . To estimate these nuisance models we again assume a logistic form for the propensity score model  $\tilde{o}(z) = \tilde{o}(z; \theta)$  and a linear form for the outcome regression model  $m_0(z) = m(z; \beta)$  as in (2.4) and use a separate estimation procedure for each basis term in our series approximation of  $\tilde{g}_0(x)$ . The estimating equations we

use to estimate each  $\hat{\beta}_j^0$  and  $\hat{\beta}_j^0$  differ from those in (2.5)-(2.6) however, and are instead given

$$\begin{aligned}\hat{\beta}_j^0 &:= \arg \min E_n[\rho_j(X)(1 - D)e^{\beta_j^0 Z} - D^0 Z g] + \lambda_j k_1 \\ \hat{\beta}_j^0 &:= \arg \min E_n[\rho_j(Z)(1 - D)e^{\beta_j^0 Z}(Y - D^0 Z)^2] + \lambda_j k_1\end{aligned}$$

which under the natural analogs of Assumptions 4.1 and 4.2 converge uniformly to population minimizers:

$$\begin{aligned}\beta_j^0 &:= \arg \min E[\rho_j(X)(1 - D)e^{\beta_j^0 Z} - D^0 Z g] \\ \beta_j^0 &:= \arg \min E[\rho_j(Z)(1 - D)e^{\beta_j^0 Z}(Y - D^0 Z)^2]\end{aligned}$$

Letting  $\bar{z}_{0,j}(z) = (z, \beta_j^0)$ , and  $\bar{m}_{0,j}(z) = m(z, \beta_j^0)$  we can repeat the decomposition of Section 3, expressing  $\tilde{Y}(\bar{z}_{0,j}; \bar{m}_{0,j})$  as functions of the parameters  $\beta_j^0$  and  $\beta_j^0$  and show that the first order conditions for  $\beta_j^0$  and  $\beta_j^0$  directly control the bias passed on to the second stage nonparametric estimator for  $\tilde{g}_0(x)$ . Convergence rates and validity of inference then follow from symmetric analysis of the results in Sections 4 and 5. Combining estimation and inference of the two conditional counterfactual outcomes then gives a doubly-robust estimator and model-assisted inference procedure for the CATE. To perform inference on the CATE we can use the variance matrix

$$\bar{\Omega} = \Omega_0 + \Omega_1 - 2\Omega_2$$

where  $\Omega_0$  is as in (2.13) but  $\Omega_1$  and  $\Omega_2$  are given

$$\begin{aligned}\Omega_1 &= Q^{-1} E[\tilde{f} \tilde{p}^k(x) \quad \tilde{g} \tilde{f} \tilde{p}^k(x) \quad \tilde{g} \tilde{g}^k] Q^{-1} \\ \Omega_2 &= Q^{-1} E[\tilde{f} \tilde{p}^k(x) \quad \tilde{g} \tilde{f} \tilde{p}^k(x) \quad \tilde{g} \tilde{g}^k] Q^{-1}\end{aligned}\tag{6.1}$$

where  $\tilde{g}_{0,j}^k = Y_0(\bar{z}_{0,j}; \bar{m}_{0,j}) - \tilde{g}_0(x)$  and  $\tilde{g}_0^k = (\tilde{g}_{0,1}^k; \dots; \tilde{g}_{0,k}^k)'$ . These matrices can be consistently estimated using their natural empirical analogs as in (5.10).

## 7 Simulation Study

The following simulations investigate the finite-sample estimation error and coverage properties, under misspecification, of the second stage estimator  $\hat{g}(x)$  described in Section 2.

### 7.1 Simulation Design

Random variables are generated i.i.d. according to the following distributions. The error term  $\epsilon$  is independently generated following  $N(0, \sigma^2)$  where  $\sigma = 1$ :  $Z = (X; \frac{1}{2}X^2; Z_1) \in \mathbb{R}^{2d_x + d_z}$  where  $d_x = 1$ ,  $d_z = 100$ ,  $X$  is independently generated following  $U(1, 2)$ , and the independent regressors  $Z_1$  are jointly centered Gaussian with a covariance matrix of the Toeplitz form

$$\text{Cov}(Z_{1,j}; Z_{1,k}) = E[Z_{1,j}Z_{1,k}] = 2^{-|j-k|}; \quad 1 \leq j, k \leq d_z:$$

These distributions and the following data-generating processes (DGP) follow the design of Tan (2020). Let  $Z^y$  be a transformation of the regressors in  $Z_1$  where  $Z_j^y = Z_j + \max(0, 1 - Z_j)^2$ ;  $j = 3; \dots; d_z$ . Let sparsity control the number of regressors in  $Z_1$ , i.e. excluding  $X$  terms and a constant term, entering the

DGP by  $\mathbf{Z} = (\underbrace{1, \dots, 1}_{\text{sparsity } d_z}, \underbrace{0, \dots, 0}_{\text{sparsity}})$ :

- (S1) *Correct specification*: Generate  $D$  given  $Z$  from a Bernoulli distribution with  $\Pr(D = 1|Z) = \frac{1}{2} + \exp(\rho_1 - X - 0.5X^2 - \theta Z)g^{-1}$  and  $Y = 1 + D + X + 0.5X^2 + \theta Z_1 + \epsilon$ .
- (S2) *Propensity score model correctly specified, but outcome regression model misspecified*: Generate  $D$  given  $Z$  as in (S1), but  $Y = 1 + D + X + 0.5X^2 + \theta Z_1^Y + \epsilon$ .
- (S3) *Propensity score model misspecified, but outcome regression model correctly specified*: Generate  $Y$  according to (S1), but generate  $D$  given  $Z$  from a Bernoulli distribution with  $\Pr(D = 1|Z) = \frac{1}{2} + \exp(\rho_2 - X - 0.5X^2 - \theta Z_1^Y)g^{-1}$ .

where the constants  $\rho_1$  and  $\rho_2$  differ in various simulation setups but are always set so that the average probability of treatment is about one half. To consider various degrees of high-dimensionality, we implement  $N \in \{500, 1,000\}$  with  $d_z = 100$ . For (S1),  $\text{sparsity} = 6$ ; for (S2),  $\text{sparsity} = 4$ ; and, for (S3),  $\text{sparsity} = 5$ . Results are reported for  $S = 1,000$  repeated simulations.

## 7.2 Estimators and Implementation

In our first-stage estimation, we need to set values for the arbitrary constants  $c_{j^*}, c_j$  for  $j \in \{1, \dots, k\}$  that appear in the pilot penalty parameters specified in (4.2). Following Tan (2020), we cross-validate to find suitable constants. These pilot penalty parameters are used to estimate  $\hat{b}_j^{\text{pilot}}$  and  $\hat{b}_j^{\text{pilot}}$  for  $j \in \{1, \dots, k\}$ . Following the multiplier bootstrap procedure outlined in (4.3) and (4.4) with  $B = 10,000$  bootstrap samples and  $c_0 = 1.1$ , the final penalty parameters  $\hat{c}_{j^*}$  and  $\hat{c}_j$  are obtained for use in the first-stage estimators (2.5) and (2.6). These  $2 \times k$  estimations for  $\hat{b}_j$  and  $\hat{b}_j$  are implemented in parallel using the `optim` function in R's `stats` package—which uses Nelder-Mead by default—with `maxit = 500` and `ndeps = 0.01` and the remaining settings using default values.

In our second-stage estimation, we follow the series estimation outlined in (2.12) to obtain  $\hat{b}^k$  which is used to estimate the target parameter  $g_0(x)$  with  $\hat{g}(x) := \hat{p}^k(x)^{\theta} \hat{b}^k$ . Results are reported for  $k = 3$ : For an approximating basis, B-splines are used from the `bSpline` function in the R package `splines2`, which uses the specification detailed in Perperoglou et al. (2019).

This Monte Carlo study aims to compare the proposed estimation and inference procedures to a state-of-the-art benchmark method targeting the CATE. Specifically, we use a benchmark method from Semenova and Chernozhukov (2021), which uses DML series estimation. This *DML* method follows Belloni et al. (2015) but uses the aIPW signal, which requires estimating the outcome regression in the first-stage. We follow the publically available code to use sample splitting in both first-stages; estimate the propensity score via logistic regression with an  $\ell_1$ -penalty; and, estimate the outcome regression via lasso.<sup>1</sup> Let *MA-DML* refer to our method, which aims to demonstrate, in relation to *DML*, the benefits of model assisted inference under misspecification.

To estimate the variance of  $\hat{g}(x)$  and form confidence intervals, we follow (2.14) to obtain the estimated covariance matrix  $\hat{\Omega}$ . Confidence bands are formed using the estimated standard error  $\hat{b}(x) := k \hat{\Omega}^{1/2} \hat{p}^k(x) k = \frac{P}{n}$  according to (2.17), which uses different bootstrapped critical values for pointwise and uniform confidence intervals.

<sup>1</sup>Vira Semenova provides several example R scripts implementing *DML*: <https://sites.google.com/view/semenovavira/research>.

Estimation error is studied for each method's  $\hat{g}(x)$  for the target parameter  $g_0(x) = \mathbb{E}[Y|D = 1; X = x]$  over a grid of 100 points spaced across  $x \in [1; 2]$ , i.e. the support of  $X$ . We study coverage across  $S$  simulations of each method's pointwise (i.e.  $b(x)$  at  $x = 1.5$ ) and uniform confidence intervals.

To compare the estimation error for the target parameter  $g(x)$  across the two different estimators  $\hat{g}_s(x)$  for each simulation  $s = 1; \dots; S$ , we utilize integrated bias, variance, and mean-squared error where  $\bar{g}(x) = S^{-1} \sum_{s=1}^S \hat{g}_s(x)$ :

$$\begin{aligned} \text{IBias}^2 &= \int_0^1 (\bar{g}(x) - g_0(x))^2 dx; \\ \text{IVar} &= S^{-1} \sum_{s=1}^S \int_0^1 (\hat{g}_s(x) - \bar{g}(x))^2 dx; \\ \text{IMSE} &= S^{-1} \sum_{s=1}^S \int_0^1 (\hat{g}_s(x) - g_0(x))^2 dx; \end{aligned}$$

### 7.3 Simulation Results

We aim to study the estimation error and coverage properties of this paper's proposed method under correctly and misspecified regimes according to S1-S3 as compared to a state-of-the-art benchmark method. Specifically, do the model-assisted inference methods improve the estimator's coverage properties under misspecification? How does estimation error change under misspecification? How does performance respond to different degrees of high-dimensionality?

Table 7.1 presents the simulation results. We present all three specifications, i.e. (S1)-(S3), in two different regimes, i.e.  $N = 500$  in the top panel and  $N = 1000$  in the lower panel, for both estimators *DML* and *MA-DML*. Estimation errors measures integrated squared bias, variance, and mean squared error are presented in columns (1)-(3). Pointwise and uniform coverage results are presented in columns (4)-(7).

First, for pointwise and uniform coverage under correct specification regime (S1), *MA-DML* has some slight improvements. Under misspecification DGPs (S2) and (S3), the pointwise coverage of *MA-DML* is closer to the targets except in the  $N = 1000$  and (S2) case where it slightly underperforms. However, *MA-DML* has a notable improvement over *DML* in the (S3) case when  $N = 1000$ : Similarly, *MA-DML* outperforms *DML* in three of the four misspecified regimes, i.e. all but (S3) when  $N = 500$  where *MA-DML* has over-coverage. Under (S2) when  $N = 1000$ ; both methods are markedly deteriorated uniform coverage, although *MA-DML* is noticeably closer to target.

Next, in regards to estimation error, in four of the six settings, *MA-DML* has a lower MSE than *DML* where regardless of sample size *MA-DML* underperforms in (S3). Notably, it does appear *MA-DML* has substantially smaller  $\text{IBias}^2$  across the DGPs.

Finally, we were surprised to find for both estimators that coverage properties, in general, improve under the higher-dimensional regime of  $N = 500$  with  $d_z = 100$  compared to  $N = 1000$  and  $d_z = 100$ : In particular, with a higher ratio of covariates to observations, the uniform coverage properties under regime (S2) were substantially better. The estimation error results were in line with our priors as the higher-dimensional regime sees in general higher estimation errors for both methods.

Table 7.1: Simulation study.

DGP	Estimator	IBias <sup>2</sup> (1)	IVar (2)	IMSE (3)	Cov90 (4)	Cov95 (5)	UCov90 (6)	UCov95 (7)
K=3, N=500, $d_Z = 100$								
(S1)	DML	0.04	0.31	0.35	0.92	0.96	1.00	1.00
	MA-DML	0.0	0.34	0.34	0.93	0.97	1.00	1.00
(S2)	DML	0.16	2.17	2.33	0.92	0.97	0.83	0.86
	MA-DML	0.03	2.12	2.15	0.90	0.94	0.88	0.91
(S3)	DML	0.03	0.55	0.59	0.87	0.93	0.95	0.97
	MA-DML	0.01	0.79	0.80	0.91	0.95	0.99	0.99
K=3, N=1000, $d_Z = 100$								
(S1)	DML	0.12	0.20	0.32	0.83	0.90	0.96	0.96
	MA-DML	0.01	0.22	0.23	0.83	0.90	0.99	0.99
(S2)	DML	0.40	2.1	2.5	0.84	0.91	0.33	0.39
	MA-DML	0.19	2.07	2.26	0.83	0.89	0.50	0.55
(S3)	DML	0.11	0.34	0.46	0.74	0.82	0.80	0.84
	MA-DML	0.01	0.53	0.54	0.84	0.89	0.89	0.91

Note: DGP refers to the three various data generating processes introduced above. IBias<sup>2</sup>, IVar, and IMSE refer to integrated squared bias, variance, and mean squared error, respectively. Cov90, Cov95, UCov90, and UCov95 refer to the coverage proportion of the 90% and 95% pointwise and uniform confidence intervals across simulations.  $K$  refers to the number of series terms,  $N$  to the sample size, and  $d_Z$  to the dimensionality of the random variable  $Z_1$ :

For coverage under correct specification, we did anticipate the underperformance of *MA-DML* given it is designed to handle misspecification with the cost of other estimators outperforming under correct specification. Additionally, We attribute the poor uniform coverage in DGP (S2) for both estimators under  $N = 1,000$  to a lack of a rich enough cross-validation given the performance was improved under a more difficult regime when the number of observations drops to  $N = 500$ :

Our findings should not be interpreted as a critique of the authors of this benchmark method, whose work we rely on and were inspired by.

## 8 Empirical Application

We apply the model assisted estimator to estimate the effect of maternal smoking on infant birthweight conditional on the age of the mother. We use the Cattaneo (2010) dataset which can be found online on the Stata website.<sup>1</sup> The dataset describes each infant's birthweight in grams,  $Y$ , whether or not the mother smoked during pregnancy,  $D = 1$  indicating smoking, and a number of covariates containing information on the mother's health and socioeconomic background,  $Z = (X; Z_1)$ , where  $X$  represents the conditioning variable, maternal age. A full summary of the data used as well as additional details/analysis from our empirical analysis can be found in Appendix F.

We compare the model assisted estimator of the CATE against one where standard MLE and OLS loss functions are used to estimate the first stage propensity score and outcome regression models. We also qualitatively compare our results to Zimmert and Lechner (2019), who use a kernel based approach to

<sup>1</sup>The dataset can be downloaded [here](#).



estimate the CATE in this setting. While this sort of comparison is not perfect since we do not know the true DGP, this setting is advantageous for analysis since we strongly expect that (i) the effect of smoking on birthweight will be negative and (ii) this effect should grow stronger in magnitude as the age of the mother increases. These hypotheses have been corroborated by other work that examines the conditional average treatment effect in this setting (Zimmert and Lechner, 2019; Abrevaya, 2006; Lee et al., 2017).

## 8.1 Empirical Results

Figure 8.1 displays our main results from implementing both the model assisted and standard MLE/OLS estimation procedures. After removing the top 3% and bottom 3% of smoker and non-smoker birthweights by maternal age, we select the penalty parameters for the first stage models via the bootstrap procedure described in Section 4. The pilot penalty parameters are uniformly taken to be equal to zero, so that the residuals used in the bootstrap procedure are generated from non-regularized estimations. We take  $c_0 = 2$  in (4.4) and select the first stage penalty parameters using the 99<sup>th</sup>, 95<sup>th</sup>, and 90<sup>th</sup> quantiles of the bootstrap distribution. For the second stage basis functions we implement second degree b-splines with 3 knots via the splines2 package in R (Wang and Yan, 2021).

Consistent with prior work, both estimators of the CATE suggest that the effect of smoking on birthweight becomes more negative with age. Both estimation procedures also generally produces negative estimates for the CATE, but it should be noted that for the lowest levels of penalization the model assisted CATE estimate suggests a slightly positive effect of smoking for particularly young mothers, though this difference is not significantly different from zero. The shapes of the estimated functions remain relatively stable under various sizes of the penalty parameter, though the model assisted procedure displays a bit more sensitivity to the level of regularization introduced.<sup>1</sup>

For the most part, the effects found here are similar to those found in Zimmert and Lechner (2019), though the effects estimated using standard first stage loss functions have somewhat larger magnitudes and in general both series estimation procedures seem to give less reasonable results on the boundaries. An advantage of using a series second stage however, compared to the kernel first stage of Zimmert and Lechner (2019), is the existence of the uniform confidence bands displayed. Reassuringly, the estimates of Zimmert and Lechner (2019) seem to be within the 95% uniform confidence bands generated by the model assisted estimator.

As a robustness check, we also try estimating the treatment effect using first degree b-splines instead of second degree splines. These results are displayed in Figure 8.2. Again, we find that the effect of smoking on child birthweight is almost uniformly negative regardless of estimation procedure used or choice of penalty parameter. The shape of the estimated CATE function using a standard MLE/OLS first stage is very stable to penalty choice here while the shape of the model assisted CATE function displays a bit more instability here at the two lower levels of regularization.

Finally, Table 8.1 reports the smoothed average treatment effect estimates taken from averaging the model assisted CATE estimates from Figure 8.1 across observations. Again, these estimates are generally in line with prior work

---

<sup>1</sup>Numerically solving the minimization problems in (2.5)-(2.6) also typically requires more iterations to converge than solving the standard MLE/OLS minimization problems.

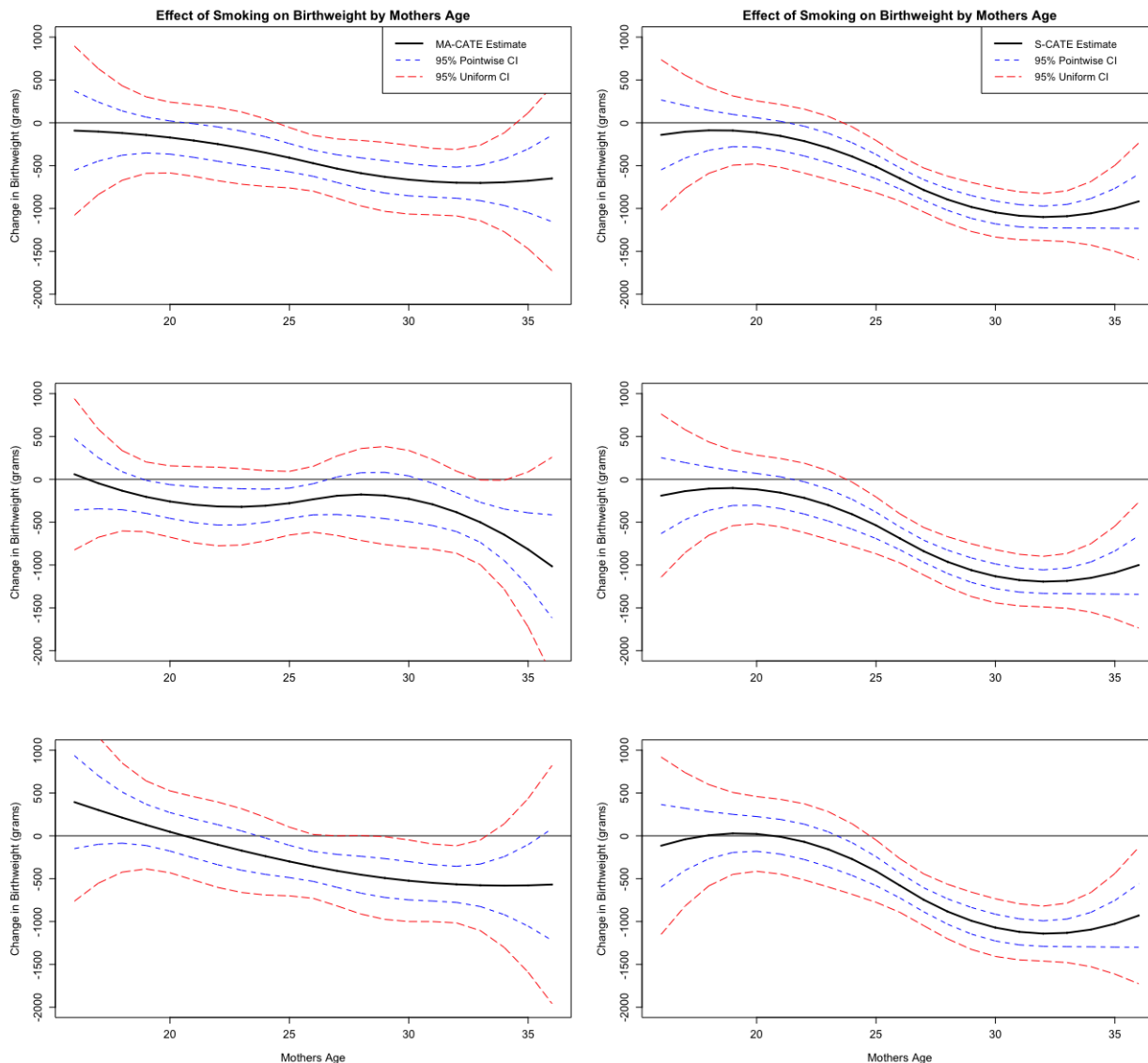


Figure 8.1: CATE of maternal smoking estimated using model assisted estimating equations (left) and standard MLE/OLS estimating equations (right). Top row uses the 99<sup>th</sup> quantile of the bootstrap distribution to select the penalty parameters, second row uses 95<sup>th</sup> quantile, and final row uses the 90<sup>th</sup> quantile. Second stage is computed using b-splines of the second degree with 3 knots. 95% pointwise confidence intervals are displayed in blue short dashes and 95% uniform confidence bands are displayed in long red dashes.

Bootstrap Penalty Qt.	99 <sup>th</sup>	95 <sup>th</sup>	90 <sup>th</sup>
Implied ATE	-295.221	-292.9086	-453.2242

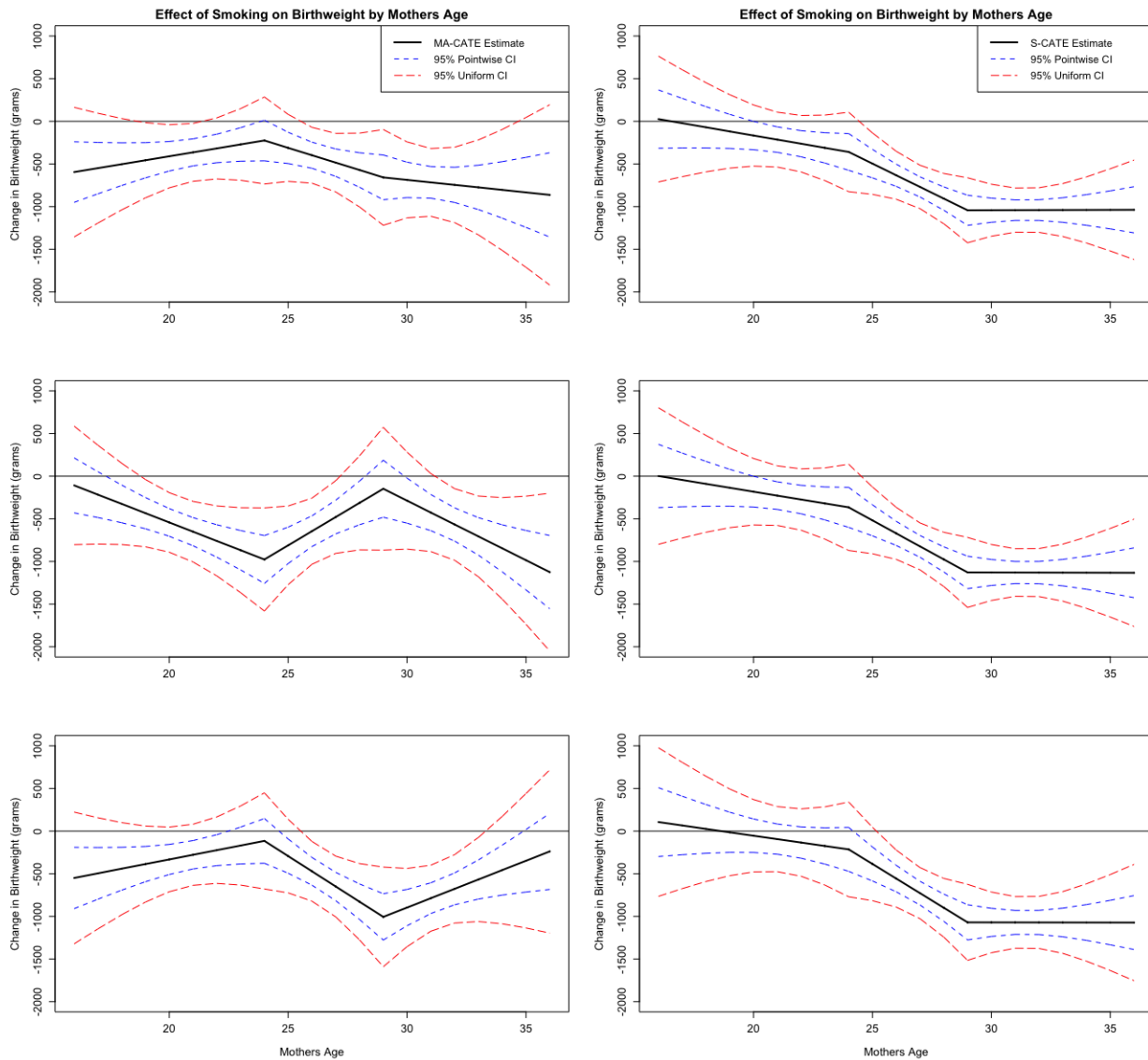


Figure 8.2: CATE of maternal smoking estimated using model assisted estimating equations (left) and standard MLE/OLS estimating equations (right). Top row uses the 99<sup>th</sup> quantile of the bootstrap distribution to select the penalty parameters, second row uses 95<sup>th</sup> quantile, and final row uses the 90<sup>th</sup> quantile. Second stage is computed using b-splines of the first degree with 3 knots. 95% pointwise confidence intervals are displayed in blue short dashes and 95% uniform confidence bands are displayed in long red dashes.

## 9 Conclusion

Estimation of conditional average treatment effects with high dimensional controls typically relies on first estimating two nuisance parameters: a propensity score model and an outcome regression model. In a high-dimensional setting, consistency of the nuisance parameter estimators typically relies on correctly specifying their functional forms. While the resulting second-stage estimator for the conditional average treatment effect typically remains consistent even if one of the nuisance parameters is inconsistent, the confidence intervals may no longer be valid.

In this paper, we consider estimation and valid inference on the conditional average treatment effect in the presence of high dimensional controls and nuisance parameter misspecification. We present a nonparametric estimator for the CATE that remains consistent at the nonparametric rate, under slightly modified conditions, even under misspecification of either the logistic propensity score model or linear outcome regression model. The resulting Wald-type confidence intervals based on this estimator also provide valid asymptotic coverage under nuisance parameter misspecification.

## References

- Abrevaya, J. (2006). Estimating the effect of smoking on birth outcomes using a matched panel data approach. *Journal of Applied Econometrics* 21(4), 489–519.
- Bauer, B. and M. Kohler (2019). On deep learning as a remedy for the curse of dimensionality in nonparametric regression. *The Annals of Statistics* 47(4), 2261 – 2285.
- Belloni, A., D. Chen, V. Chernozhukov, and C. Hansen (2012). Sparse models and methods for optimal instruments with an application to eminent domain. *Econometrica* 80(6), 2369–2429.
- Belloni, A. and V. Chernozhukov (2013). Least squares after model selection in high-dimensional sparse models. *Bernoulli* 19(2), 521 – 547.
- Belloni, A., V. Chernozhukov, D. Chetverikov, C. Hansen, and K. Kato (2018). High-dimensional econometrics and regularized gmm.
- Belloni, A., V. Chernozhukov, D. Chetverikov, and K. Kato (2015). Some new asymptotic theory for least squares series: Pointwise and uniform results. *Journal of Econometrics* 186(2), 345–366. High Dimensional Problems in Econometrics.
- Belloni, A., V. Chernozhukov, and C. Hansen (2013, 11). Inference on treatment effects after selection among high-dimensional controls. *The Review of Economic Studies* 81(2), 608–650.
- Bickel, P., C. Klaassen, Y. Ritov, and J. Wellner (1993). *Efficient and Adaptive Estimation for Semiparametric Models*. Johns Hopkins series in the mathematical sciences. Johns Hopkins University Press.
- Bickel, P. J., Y. Ritov, and A. B. Tsybakov (2009). Simultaneous analysis of Lasso and Dantzig selector. *The Annals of Statistics* 37(4), 1705 – 1732.
- Bühlmann, P. and S. van de Geer (2011). *Statistics for high-dimensional data*. Springer Series in Statistics. Springer, Heidelberg. Methods, theory and applications.
- Cattaneo, M. (2010). Efficient semiparametric estimation of multi-valued treatment effects under ignorability. *Journal of Econometrics* 155(2), 138–154.
- Chen, X. (2007). Large sample sieve estimation of semi-nonparametric models. In J. Heckman and E. Leamer (Eds.), *Handbook of Econometrics* (1 ed.), Volume 6B, Chapter 76, pp. 5549–5632. Elsevier.
- Chernozhukov, V., D. Chetverikov, M. Demirer, E. Duflo, C. Hansen, W. Newey, and J. Robins (2018, 01). Double/debiased machine learning for treatment and structural parameters. *The Econometrics Journal* 21(1), C1–C68.
- Chernozhukov, V., D. Chetverikov, and K. Kato (2017). Central limit theorems and bootstrap in high dimensions. *The Annals of Probability* 45(4), 2309–2352.
- Chetverikov, D., Z. Liao, and V. Chernozhukov (2021). On cross-validated Lasso in high dimensions. *The Annals of Statistics* 49(3), 1300 – 1317.
- Chetverikov, D. and J. R.-V. Sørensen (2021). Analytic and bootstrap-after-cross-validation methods for selecting penalty parameters of high-dimensional m-estimators. *ArXiv NA*, 1–50.

- De Boor, C. (2001). *A practical guide to splines; rev. ed.* Applied mathematical sciences. Berlin: Springer.
- der Vaart, A. V. and J. Wellner (1996). *Weak Convergence and Empirical Processes* (1 ed.). Springer Series in Statistics. Springer, New York, NY.
- Dudley, R. (1967). The sizes of compact subsets of hilbert space and continuity of gaussian processes. *Journal of Functional Analysis* 1(3), 290–330.
- Fan, Q., Y.-C. Hsu, R. P. Lieli, and Y. Zhang (2022). Estimation of conditional average treatment effects with high-dimensional data. *Journal of Business & Economic Statistics* 40(1), 313–327.
- Giné, E. and V. Koltchinskii (2006). Concentration inequalities and asymptotic results for ratio type empirical processes. *The Annals of Probability* 34(3), 1143 – 1216.
- Hlavac, M. (2022). *stargazer: Well-Formatted Regression and Summary Statistics Tables*. Bratislava, Slovakia: Social Policy Institute. R package version 5.2.3.
- Lee, S., R. Okui, and Y.-J. Whang (2017). Doubly robust uniform confidence band for the conditional average treatment effect function. *Journal of Applied Econometrics* 32(7), 1207–1225.
- Newey, W. (1997). Convergence rates and asymptotic normality for series estimators. *Journal of Econometrics* 79(1), 147–168.
- Newey, W. K. and D. McFadden (1994). Chapter 36 large sample estimation and hypothesis testing. *Handbook of Econometrics* 4, 2111–2245.
- Perperoglou, A., W. Sauerbrei, M. Abrahamowicz, and M. Schmid (2019). A review of spline function procedures in r. *BMC medical research methodology* 19(1), 1–16.
- Pollard, D. (2001). *A User’s Guide to Measure Theoretic Probability*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press.
- Rubin, D. B. (1974). Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of Educational Psychology* 66, 688–701.
- Rubin, D. B. (1978). Bayesian inference for causal effects. *The Annals of Statistics* 6(1), 34–58.
- Rudelson, M. (1999). Random vectors in the isotropic position. *J. Funct. Anal* 164, 60–72.
- Schmidt-Hieber, J. (2020, 08). Nonparametric regression using deep neural networks with relu activation function. *Annals of Statistics* 48, 1875–1897.
- Semenova, V. and V. Chernozhukov (2021, 08). Debiased machine learning of conditional average treatment effects and other causal functions. *The Econometrics Journal* 24, 264–289. utaa027.
- Smucler, E., A. Rotnitzky, and J. M. Robins (2019). A unifying approach for doubly-robust  $\lambda_1$  regularized estimation of causal contrasts. *ArXiv NA*, 1–125.
- Stone, C. J. (1985). Additive Regression and Other Nonparametric Models. *The Annals of Statistics* 13(2), 689 – 705.
- Tan, Z. (2017). Regularized calibrated estimation of propensity scores with model misspecification and high-dimensional data. *ArXiv NA*, 1–60.

- Tan, Z. (2020). Model-assisted inference for treatment effects using regularized calibrated estimation with high-dimensional data. *The Annals of Statistics* 48(2), 811 – 837.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)* 58(1), 267–288.
- van der Greer, S. (2016). *Estimation and Testing under Sparsity*. Lecture Notes in Mathematics. Springer, New York, NY.
- Wang, W. and J. Yan (2021). Shape-restricted regression splines with R package splines2. *Journal of Data Science* 19(3), 498–517.
- Wu, P., Z. Tan, W. Hu, and X.-H. Zhou (2021). Model-assisted inference for covariate-specific treatment effects with high-dimensional data.
- Zimmert, M. and M. Lechner (2019). Nonparametric estimation of causal heterogeneity under high-dimensional confounding.

## A Proofs for Results in Main Text

Here we provide proofs of the main results in Sections 3-5. The proofs for Section 4 rely on an assortment of supporting lemmas proved in Appendix B.

### A.1 Proof of Proposition 3.1

Immediately follows from Theorems 4.3, 4.4, 5.1 and 5.4 under the sparsity bounds in (4.10) and (4.13).

### A.2 Proof of Proposition 3.2

Immediately follows from Theorems 4.3, 4.4, 5.4 and 5.5 under the sparsity bounds in (4.10) and (4.13).

### A.3 Proofs of Results in Section 4

#### Proof of Theorem 4.1

Lemma B.7 shows that there is a constant  $M$  that does not depend on  $k$  such that under the conditions of Assumption 4.1 and Theorem 4.1,  $\Pr(\Omega_{k;1} \setminus \Omega_{k;2}) \leq o(1)$  where

$$\Pr(\Omega_{k;1} \setminus \Omega_{k;2}) \leq M \sqrt{\frac{\log(d_{z^*})}{n}}.$$

Lemma B.1 provides the bound under this event.

#### Proof of Theorem 4.2

Lemmas B.7-B.10 show that there is a constant  $M$  that does not depend on  $k$  such that under the conditions of Assumptions 4.1 and 4.2 and Theorem 4.2,  $\Pr(\bigcap_{m=1}^6 \Omega_{k;m}) \leq o(1)$  where

$$\Pr(\bigcap_{m=1}^6 \Omega_{k;m}) \leq M \sqrt{\frac{\log(d_{z^*})}{n}}.$$

Lemma B.2 provides the bound under this event.

#### Proof of Theorem 4.3

Lemmas B.7-B.10 show that there is a constant  $M$  that does not depend on  $k$  such that under the conditions of Assumptions 4.1 and 4.2 and Theorem 4.3,  $\Pr(\bigcap_{m=1}^6 \Omega_{k;m}) \leq o(1)$  where

$$\Pr(\bigcap_{m=1}^6 \Omega_{k;m}) \leq M \sqrt{\frac{\log(d_{z^*})}{n}}.$$

Lemma B.3 provides the bound under this event.

#### Proof of Theorem 4.4

Lemmas B.7-B.11 show that there is a constant  $M$  that does not depend on  $k$  such that under the conditions of Assumptions 4.1 and 4.2 and Theorem 4.4,  $\Pr(\bigcap_{m=1}^7 \Omega_{k;m}) \leq o(1)$  where

$$\Pr(\bigcap_{m=1}^7 \Omega_{k;m}) \leq M \sqrt{\frac{\log(d_{z^*})}{n}}.$$



Lemma B.4 provides the bound under this event.

### A.4 Proofs of Results in Section 5

The proofs for Section 5 closely follow those of Belloni et al. (2015) with some modifications to deal with the various error terms. They also rely on some additional second stage results proved in Appendix C.

#### Proof of Theorem 5.1

Equation (3.2) follows from applying (3.1) with  $\rho(x) = k\rho(x)k$  and (3.3) follows from (3.2). So it suffices to prove (3.1).

For any  $\rho \geq S^{k-1}$ ,  $1 \leq k \leq \Omega^{1-2k}$  because of the conditional variance of  $\rho_j^2$  is bounded from below and from above and under the positive semidefinite ranking

$$\Omega \leq \Omega_0 \leq \rho^2 Q^{-1};$$

Moreover, by condition (ii) of the theorem and Lemma C.2,  $R_{1n}(\rho) = o_p(1)$ . So we can write

$$\begin{aligned} \rho_{\bar{n}}^{-1}(\rho) &= \frac{\rho_{\bar{n}}^{-1}}{k \Omega^{1-2k}} G_n[\rho^k(x) \cdot (k + r_k)] + o_p(1) \\ &= \sum_{i=1}^k \frac{\rho_{\bar{n}}^{-1}}{\bar{n} k \Omega^{1-2k}} f \rho^k(x) \cdot (k + r_k) g_i \end{aligned}$$

Goal will be to verify Lindberg's condition for the CLT. Throughout the rest of the proof, it will be helpful to make the following notations. First, for any vector  $a = (a_1, \dots, a_k) \geq S^{k-1}$ , let  $|a| = (|a_1|, \dots, |a_k|)$  and note that  $|a| \geq S^{k-1}$  as well:

$$\rho_n^{-1} = \frac{\rho_{\bar{n}}^{-1}}{\bar{n} k \Omega^{1-2k}}; \quad \rho_n := \sum_{j=1}^k \rho^k(x); \quad \text{and} \quad \bar{\rho}_k := \sup_{j=1, \dots, k} |a_j|$$

Now, by the definition of  $\Omega$  we have that

$$\text{Var} \left[ \sum_{i=1}^k \frac{\rho_{\bar{n}}^{-1}}{\bar{n} k \Omega^{1-2k}} f \rho^k(x) \cdot (k + r_k) g_i \right] = 1;$$

Second for each  $\epsilon > 0$

$$\begin{aligned} & \sum_{i=1}^k \mathbb{E} \left[ \left( \frac{\rho_{\bar{n}}^{-1}}{\bar{n} k \Omega^{1-2k}} f \rho^k(x) \cdot (k + r_k) g_i \right)^2 \mathbf{1}_{\left| \sum_{j=1}^k \rho^k(x) \cdot (k + r_k) g_j \right| > \epsilon} \right] \\ & \leq \sum_{i=1}^k \mathbb{E} \left[ \frac{\rho_{\bar{n}}^{-1}}{\bar{n} k \Omega^{1-2k}} \mathbf{1}_{\left| \sum_{j=1}^k \rho^k(x) \cdot (k + r_k) g_j \right| > \epsilon} \right] \leq \epsilon \end{aligned} \tag{A.1}$$

What we are using here is the following. Suppose  $\mu$  is a nonrandom vector in  $\mathbb{R}^k$ ,  $a$  is a (positive) random vector in  $\mathbb{R}^k$  and  $b$  is a random vector in  $\mathbb{R}^k$ . Then,

$$f^\theta(a, b)g = \sum_{j=1}^k a_j b_j - \mu^T \sum_{j=1}^k a_j = (j^\theta a) \mu^T : \quad (\text{A.2})$$

To bound the right hand side of (A.1) use the fact that  $\| \cdot \|_{\Omega} \leq k \|\cdot\|_{\Omega^{-1}}$  because  $\| \cdot \|_{\Omega} \leq \|\cdot\|_{\Omega^{-1}}$  and

$$\Omega \preceq \Omega_0 - \epsilon^2 Q^{-1}$$

in the positive semidefinite sense. Using these two we have

$$n E \|j^\theta p^k(x)\|^2 = \text{trace}(\Omega^{-1} \Omega) \leq k:$$

By the bounded eigenvalue condition and using the trace operator:

$$E[(j^\theta p^k(x))^2] = \text{trace}(E[j^\theta p^k(x)^\theta p^k(x)^\theta]) = \sum_{k=1}^k \lambda_k$$

Further note,  $\|j^\theta p^k(x)\| \leq \frac{c}{n}$ . Using  $(a+b)^2 \leq 2a^2 + 2b^2$ , the right hand side of (A.1) is bounded by

$$2n E \|j^\theta p^k(x)\|^2 + \frac{c^2}{n} \leq 2n E \|j^\theta p^k(x)\|^2 + \frac{c^2}{n} \leq 2n E \|j^\theta p^k(x)\|^2 + \frac{c^2}{n}$$

and both terms converge to zero. Indeed, to bound the first term note that, for some  $c > 0$ :

$$2n E \|j^\theta p^k(x)\|^2 + \frac{c^2}{n} \leq 2n E \|j^\theta p^k(x)\|^2 + \frac{c^2}{n} \leq 2n E \|j^\theta p^k(x)\|^2 + \frac{c^2}{n} = o(1)$$

where here we use the first part of Assumption 5.2. To show the second term converges to zero, follow the same steps as for the first term, but apply the second part of Assumption 5.2.

**Proof of Theorem 5.2**

The goal will be to apply the following two theorems from [Giné and Koltchinskii \(2006\)](#) and [der Vaart and Wellner \(1996\)](#).

Preliminaries for Proof of Theorem 5.2

**Theorem** (Gine and Koltchinskii, 2006). Let  $\{f_i\}_{i=1}^n$  be i.i.d random variables taking values in a measurable space  $(S; \mathcal{S})$  with a common distribution  $P$  defined on the underlying  $n$ -fold product space. Let  $F$  be a measurable class of functions mapping  $S \rightarrow \mathbb{R}$  with a measurable envelope  $F$ . Let  $V$  be a constant such that  $\sup_{f \in F} \text{Var}(f) \leq V k_{L^2(P)}^2$ . Suppose there exist constants  $A > e^2$  and  $V > 2$  such that  $\sup_Q N(F; L^2(Q); \|k_{L^2(Q)}\|) \leq (A\|k_{L^2(Q)}\|)^V$  for all  $0 < \|k_{L^2(Q)}\| \leq 1$ . Then

$$E \sum_{i=1}^n (f_i - E[f_i])^2 \leq C \frac{V}{n} \log \frac{A k_{L^2(P)}}{V} + V k_{L^2(P)} \log \frac{A k_{L^2(P)}}{V} \quad (\text{GK})$$

where  $C$  is a universal constant.

**Theorem** (VdV&W 2.14.1). Let  $F$  be a  $P$ -measurable class of measurable functions with a measurable envelope function  $F$ . Then for any  $p \geq 1$ ,

$$k_{G_n, F, P, p} \leq k_J(n; F) k_{F, k_{P, p}} \leq J(1; F) k_{F, k_{P, 2-p}} \quad (\text{VW})$$

where  $k_{G_n, F, P, p} = k_{F, k_{P, p}}$ , where  $k_{F, k_{P, p}}$  is the  $L_2(P_n)$  seminorm and the inequalities are valid up to constants depending only on the  $p$  in the statement. The term  $J(\cdot; \cdot)$  is given

$$J(\cdot; F) = \sup_Q \int_0^1 \frac{1}{1 + \log N(F; k_{L^2(Q)}; \|k_{L^2(Q)}\|)} d\|k_{L^2(Q)}\|$$

We would like to apply these theorems to bound  $\sup_{x \in X} j(x) G_n[p^k(x) - k]$  and thus show (5.4). The other two statements of Theorem 5.2 follow from this. To this end, let's consider the class of functions

$$G := f(\cdot; x) - (v)^0(p^k(x) - k); v \in X$$

Let's note that  $j(v)^0(p^k(x) - k)$ ,  $\text{Var}((v)^0(p^k(x) - k)) = 1$ , and for any  $v, \tilde{v} \in X$

$$j(v)^0(p^k(x) - k) - j(\tilde{v})^0(p^k(x) - k) \leq |j(v) - j(\tilde{v})| \frac{1}{k} \leq \tilde{v} k$$

where  $\tilde{v} k = k^k k_1$ . Then, taking  $G(\cdot; x) = \tilde{v} k$  we have that

$$\sup_Q N(G; L^2(Q); \|k_{L^2(Q)}\|) \leq \frac{C}{\tilde{v} k} \frac{1}{k} \quad (\text{A.3})$$

Now, for a  $\delta > 0$  specified later define  $k^- = \mathbf{1}_{f_j^- > g} E[\mathbf{1}_{f_j^- > g} j^- X]$  and  $k^+ = \mathbf{1}_{f_j^- > g} E[\mathbf{1}_{f_j^- > g} j^- X]$ . Since  $E[j^- X] = 0$  we have that  $k^- = k^- + k^+$ . Using this decompose:

$$\frac{1}{n} \sum_{i=1}^n (v)^0(p^k(x) - k) = \frac{1}{n} \sum_{i=1}^n (v)^0(p^k(x) - k^-) + \frac{1}{n} \sum_{i=1}^n (v)^0(p^k(x) - k^+)$$

We deal with each of these terms individually, in two steps.

**Step 1:** For the first term, we set up for an application of (GK). Equation (A.3) gives us the constants

$A = C \frac{L}{k}$  and  $V = d_{X-2}$ . To get  $\frac{2}{k}$  note that for any  $v \geq X$ ,

$$\begin{aligned} \text{Var}((v)^\theta(p^k(x) - \rho_{\bar{n}})) &= E[(v)^\theta(p^k(x) - \rho_{\bar{n}})^2] \\ &= \frac{1}{n} E[(v)^\theta(p^k(x))^2] \sup_{x \in X} E[k_k k_1^2 | \mathcal{X} = x] \\ &= \frac{\frac{2}{k} \wedge 2}{n} \end{aligned}$$

Finally note that we can take the envelope  $G = k_k k_1 \rho_{\bar{n}}$  where  $k_k k_1 \rho_{\bar{n}} \in L^2(P)$  and  $k_k k_1 \rho_{\bar{n}} \leq \frac{L}{n}$ .

We can now apply (GK) to get that

$$E[\sup_{x \in X} j(x)^\theta G_n[p^k(x) - \rho_{\bar{n}}]] \leq \frac{1}{k} \frac{2 \log(\frac{L}{k})}{\rho_{\bar{n}}} + \frac{k \log(\frac{L}{k})}{\rho_{\bar{n}}}$$

**Step 2:** For the second term, we set up for an application of (VW) with the envelope function  $G = k_k k_1 \rho_{\bar{n}}$  and note that

$$E[k_k k_1^2] = E[\frac{2}{k} \mathbf{1}_{\{j-k_j > g\}}] \leq \frac{m+2}{k} E[j-k_j^m]$$

We can now use (VW) to bound

$$\begin{aligned} E[\sup_{x \in X} j(x)^\theta G_n[p^k(x) - \rho_{\bar{n}}]] &\leq \frac{1}{k} \frac{2 \log(\frac{L}{k})}{\rho_{\bar{n}}} + \frac{k \log(\frac{L}{k})}{\rho_{\bar{n}}} \\ &\leq \frac{1}{k} \frac{2 \log(\frac{L}{k})}{\rho_{\bar{n}}} + \frac{k \log(\frac{L}{k})}{\rho_{\bar{n}}} \end{aligned}$$

**Step 3:** Let  $m = \frac{2(m-2)}{k}$  and apply Markov's inequality. The bounds from step one and two become

$$\begin{aligned} \sup_{x \in X} j(x)^\theta G_n[p^k(x) - \rho_{\bar{n}}] &\leq \frac{1}{k} \frac{2 \log(\frac{L}{k})}{\rho_{\bar{n}}} + \frac{k \log(\frac{L}{k})}{\rho_{\bar{n}}} \\ \sup_{x \in X} j(x)^\theta G_n[p^k(x) - \rho_{\bar{n}}] &\leq \frac{1}{k} \frac{2 \log(\frac{L}{k})}{\rho_{\bar{n}}} \end{aligned}$$

Applying Assumption 5.3 along with the inequality

$$\frac{k \log k}{\rho_{\bar{n}}} \leq \frac{1}{\log k} \frac{2m(m-2) \log k}{n} \leq \log k$$

completes the proof.

**Proof of Theorem 5.3**

We apply Yurinskii's coupling lemma (Pollard, 2001)

**Yurinskii's Coupling Lemma**

Let  $x_1, \dots, x_n$  be independent random  $k$ -vectors with  $E[x_i] = 0$  and  $\sigma := \sum_{i=1}^n E[|x_i|^3]$  finite. Let  $S := x_1 + \dots + x_n$ . For each  $\epsilon > 0$  there exists a random vector  $T$  with a  $N(0; \text{var}(S))$  distribution such that

$$P(|S - T| > \epsilon) \leq C_0 B \left(1 + \frac{j \log(1+B)j}{k}\right) \quad \text{where } B := k^{-3} \quad \text{(YC)}$$

for some universal constant  $C_0$ .

In order to apply the coupling, we want to consider a first order approximation to the estimator

$$\frac{1}{n} \sum_{i=1}^n x_i \quad x_i = \Omega^{-1/2} p^k(x) \quad (k + r_k):$$

When  $\bar{R}_{2n} = o_p(a_n^{-1})$  a similar argument can be used with  $x_i = \Omega^{-1/2} p^k(x) \quad (k + r_k)$  replaced with  $\Omega^{-1/2} p^k(x) \quad k$ . As before, the eigenvalues of  $\Omega$  are bounded away from zero, therefore

$$\begin{aligned} E|k|^3 &\leq E[k p^k(x) \quad (k(x) + r_k) k^3] \\ &\leq E[k p^k(x) k^3 (j^{-3} + j r_k^3)] \\ &\leq E[k p^k(x) k^3] (k^{-3} + k^3 c_k^3) \\ &\leq E[k p^k(x) k^3] k (k^{-3} + k^3 c_k^3) \\ &\leq k k (k^{-3} + k^3 c_k^3) \end{aligned}$$

Therefore, by Yurinskii's coupling lemma (YC), for each  $\epsilon > 0$ ,

$$\begin{aligned} \Pr \left\{ \left| \frac{1}{n} \sum_{i=1}^n x_i - \Omega^{-1/2} N_k k \right| > \epsilon a_n^{-1} \right\} &\leq \frac{n k^2 k (k^{-3} + k^3 c_k^3)}{(a_n^{-1} n)^3} \left( 1 + \frac{\log(k^3 k (k^{-3} + k^3 c_k^3))}{k} \right) \\ &\leq \frac{\epsilon^3 k^2 k (k^{-3} + k^3 c_k^3)}{3 n^{1/2}} \left( 1 + \frac{\log n}{k} \right) \rightarrow 0 \end{aligned}$$

because  $\epsilon^3 k^2 k (k^{-3} + k^3 c_k^3) \log^2 n = o(n) \rightarrow 0$ . Using the first two results from Lemma C.3, (C.6)-(C.7), we obtain that

$$k \frac{1}{n} \sum_{i=1}^n x_i \quad (x)^\theta (b_k \quad k) \quad (x)^\theta \Omega^{1/2} N_k k \quad k \frac{1}{n} \sum_{i=1}^n x_i \quad (x)^\theta \Omega^{1/2} \quad (x)^\theta \Omega^{1/2} N_k k + \bar{R}_{1n} = o_p(a_n^{-1}):$$

uniformly over  $x \in X$ . Since  $k (x)^\theta \Omega^{1/2} k$  is bounded from below uniformly over  $x \in X$  we obtain the first statement of Theorem 5.2 from which the second statement directly follows.

Finally, under the assumption that  $\sup_{x \in X} n^{1/2} j r(x) = o_p(a_n^{-1})$ ,

$$\frac{n p(x)^\theta (b_k \quad k)}{k s(x) k} \quad \frac{1}{n} \sum_{i=1}^n g(x) \quad g_0(x) = o_p(a_n^{-1})$$

so that the third statement, (5.9) holds.

### Proof of Theorem 5.4

Preliminaries for Proof of Theorem 5.4

**Lemma** (Symmetrization). *Let  $Z_1, \dots, Z_n$  be independent stochastic processes with mean zero and let  $\epsilon_1, \dots, \epsilon_n$  be independent Rademacher random variables generated independently of the data. Then*

$$\mathbb{E} \Phi \left[ \frac{1}{2} \sum_{i=1}^n \epsilon_i Z_i \right] \leq \mathbb{E} \Phi \left[ \sum_{i=1}^n Z_i \right] \leq 2 \mathbb{E} \Phi \left[ \sum_{i=1}^n |Z_i| \right] \quad (\text{SI})$$

for every nondecreasing, convex  $\Phi : \mathbb{R} \rightarrow \mathbb{R}$  and arbitrary functions  $Z_i : F \rightarrow \mathbb{R}$ .

For  $p \geq 1$  consider the Schatten norm  $S_p$  on symmetric  $k \times k$  matrices  $Q$  defined by  $\|Q\|_{S_p} = \left( \sum_{j=1}^k | \lambda_j(Q) |^p \right)^{1/p}$  where  $\lambda_1(Q), \dots, \lambda_k(Q)$  are the eigenvalues of  $Q$ . The case  $p = 1$  recovers the operator norm and  $p = 2$  recovers the Frobenius norm.

**Lemma** (Khinchin's Inequality for Matrices). *For symmetric  $k \times k$  matrices  $Q_i$ ,  $i = 1, \dots, n$ ,  $2 \leq p \leq \infty$ , and an i.i.d. sequence of Rademacher random variables  $\epsilon_1, \dots, \epsilon_n$  we have*

$$\mathbb{E} \| \sum_{i=1}^n \epsilon_i Q_i \|_{S_p}^{1-2/p} \leq \mathbb{E} \| \sum_{i=1}^n Q_i \|_{S_p}^{1-2/p} \leq C^{D/p} \mathbb{E} \| \sum_{i=1}^n Q_i \|_{S_p}^{1-2/p} \quad (\text{KI-1})$$

where  $C$  is an absolute constant. So, for  $k \geq 2$ ,

$$\mathbb{E} \| \sum_{i=1}^n \epsilon_i Q_i \|_{S_p} \leq C^{D/p} \sqrt{\log k} (\mathbb{E} \| \sum_{i=1}^n Q_i \|_{S_p}^2)^{1/2} \quad (\text{KI-2})$$

for some (possibly different) absolute constant  $C$ .

We will establish consistent estimation of

$$\Sigma = \mathbb{E} [ \tilde{r} p^k(x) - (r^k + r_k) g \tilde{r} p^k(x) - (r^k + r_k) g^{\tilde{r}} ]$$

using

$$\hat{\Sigma} = \mathbb{E}_n [ \tilde{r} p^k(x) - b^k g \tilde{r} p^k(x) - b^k g^{\tilde{r}} ]$$

Consistency of  $\hat{\Sigma}$  will then follow from the consistency of  $\hat{\mathcal{Q}}$  established by Lemma C.1. To save notation, define the vectors

$$\tilde{\mathcal{V}} := \begin{pmatrix} Y(b_1; \mathbf{m}_1) \\ \vdots \\ Y(b_k; \mathbf{m}_k) \end{pmatrix} \quad \text{and} \quad \mathcal{V} := \begin{pmatrix} Y(b_1; \mathbf{m}_1) \\ \vdots \\ Y(b_k; \mathbf{m}_k) \end{pmatrix} \quad (\text{A.4})$$

Also define  $\tilde{r}^k := (\tilde{r}_1^k, \dots, \tilde{r}_k^k)$  so that  $\tilde{r}^k := Y(\tilde{r}; \mathbf{m}_j) - \mathbf{g}(x)$ . Ideally, we would like to use  $\tilde{r}^k$  to estimate  $\hat{\Sigma}$ , but we don't observe  $\tilde{r}^k$ . Define  $\Delta := b^k - \tilde{r}^k = \tilde{\mathcal{V}}^k - \mathcal{V}^k \in \mathbb{R}^k$ .

Using this, we can decompose

$$\begin{aligned}
\mathfrak{B} &= E_n[\bar{f}p^k(x) (\Delta + \cdot^k)g\bar{f}p^k(x) (\Delta + \cdot^k)g] \\
&= \underbrace{E_n[\bar{f}p^k(x) \frac{\Delta g\bar{f}p^k(x)}{Z} \Delta g^{\theta}]}_1 + \underbrace{E_n[\bar{f}p^k(x) \frac{\cdot^k g\bar{f}p^k(x)}{Z} \Delta g^{\theta}]}_2 \\
&\quad + \underbrace{E_n[\bar{f}p^k(x) \frac{\Delta g\bar{f}p^k(x)}{Z} \cdot^k g^{\theta}]}_3 + \underbrace{E_n[\bar{f}p^k(x) \frac{\cdot^k g\bar{f}p^k(x)}{Z} \cdot^k g^{\theta}]}_4
\end{aligned} \tag{A.5}$$

We first show that  $k\Sigma_4 = \Sigma k \neq 0$ . This is nonstandard because of the Hadamard product.

**Lemma A.1** (Pseudo-Variance Estimator Consistency). *Suppose Assumption 5.1 and Assumption 5.3 hold. Further, de  $n \rightarrow \infty$ ,  $v_n = E[\max_{1 \leq i \leq n} j_{k,i}^2]^{1/2} = o(1)$ . In addition, assume that  $\bar{R}_{1n} + \bar{R}_{2n} = o(\log k)^{1/2}$ . Then,*

$$\begin{aligned}
k\Sigma_4 &= O(k \cdot P \frac{\sqrt{\frac{2 \log k}{n}}}{n}) = o(1) \\
\text{and } k\Sigma_4 &= \Sigma k \cdot P (v_n - 1 + \cdot^k c_k) \frac{\sqrt{\frac{2 \log k}{n}}}{n}
\end{aligned}$$

*Proof.* The first result is established by Lemma C.1 (Matrix LLN). Rest of proof will follow proof of Theorem 4.6 in Belloni et al. (2015). Like in (A.5) we can define  $\Delta = \cdot^k = g_0(x) - \bar{g}(x)$ <sup>1</sup> and decompose

$$\begin{aligned}
\Sigma_4 &= \underbrace{E_n[p^k(x)p^k(x)^{\theta} \Delta^2]}_{41} + \underbrace{E_n[\bar{f}p^k(x) (\cdot^k + r_k)g\bar{f}p^k(x) \Delta g^{\theta}]}_{42} \\
&\quad + \underbrace{E_n[\bar{f}p^k(x) \frac{\Delta g\bar{f}p^k(x)}{Z} (\cdot^k + r_k)g^{\theta}]}_{43} + \underbrace{E_n[\bar{f}p^k(x) (\cdot^k + r_k) \frac{g\bar{f}p^k(x)}{Z} (\cdot^k + r_k)g]}_{44}
\end{aligned}$$

The terms  $\Sigma_{41}, \Sigma_{42}$  and  $\Sigma_{43}$  are simple to show are negligible.

$$\begin{aligned}
&k\Sigma_{41} + \Sigma_{42} + \Sigma_{43}k \\
&= kE_n[\bar{f}p^k(x)^{\theta} (\cdot^k + r_k)g\bar{f}p^k(x) (\cdot^k + r_k)g] + kE_n[\bar{f}p^k(x) (\cdot^k + r_k)g\bar{f}p^k(x)^{\theta} (\cdot^k + r_k)g] \\
&\quad + kE_n[\bar{f}p^k(x) \bar{f}p^k(x)^{\theta} (\cdot^k + r_k)g\bar{f}p^k(x) (\cdot^k + r_k)g^{\theta}] \\
&\quad + 2 \max_{1 \leq i \leq n} j_{k,i} p^k(x) (\cdot^k + r_k)^2 kE_n[p^k(x)p^k(x)^{\theta}]k \\
&\quad + 2 \max_{1 \leq i \leq n} j_{k,i} + j_{r_k,i} \max_{1 \leq i \leq n} j_{k,i} p^k(x)^{\theta} (\cdot^k + r_k) kE_n[p^k(x)p^k(x)^{\theta}]k
\end{aligned}$$

By Theorem 5.2  $j_{k,i} \max_{1 \leq i \leq n} j_{k,i} p^k(x)^{\theta} (\cdot^k + r_k) j_{k,i} \cdot P \frac{2(\sqrt{\log k} + \bar{R}_{1n} + \bar{R}_{2n})^2}{n} = o(1)$ , by Assumption 5.1 the approximation error is bounded  $\max_{1 \leq i \leq n} j_{r_k,i} \cdot P \cdot c_k$ , by Assumption 5.3 and Markov's inequality the errors are bounded  $\max_{1 \leq i \leq n} j_{k,i} \cdot P \cdot v_n^2$ . Finally, by the first part of Lemma A.1  $k\Sigma_4 = O(k \cdot P \cdot 1)$ . Putting this all together with  $\bar{R}_{1n} + \bar{R}_{2n} = o(\log k)^{1/2}$  and  $\frac{2 \log k}{n} = o(1)$  gives

$$k\Sigma_{41} + \Sigma_{42} + \Sigma_{43}k \cdot P (v_n - 1 + \cdot^k c_k) \frac{\sqrt{\frac{2 \log k}{n}}}{n}$$

Next, we want to control  $\Sigma_{44} = \Sigma$ . To do this, let  $\epsilon_1, \dots, \epsilon_n$  be independent Rademacher random variables

<sup>1</sup>It is useful to recall that  $\cdot^k = \bar{y}^k - \bar{g}(x)$  and  $\cdot^k = \bar{y}^k - g_0(x)$

generated independently from the data. Then for  $\mathbf{x} = (x_1, \dots, x_n)$

$$\begin{aligned} & E[kE_n[\hat{f}p^k(x) - (k+r_k)gf\hat{p}^k(x) - (k+r_k)g^{\hat{\eta}}] - \Sigma k] \\ & \leq E[E_n[k\hat{f}p^k(x) - (k+r_k)gf\hat{p}^k(x) - (k+r_k)g^{\hat{\eta}}k]] \\ & \leq \frac{\log k}{n} E[(kE_n[kp^k(x)k^2 - (k+r_k)^2\hat{f}p^k(x) - (k+r_k)gf\hat{p}^k(x) - (k+r_k)g^{\hat{\eta}}k]^{1=2})] \\ & \leq \frac{\frac{2}{k}\log k}{n} E[\max_{i=1, \dots, n} |j_{k,i} - r_k| (kE_n[\hat{f}p^k(x) - (k+r_k)gf\hat{p}^k(x) - (k+r_k)g^{\hat{\eta}}k]^{1=2})] \\ & \leq \frac{\frac{2}{k}\log k}{n} (E[\max_{i=1, \dots, n} |j_{k,i} - r_k|^2])^{1=2} (E[kE_n[\hat{f}p^k(x) - (k+r_k)gf\hat{p}^k(x) - (k+r_k)g^{\hat{\eta}}k]^{1=2}) \end{aligned}$$

where the first inequality holds from Symmetrization (SI), the second from Khinchin's inequality (KI-1), the third by  $\max_{i=1, \dots, n} |k\hat{p}^k(x)k - k|$  and the fourth by Cauchy-Schwarz inequality.

Since for any positive numbers  $a, b$  and  $R$ ,  $a \leq R(a+b)^{1=2}$  implies  $a \leq R^2 + R^{\frac{1}{2}}b$ , the expression above and the triangle inequality yields

$$\begin{aligned} & E[kE_n[\hat{f}p^k(x) - (k+r_k)gf\hat{p}^k(x) - (k+r_k)g^{\hat{\eta}}] - \Sigma k] \\ & \leq \frac{\frac{2}{k}\log k}{n} (v_n^2 + \frac{2}{k}c_k^2) + \frac{\frac{2}{k}\log k}{n} r_n v_n + \frac{2}{k}c_k^2 g \leq k\Sigma k^{1=2} \end{aligned}$$

and so, because  $k\Sigma k \leq 1$  and  $(v_n^2 + \frac{2}{k}c_k^2) \frac{2}{k}\log k = n! \rightarrow 0$  we have

$$E[kE_n[\hat{f}p^k(x) - (k+r_k)gf\hat{p}^k(x) - (k+r_k)g^{\hat{\eta}}] - \Sigma k] \leq (v_n - 1 + \frac{2}{k}c_k) \frac{\frac{2}{k}\log k}{n}.$$

The second result of Lemma A.1 follows from Markov's inequality.  $\square$

Now, we need to take care of the terms

$$\begin{aligned} \Sigma_1 &= E_n[\hat{f}p^k(x) - \Delta gf\hat{p}^k(x) - \Delta g^{\hat{\eta}}] \\ \Sigma_2 &= E_n[\hat{f}p^k(x) - \cdot^k gf\hat{p}^k(x) - \Delta g^{\hat{\eta}}] \\ \Sigma_3 &= E_n[\hat{f}p^k(x) - \Delta gf\hat{p}^k(x) - \cdot^k g^{\hat{\eta}}] \end{aligned}$$

where  $\Delta = \hat{p}^k - \bar{Y}^k$  and  $\cdot^k = \bar{Y}^k - \mathbf{g}(x) = \mathbf{g}^k(x) - \mathbf{g}(x) + \cdot^k$ . To do so we will use Condition 2.

**Lemma A.2** (Negligible Variance Bias). *Suppose that Condition 2, Assumption 5.1 and Assumption 5.3 hold. Then*

$$k\Sigma_1 + \Sigma_2 + \Sigma_3 k = o_p(1):$$

*Proof.* From Condition 2, the term  $\Sigma_1$  being negligible immediately follows from Cauchy-Schwarz. Notice



that

$$\begin{aligned}
k\Sigma_1 & \leq k \sup_{1 \leq j \leq k} E_n[\rho_l(X)(Y(\hat{\cdot}_l; \hat{m}_l) - Y(\cdot_j; \bar{m}_j))\rho_l(X)(Y(b_j; \bar{m}_j) - Y(\cdot_j; \bar{m}_j))] \\
& \leq k \sup_{1 \leq j \leq k} (E_n[\rho_j(X)^2(Y(b_j; \bar{m}_j) - Y(\cdot_j; \bar{m}_j))^2])^{1/2} \sup_{1 \leq j \leq k} (E_n[\rho_j(X)^2(Y(b_j; \bar{m}_j) - Y(\cdot_j; \bar{m}_j))^2])^{1/2} \\
& = o_p(1).
\end{aligned}$$

To see that  $\Sigma_2$  is negligible notice that

$$\begin{aligned}
k\Sigma_2 & \leq k \sup_{1 \leq j \leq k} E_n[\rho_l(X)(\cdot_l + \rho^k(X)^{\theta(b^k \cdot)})\rho_j(X)(Y(b_j; \bar{m}_j) - Y(\cdot_j; \bar{m}_j))] \\
& \leq k \sup_{1 \leq j \leq k} E_n[\rho_l(X)^2(\cdot_l + \rho^k(X)^{\theta(b^k \cdot)})^2]^{1/2} E_n[\rho_j(X)^2(Y(b_j; \bar{m}_j) - Y(\cdot_j; \bar{m}_j))^2]^{1/2} \\
& \leq k; \gamma (\max_{1 \leq i \leq n} j^{-k} + \max_{1 \leq i \leq n} \rho^k(X)^{\theta(b^k \cdot)}) E_n[\rho_j(X)^2(Y(b_j; \bar{m}_j) - Y(\cdot_j; \bar{m}_j))^2]^{1/2}
\end{aligned}$$

Applying Assumption 5.3 and Theorem 5.2 gives

$$\leq P k; \gamma n^{1-m} E[\rho_j(X)^2(Y(b_j; \bar{m}_j) - Y(\cdot_j; \bar{m}_j))^2]^{1/2} = o_p(1)$$

where the final line is via Condition 2. Showing negligibility of  $\Sigma_3$  follows the same steps.  $\square$

### Proof of Theorem 5.5

Follows from the exact same steps as Theorem 3.5 in [Semenova and Chernozhukov \(2021\)](#) after establishing strong approximation by a gaussian process as in Theorem 5.3 and consistent variance estimation as in Theorem 5.4.

## B Supporting Lemmas for First Stage

Here we provide supporting lemmas and their proofs. We start off with non-asymptotic bounds for first stage parameters and means.

### B.1 Nonasymptotic Bounds for the First Stage

The nonasymptotic bounds for the first stage will depend on certain events. In Appendix B.3 we will show that under Assumptions 4.1 and 4.2 these events happen with probability approaching one. Define the scores

$$\begin{aligned}
S_{\cdot j} & := E_n[U_{\cdot j} Z] \\
S_{\cdot j} & := E_n[U_{\cdot j} Z]
\end{aligned} \tag{B.1}$$

With these in mind, we will consider nonasymptotic bounds under the events:

$$\begin{aligned}
\Omega_{k;1} & := f_{\cdot j} \leq c_0 k S_{\cdot j} k_1; \delta_j \leq k g \\
\Omega_{k;2} & := f_{\cdot j} \leq \bar{c}_k; \delta_j \leq k g
\end{aligned} \tag{B.2}$$

Following Chetverikov and Sørensen (2021), the first event is referred to as “score domination” while the second event is referred to as “penalty majorization”.

**Lemma B.1** (Nonasymptotic Bounds for Logistic Model). *Suppose that Assumption 4.1 holds with  $c_0 > (c_0 + 1) = (c_0 - 1)$  and  $2C_0 \sigma_k^{-2} < 1$ . Then, under the events  $\Omega_{k,1} \setminus \Omega_{k,2}$  defined in (B.2), there exists a finite constant  $M_0$  that does not depend on  $k$  such that*

$$\max_{1 \leq j \leq k} D^z(\hat{g}; \mathbf{g}) \leq M_0 \sigma_k^{-2} \text{ and } \max_{1 \leq j \leq k} k b_j^{-1} k_1 \leq M_0 \sigma_k^{-1} \quad (\text{B.3})$$

*Proof.* We show that the bound of (B.3) holds for each  $j = 1; \dots; k$ . For any  $u \geq \mathbb{R}^d$  define  $\tilde{g}_j(\cdot) := E_n[\rho_j(X) D e^{-\langle \cdot, Z \rangle} + (1 - D) e^{-\langle \cdot, Z \rangle}]$ . By optimality of  $b_j$  we must have, for any  $u \geq (0; 1]$ :

$$\tilde{g}_j(b_j) + \sum_{l \in S_j} k b_l k_1 \leq \tilde{g}_j((1 - u)b_j + u^{-1}k) + \sum_{l \in S_j} k(1 - u)b_l + u^{-1}k_1$$

Using convexity of the  $\ell_1$  norm  $k \cdot k_1$ , this gives after rearrangement

$$\tilde{g}_j(b_j) \leq \tilde{g}_j((1 - u)b_j + u^{-1}k) + \sum_{l \in S_j} u k b_l k_1 - \sum_{l \in S_j} u k^{-1} k_1$$

Divide both sides by  $u$  and let  $u \downarrow + 0$

$$E_n[\rho_j(X) D e^{-\langle b_j, Z \rangle} + (1 - D) g_j^0 Z^{-1} Z g_j] + \sum_{l \in S_j} k b_l k_1 \leq \sum_{l \in S_j} k^{-1} k_1$$

By direct calculation, we have that  $D^z_{\cdot; j}(b_j; \bar{b}_j)$  from (4.6) can be expressed

$$D^z_{\cdot; j}(b_j; \bar{b}_j) = E_n[\rho_j(X) D e^{-\langle b_j, Z \rangle} + (1 - D) g_j^0 Z^{-1} Z g_j] - E_n[\rho_j(X) D e^{-\langle \bar{b}_j, Z \rangle} + (1 - D) g_j^0 Z^{-1} Z g_j]$$

Combining the last two displays yields

$$D^z_{\cdot; j}(b_j; \bar{b}_j) + E_n[\rho_j(X) D e^{-\langle \bar{b}_j, Z \rangle} + (1 - D) g_j^0 Z^{-1} Z g_j] + \sum_{l \in S_j} k b_l k_1 \leq \sum_{l \in S_j} k^{-1} k_1 \quad (\text{L.1})$$

In the event  $\Omega_{k,1}$  we have that

$$\sum_{l \in S_j} E_n[\rho_j(X) D e^{-\langle \bar{b}_l, Z \rangle} + (1 - D) g_l^0 Z^{-1} Z g_l] \leq c_0^{-1} \sum_{l \in S_j} k b_l k_1 \quad (\text{L.2})$$

Combining (L.1) and (L.2) yields

$$D^z_{\cdot; j}(b_j; \bar{b}_j) + \sum_{l \in S_j} k b_l k_1 \leq \sum_{l \in S_j} k^{-1} k_1 + c_0^{-1} \sum_{l \in S_j} k b_l k_1$$

Expanding  $k \cdot k_1 = \sum_{l \in S_j} j \cdot l + \sum_{l \in S_j} j \cdot l$  for  $\cdot = b_j; \bar{b}_j$  and applying the triangle inequalities  $j b_{j;l} j \bar{b}_{j;l} j b_{j;l} - j \bar{b}_{j;l} j$  for  $l \in S_j$  and the equality  $b_{j;l} = b_{j;l} - \bar{b}_{j;l}$  gives

$$D^z_{\cdot; j}(b_j; \bar{b}_j) + \sum_{l \in S_j} \sum_{l \in S_j} j \cdot l + \sum_{l \in S_j} \sum_{l \in S_j} j \cdot l + \sum_{l \in S_j} \sum_{l \in S_j} j \cdot l + c_0^{-1} \sum_{l \in S_j} \sum_{l \in S_j} j \cdot l + c_0^{-1} \sum_{l \in S_j} \sum_{l \in S_j} j \cdot l$$

Rearrange to get

$$D^z_{\cdot j}(b_{j\cdot}; \bar{y}_j) + (1 - c_0) \prod_{l \in S} b_{j:l} \bar{y}_{j:l} = (1 + c_0) \prod_{l \in S} b_{j:l} \bar{y}_{j:l}$$

Adding  $(1 - c_0) \prod_{l \in S} b_{j:l} \bar{y}_{j:l}$  gives

$$D^z_{\cdot j}(b_{j\cdot}; \bar{y}_j) + (1 - c_0) \prod_{l \in S} b_{j:l} \bar{y}_{j:l} = 2 \prod_{l \in S} b_{j:l} \bar{y}_{j:l} \quad (\text{L.3})$$

By Lemma 4 in Appendix V.3 of Tan (2017) we have that for  $y_j := b_{j\cdot} \bar{y}_j$

$$D^z_{\cdot j}(b_{j\cdot}; \bar{y}_j) = \frac{1 - e^{-c_0 k_j k_1}}{c_0 k_j k_1} \tilde{\Sigma}_{\cdot j} \quad (\text{L.4})$$

By (L.3) and  $c_0 > (c_0 + 1) = (c_0 - 1)$  we have that  $\prod_{l \in S} b_{j:l} \bar{y}_{j:l} \leq \prod_{l \in S} b_{j:l} \bar{y}_{j:l}$ . Applying the empirical compatibility condition from Assumption 4.1 to (L.3) then yields

$$D^z_{\cdot j}(b_{j\cdot}; \bar{y}_j) + (1 - c_0) \prod_{l \in S} b_{j:l} \bar{y}_{j:l} = 2 \prod_{l \in S} b_{j:l} \bar{y}_{j:l}^{1-2} (\tilde{\Sigma}_{\cdot j})^{1-2} \quad (\text{L.5})$$

Combining (L.4) and (L.5) to get an upper bound on  $(\tilde{\Sigma}_{\cdot j})^{1-2}$  gives

$$\prod_{l \in S} b_{j:l} \bar{y}_{j:l} \leq (\tilde{\Sigma}_{\cdot j})^{1-2} = 2 \prod_{l \in S} b_{j:l} \bar{y}_{j:l}^{1-2} \frac{c_0 k_j k_1}{1 - e^{-c_0 k_j k_1}}.$$

Plugging the second bound into (L.5) gives

$$D^z_{\cdot j}(b_{j\cdot}; \bar{y}_j) + (1 - c_0) \prod_{l \in S} b_{j:l} \bar{y}_{j:l} = 2 \prod_{l \in S} b_{j:l} \bar{y}_{j:l}^{1-2} \frac{c_0 k_j k_1}{1 - e^{-c_0 k_j k_1}}.$$

The second inequality and  $\prod_{l \in S} b_{j:l} \bar{y}_{j:l} \leq \prod_{l \in S} b_{j:l} \bar{y}_{j:l}$  imply  $1 - e^{-c_0 k_j k_1} \geq 2c_0 \prod_{l \in S} b_{j:l} \bar{y}_{j:l}$  so,

$$\frac{1 - e^{-c_0 k_j k_1}}{c_0 k_j k_1} = \int_0^1 e^{-c_0 k_j k_1 u} du \leq e^{-c_0 k_j k_1} = 1 - e^{-c_0 k_j k_1}.$$

Combining the last two displays gives

$$D^z_{\cdot j}(b_{j\cdot}; \bar{y}_j) + (1 - c_0) \prod_{l \in S} b_{j:l} \bar{y}_{j:l} = 4 \prod_{l \in S} b_{j:l} \bar{y}_{j:l}^2 (1 - e^{-c_0 k_j k_1}) \quad (\text{L.6})$$

Applying  $\Omega_{k,2}$  to bound  $\prod_{l \in S} b_{j:l} \bar{y}_{j:l}$  and noting that  $\prod_{l \in S} b_{j:l} \bar{y}_{j:l} \leq s_k$  by definition gives (B.3) with  $M_0 = \frac{4}{1 - c_0} \frac{1 - e^{-c_0 k_j k_1}}{c_0 k_j k_1}$ .  $\square$

For each  $j$ , consider the matrices,

$$\begin{aligned} \tilde{\Sigma}_{\cdot j} &:= E_n[\rho_j(X) D e^{\tilde{y}_j^Z (Y - \tilde{y}_j^Z)^2 Z Z^T}] \\ \tilde{\Sigma}_{\cdot j} &:= E_n[\rho_j(X) D e^{\tilde{y}_j^Z Z Z^T}] \end{aligned} \quad (\text{B.4})$$

In addition define  $\Sigma_{\cdot j} := E \tilde{\Sigma}_{\cdot j}$  and  $\Sigma_{\cdot j} := E \tilde{\Sigma}_{\cdot j}$ . For the outcome regression model, we will consider



In the event  $\Omega_{k,1} \setminus \Omega_{k,2}$  using the bound in (B.3) and the fact that  $C_0 \rho_0^2 s_k^{-k} < 1$  gives us that

$$C_0 k b_j^{-j} k_1^{-k} C_0 M_0 s_k^{-k} M_0 : \quad (O.3)$$

In the event  $\Omega_{k,1} \setminus \Omega_{k,2}$  the bound in (L.6) also gives us that  $D^z_{\cdot j}(b_j; \bar{j}) \leq M_0 s_k^{-2j}$ . Combining the above displays then yields

$$M_0 s_k^{-2j} D^z_{\cdot j}(b_j; \bar{j}) \leq e^{M_0} E_n[\rho_j(X) D e^{j^z f b_j^0 Z^{-j} Z g^2}]: \quad (O.4)$$

Again applying the bound on  $C_0 k b_j^{-j} k_1^{-k}$  (O.3) gives

$$\begin{aligned} D^z_{\cdot j}(b_j; \bar{j}; b_j) &= E_n[\rho_j(X) D e^{b_j^z (b_j^0 Z^{-j} Z)^2}] \\ &= E_n[\rho_j(X) D e^{(b_j - j)^0 Z e^{j^z (b_j^0 Z^{-j} Z)^2}}] \\ &\leq e^{M_0} D^z_{\cdot j}(b_j; \bar{j}; \bar{j}) \end{aligned} \quad (O.5)$$

Decomposing the empirical expectation on the RHS of (O.1) gives

$$\begin{aligned} (b_j - j)^0 E_n[\rho_j(X) D e^{b_j^z f Y^{-j} Z g Z}] &= \underbrace{(b_j - j)^0 E_n[\rho_j(X) D e^{j^z f Y^{-j} Z g Z}]}_{1:j} \\ &\quad + \underbrace{E_n[\rho_j(X) D f e^{b_j^z} e^{j^z g f Y^{-j} Z g f b_j^0 Z^{-j} Z g}]}_{2:j} \end{aligned}$$

By Hölder's inequality, in the event  $\Omega_{k,3}$ ,  $1:j$  is bounded

$$1:j \leq C_0^1 k b_j^{-j} k_1^{-k} :j \quad (O.6)$$

By the mean value equation (O.2) and the Cauchy-Schwarz inequality,  $2:j$  can be bounded from above by

$$\begin{aligned} 2:j &\leq e^{C_0 k b_j^{-j} k_1^{-k}} E_n^{1=2}[\rho_j(X) D e^{j^z f b_j^0 Z^{-j} Z g^2}] \\ &\quad E_n^{1=2}[\rho_j(X) D e^{j^z f Y^{-j} Z g^2 f b_j^0 Z^{-j} Z g^2}] \end{aligned} \quad (O.7)$$

Using (O.3) the first term in (O.7) can be bounded by  $e^{M_0}$ . The second term is exactly the square root of  $D^z_{\cdot j}(b_j; \bar{j}; \bar{j})$ . The third term is bounded in a few steps. First, in the event  $\Omega_{k,5}$  we have that

$$(E_n - E)[\rho_j(X) D e^{j^z f Y^{-j} Z g^2 f b_j^0 Z^{-j} Z g} \bar{k} k b_j^{-j} k_1^{-k}]:$$

By Assumption 4.2 and Lemma E.6 we have that  $E[D f Y^{-j} Z g^2] \leq D_0^2 + D_1^2$  so that:

$$E[\rho_j(X) D e^{j^z f Y^{-j} Z g^2 f b_j^0 Z^{-j} Z g^2}] \leq (D_0^2 + D_1^2) E[\rho_j(X) D e^{j^z f b_j^0 Z^{-j} Z g^2}]:$$

In the event  $\Omega_{k,6}$  we have that

$$(E_n - E)[\rho_j(X) D e^{j^z f b_j^0 Z^{-j} Z g^2}] \leq \bar{k} k b_j^{-j} k_1^{-k}:$$

and we can bound  $E_n[\rho_j(X)De^{jZ}f\mathbb{b}_j^0Z^{-j}Zg^2]$  using (O.4). Putting this together gives

$$E_n[\rho_j(X)De^{jZ}fY^{-j}Zg^2f\mathbb{b}_j^0Z^{-j}Zg^2]^{-k}kb_j^{-j}k_1^2 + (D_0^2 + D_1^2)^{-k}kb_j^{-j}k_1^2 + (D_0^2 + D_1^2)e^{M_0}M_0s_k^{-2}j \quad (\text{O.8})$$

Applying convexity of  $f^{\rho_-}$  and the bounds on  $kb_j^{-j}k_1^2$  in the event  $\Omega_{k,1} \setminus \Omega_{k,2}$  from (L.6) gives

$$2j \quad f e^{M_0} (1 + (D_0^2 + D_1^2)^{1=2})(M_0^{-k} j s_k)^{1=2} + (D_0^2 + D_1)^2 (M_0 s_k^{-2} j)^{1=2} g D^Z j (b_j;^{-j};^{-j})^{1=2} + \tilde{C} f (s_k^{-k} j s_k)^{1=2} + (s_k^{-j})^{1=2} g D^Z j (b_j;^{-j};^{-j})^{1=2} \quad (\text{O.9})$$

where  $\tilde{C} = \max f e^{M_0} M_0^{1=2} (1 + D_0 + D_1); (D_0^2 + D_1^2) M_0^{1=2} g$ . Combining (O.6) and (O.9) gives a bound on the empirical expectation on the RHS of (O.1).

$$(b_j;^{-j})^0 E_n[\rho_j(X)De^{b_j^0 Z} f Y^{-j} Z g Z] \underbrace{\{c_0^{-1} kb_j^{-j} k_1^{-j}\}}_{\text{Bound on } 1:j \text{ from (O.6)}} + \underbrace{\{\tilde{C} f (s_k^{-k} j s_k)^{1=2} + (s_k^{-2} j)^{1=2} g D^Z j (b_j;^{-j};^{-j})^{1=2}\}}_{\text{Bound on } 2:j \text{ from (O.9)}} \quad (\text{O.10})$$

For convenience, we will sometimes continue to refer to the bound on  $2:j$  from (O.9) as simply  $2:j$ .

*Step 3: Express Minimization Constraint in Terms of  $^{-j}$  and Simplify.* We use the results from Step 2 to rewrite the minimization bound (O.1) from Step 1. Using (O.5) and (O.10) together with the minimization bound (O.1) yields

$$e^{M_0} D^Z j (b_j;^{-j};^{-j}) + j kb_j k_1 c_0^{-1} j kb_j^{-j} k_1 + j k^{-j} k_1 + 2:j \quad (\text{O.11})$$

Apply the triangle inequality  $j|b_{j;l}| - j|^{-j;l}| - j|b_{j;l}^{-j;l}|$  for  $l \in S_j$  and  $j|b_{j;l}| = j|b_{j;l}^{-j;l}|$  for  $l \notin S_j$  to the above to obtain

$$e^{M_0} D^Z j (b_j;^{-j};^{-j}) + (1 - c_0^{-1}) kb_j^{-j} k_1 - 2 \sum_{l \in S_j} j |b_{j;l}^{-j;l}| + 2:j$$

Let  $j = b_j^{-j}$ . We use the form  $D^Z j (b_j;^{-j}) = E_n[\rho_j(X)De^{jZ}f\mathbb{b}_j^0Z^{-j}Zg^2] = j \tilde{\Sigma} j j$  to expand out

$$e^{M_0} (j \tilde{\Sigma} j j) + (1 - c_0^{-1}) j k k_1 - 2 \sum_{l \in S_j} j |j;l| + \tilde{C} f (s_k^{-k} j)^{1=2} + (s_k^{-j})^{1=2} g (j \tilde{\Sigma} j j)^{1=2} \quad (\text{O.12})$$

*Step 4: Apply Empirical Compatability Condition.* Let  $3:j := \tilde{C} f (s_k^{-k} j)^{1=2} + (s_k^{-j})^{1=2} g$  and  $D^Z j := e^{M_0} (j \tilde{\Sigma} j j) + (1 - c_0^{-1}) j k k_1$ . In the even  $\Omega_{k,1} \setminus \Omega_{k,2} \setminus \Omega_{k,3} \setminus \Omega_{k,5} \setminus \Omega_{k,6}$  that (O.12) holds, there are two possibilities. For  $2 = 1 - 2c_0 = f(1 + 1)(c_0 - 1)g \geq (0; 1]$  either

$$2 D^Z j - 3:j (j \tilde{\Sigma} j j)^{1=2} \quad (\text{O.13})$$

or  $(1 - c_0)D_{j,j}^2 \leq 2 \sum_{l \in S_{j,j}} P_{j,j;l}$ , that is

$$D_{j,j}^2 \leq (1 + 1)(c_0 - 1)c_0^{-1} \sum_{l \in S_{j,j}} P_{j,j;l} \quad (\text{O.14})$$

We deal with these two cases separately. First, if (O.14) holds, then  $\sum_{l \in S_{j,j}} P_{j,j;l} \leq \sum_{l \in S_{j,j}} P_{j,j;l}$ . We can apply the empirical compatibility of Assumption 4.2 to (O.14) to obtain.

$$e^{-M_0} \left( \sum_{j,j} \tilde{\Sigma}_{j,j} \right) + (1 - c_0^{-1}) \sum_{j,k} k_{j,k} \leq (1 + 1)(1 - 1) \sum_{j,j} (S_{j,j} \tilde{\Sigma}_{j,j})^{1=2}.$$

Inverting for  $(\sum_{j,j} \tilde{\Sigma}_{j,j})^{1=2}$  and plugging in gives

$$e^{-M_0} D_{j,j}^z(b_j; \bar{;} \bar{;} + (1 - c_0^{-1}) \sum_{j,k} k_{j,k} \bar{;} \bar{;} \tilde{M} s_{k,j}^2 \quad (\text{O.15})$$

where  $\tilde{M} = e^{M_0} (1 + 1)(c_0 - 1)c_0^{-1}$ . Next, assume that (O.13) holds. In this case, we can directly invert for  $(\sum_{j,j} \tilde{\Sigma}_{j,j})^{1=2}$  to get that

$$e^{-M_0} D_{j,j}^z(b_j; \bar{;} \bar{;} + (1 - c_0^{-1}) \sum_{j,k} k_{j,k} \bar{;} \bar{;} \tilde{C} f(s_{k,j}^-)^{1=2} + (s_{k,j}^2)^{1=2} g^2 \quad (\text{O.16})$$

Combining (O.15) and (O.16) gives

$$e^{-M_0} D_{j,j}^z(b_j; \bar{;} \bar{;} + (1 - c_0^{-1}) \sum_{j,k} k_{j,k} \bar{;} \bar{;} \tilde{M} s_{k,j}^2 + \tilde{C} f(s_{k,j}^-)^{1=2} + (s_{k,j}^2)^{1=2} g^2 \quad (\text{O.17})$$

*Step 5: Apply Penalty Majorization and Bounded Penalty Ratio.* Use the fact that  $\bar{;} = \bar{;} c^{-1}$  to express (O.17) as

$$D_{j,j}^z(b_j; \bar{;} \bar{;} = e^{M_0} \tilde{M} s_{k,j}^2 + e^{M_0} \tilde{C} f(s_{k,j}^-)^{1=2} + (s_{k,j}^2)^{1=2} g^2 + \sum_{j,k} k_{j,k} (1 - c_0^{-1})^{-1} \tilde{M} s_{k,j} + (1 - c_0^{-1})^{-1} c^{-1} \tilde{C} f(s_{k,j}^-)^{1=2} + (s_{k,j}^2)^{1=2} g^2$$

In the event  $\Omega_{k,2} \setminus \Omega_{k,3}$  we have that  $\bar{;} = \bar{;} \bar{;}_{k,j}$ , so that the above simplifies to

$$D_{j,j}^z(b_j; \bar{;} \bar{;} = M_1 s_{k,j}^{-2} + \sum_{j,k} k_{j,k} \bar{;} \bar{;} M_1 s_{k,j}^{-2} \quad (\text{O.18})$$

for  $M_1 = \max f e^{M_0} ; c^{-1} (1 - c_0^{-1})^{-1} g (\tilde{M} + 2e^{M_0} \tilde{C})$ . This completes the result (B.6).  $\square$

Let  $\tilde{\Sigma}_{j,j}^1 := E_n[\rho_j(X) D e^{-\sum_{j,j} Z_j Y_j - \sum_{j,j} Z_j Z_j^0}]$  and  $\Sigma_{j,j}^1 := E \tilde{\Sigma}_{j,j}^1$ . Define the event

$$\Omega_{k,7} := f k \tilde{\Sigma}_{j,j}^1 \leq \Sigma_{j,j}^1 ; k_{j,k} \bar{;}_{k,j} \leq g \quad (\text{B.7})$$

**Lemma B.3** (Nonasymptotic Bounds for Weighted Means). *Suppose that Assumptions 4.1 and 4.2 hold,  $c_0 > (c_0 + 1)(c_0 - 1)$ , and  $2C_0 c_0^{-2} s_{k,j}^- < 1$ . In addition, assume there is a constant  $c > 0$  such that  $\bar{;} = \bar{;} c$  for all  $j, k$ . Then, under the event  $\bigcap_{m=1}^7 \Omega_{k,m}$ , there is a constant  $M_2$  that does not depend*

on  $k$  such that

$$\max_{j=1, \dots, k} |E_n[\rho_j(X)Y(b_j; \mathbf{m}_j)] - E_n[\rho_j(X)Y(\bar{b}_j; \bar{\mathbf{m}}_j)]| \leq M_2 S_k^{-2} \quad (\text{B.8})$$

*Proof.* We show that the bound of (B.8) holds for any  $j = 1, \dots, k$  in a couple steps. To save notation, define

$$\begin{aligned} \rho_j(\cdot; \mathbf{m}) &:= E_n[\rho_j(X)Y(\cdot; \mathbf{m})] \\ &= E_n[\rho_j(X) \frac{DY}{(Z)} + \frac{D}{(Z)} \mathbb{1}(\cdot = \mathbf{m}(Z))] \end{aligned}$$

*Step 1: Decompose Difference and Use Logistic FOCs.* Consider the following decomposition

$$\begin{aligned} \rho_j(b_j; \mathbf{m}_j) - \rho_j(\bar{b}_j; \bar{\mathbf{m}}_j) &= E[\rho_j(X)f(\mathbf{m}_j(Z)) - \rho_j(X)f(\bar{\mathbf{m}}_j(Z))] + \frac{D}{\bar{b}_j(X)} \\ &\quad + E_n[\rho_j(X)DfY - \rho_j(X)Df\bar{Y}] + \frac{1}{b_j(Z)} - \frac{1}{\bar{b}_j(Z)} \\ &\quad + E_n[\rho_j(X)f(\mathbf{m}_j(Z)) - \rho_j(X)f(\bar{\mathbf{m}}_j(Z))] + \frac{D}{\bar{b}_j(Z)} - \frac{D}{b_j(Z)} \\ &:= \rho_{1j} + \rho_{2j} + \rho_{3j} \end{aligned}$$

Notice that  $\rho_{1j} + \rho_{3j} = (b_j - \bar{b}_j)^0 E_n[\rho_j(X)(1 - D=b_j(Z))Z]$ . By the first order conditions for  $b_j$  we have that

$$|E_n[\rho_j(X)f(Z) - DZ=b_j(Z)g]| \leq \rho_{1j} \leq \rho_{3j} \leq k E_n[\rho_j(X)f(Z) - DZ=b_j(Z)g] \leq k \rho_{1j} \leq \rho_{3j}$$

Applying Hölder's inequality to  $\rho_{1j} + \rho_{3j}$  then gives us that on the event  $\Omega_{k,2}$

$$|\rho_{1j} + \rho_{3j}| \leq k |b_j - \bar{b}_j| k_1 \leq k |b_j - \bar{b}_j| k^{-k}$$

By Lemma B.2 on the event  $\bigcap_{m=1}^6 \Omega_{k,m}$  and under the conditions of Lemma B.3,  $k|b_j - \bar{b}_j| \leq M_1 S_k^{-k}$  where  $M_1$  is a constant that does not depend on  $k$ . So

$$|\rho_{1j} + \rho_{3j}| \leq M_1 S_k^{-2} \quad (\text{M.1})$$

*Step 2: Use Outcome Regression Score Domination to Bound  $\rho_{2j}$ .* Now deal with the term  $\rho_{2j}$ . By first order Taylor expansion, for some  $u \in (0, 1)$

$$\begin{aligned} \rho_{2j} &= (b_j - \bar{b}_j)^0 E_n[\rho_j(X)DfY - \rho_j(X)Df\bar{Y}] + \frac{1}{b_j(Z)} - \frac{1}{\bar{b}_j(Z)} \\ &\quad + (b_j - \bar{b}_j)^0 E_n[\rho_j(X)DfY - \rho_j(X)Df\bar{Y}] + \frac{1}{b_j(Z)} - \frac{1}{\bar{b}_j(Z)} \\ &:= \rho_{21j} + \rho_{22j} \end{aligned}$$

In the event  $\Omega_{k,1} \setminus \Omega_{k,2} \setminus \Omega_{k,3} \setminus \Omega_{k,4}$  we have by score domination of the linear outcome regression model and Lemma B.1 that  $\rho_{21j} \leq M_0 S_k^{-2}$ .

The term  $\rho_{22j}$  is second order. On the event  $\Omega_{k,0} \setminus \Omega_{k,1}$  where  $k|b_j - \bar{b}_j| \leq M_0 S_k^{-k}$   $M_0 = C_0$  it can be



bounded with

$$\begin{aligned} & e^{C_0 k b_j} \mathbb{E}_{n, j, k_1} [\rho_j(X) D e^{j Z} j Y \bar{m}_j(Z) j \bar{f}_j^\ell Z^{-\ell} Z g^2] \\ & e^{M_0} \mathbb{E}_n [\rho_j(X) D e^{j Z} j Y \bar{m}_j(Z) j \bar{f}_j^\ell Z^{-\ell} Z g^2]: \end{aligned}$$

This in turn is bounded in a few steps. First note on the event  $\Omega_{k,7}$

$$(\mathbb{E}_n - \mathbb{E})[\rho_j(X) D e^{j Z} j Y \bar{m}_j(Z) j \bar{f}_j^\ell Z^{-\ell} Z g^2] \leq k b_j \bar{m}_j k_1^2:$$

By Assumption 4.2 and Lemma E.8 we have that  $D_0^2 E[D_j Y \bar{m}_j(Z) j j Z] = D_1^2 = D_0 + D_0$  so that,

$$\mathbb{E}[\rho_j(X) D e^{j Z} j Y \bar{m}_j(Z) j \bar{f}_j^\ell Z^{-\ell} Z g^2] = (D_1^2 = D_0 + D_0) \mathbb{E}[\rho_j(X) D e^{j Z} \bar{f}_j^\ell Z^{-\ell} Z g^2].$$

On the event  $\Omega_{k,6}$  we have that

$$(\mathbb{E}_n - \mathbb{E})[\rho_j(X) D e^{j Z} \bar{f}_j^\ell Z^{-\ell} Z g^2] \leq k b_j \bar{m}_j k_1:$$

Putting these all together gives

$$\begin{aligned} & \mathbb{E}_n[\rho_j(X) D e^{j Z} j Y \bar{m}_j(Z) j \bar{f}_j^\ell Z^{-\ell} Z g^2] \leq k b_j \bar{m}_j k_1^2 + (D_1^2 = D_0 + D_0) k b_j \bar{m}_j k_1^2 \\ & + (D_1^2 = D_0 + D_0) \mathbb{E}_n[\rho_j(X) D e^{j Z} \bar{f}_j^\ell Z^{-\ell} Z g^2] \end{aligned} \quad (\text{M.2})$$

To bound (M.2) note again that in the event  $\Omega_{k,1} \setminus \Omega_{k,2}$ ,  $k b_j \bar{m}_j k_1 \leq M_0 s_k^{-k}$  and that using (O.4)

$$\mathbb{E}_n[\rho_j(X) D e^{j Z} \bar{f}_j^\ell Z^{-\ell} Z g^2] \leq e^{M_0} M_0 s_k^{-2k}.$$

Plugging these into (M.2) gives

$$22:j \leq e^{M_0} M_0^2 s_k^{-3k} + e^{M_0} (D_1^2 = D_0 + D_0) M_0^2 s_k^{-3k} + (D_1^2 = D_0 + D_0) M_0 s_k^{-2k} \quad (\text{M.3})$$

so that in total  $2:j = 21:j + 22:j$  is bounded

$$2:j \leq M_0 s_k (D_1^2 = D_0 + D_0 + 1)^{-2k} + e^{M_0} M_0^2 s_k^2 (D_1^2 = D_0 + D_0 + 1)^{-3k} \quad (\text{M.4})$$

*Step 3: Combine Terms.* Putting this together yields

$$\begin{aligned} & j 1:j + 2:j + 3:j \leq f M_1 + M_0 (D_1^2 = D_0 + D_0 + 1) g s_k^{-2k} \\ & + e^{M_0} (D_1^2 = D_0 + D_0) M_0^2 s_k^{-3k} \end{aligned} \quad (\text{M.5})$$

Use the fact that  $s_k^{-k} < 1$  to simplify the last term of this expression

$$\begin{aligned} & j 1:j + 2:j + 3:j \leq f M_1 + M_0 (D_1^2 = D_0 + D_0 + 1) g s_k^{-2k} \\ & + e^{M_0} (D_1^2 = D_0 + D_0) M_0^2 s_k^{-k} \end{aligned} \quad (\text{M.6})$$

This gives the result (B.8) after taking  $M_2 = M_1 + M_0 (D_1^2 = D_0 + D_0 + 1) + e^{M_0} (D_1^2 = D_0 + D_0) M_0^2$ .

□

**Lemma B.4** (Nonasymptotic Bounds for Variance Estimation). *Suppose that Assumptions 4.1 and 4.2 hold,  $c_0 > (c_0 + 1) = (c_0 - 1)$ , and  $2C_0 \sigma_k^{-2} s_k^{-k} < 1$ . In addition, assume there is a constant  $c > 0$  such that  $\bar{y}_j = \bar{y}_j - c$  for all  $j = 1, \dots, k$ . Then, under the event  $\bigcap_{m=1}^7 \Omega_{k,m}$ , there is a constant  $M_3$  that does not depend on  $k$  such that*

$$\max_{j=1, \dots, k} \mathbb{E}_n[\rho_j^2(X)(Y(b_j; \bar{m}_j) - Y(\bar{y}_j; \bar{m}_j))^2] \leq M_3 \frac{2}{k;1} s_k^{-2} \quad (\text{B.9})$$

*Proof.* We show the bound holds for each  $j = 1, \dots, k$ . We start by decomposing

$$\begin{aligned} \rho_j(X)(Y(b_j; \bar{m}_j) - Y(\bar{y}_j; \bar{m}_j)) &= \rho_j(X) f \bar{m}_j(Z) - \bar{m}_j(Z) g^{-1} \frac{D}{\bar{y}_j(X)} \\ &\quad + \rho_j(X) D f Y - \bar{m}_j(Z) g \frac{1}{b_j(Z)} - \frac{1}{\bar{y}_j(Z)} \\ &\quad + \rho_j(X) f \bar{m}_j(Z) - \bar{m}_j(Z) g \frac{D}{\bar{y}_j(Z)} - \frac{D}{b_j(Z)} \\ &:= \tilde{1}_{;j} + \tilde{2}_{;j} + \tilde{3}_{;j} \end{aligned}$$

We will use the fact that  $(a + b + c)^2 \leq 4a^2 + 4b^2 + 4c^2$  to bound

$$\mathbb{E}_n[\rho_j^2(X)(Y(b_j; \bar{m}_j) - Y(\bar{y}_j; \bar{m}_j))^2] \leq 4\mathbb{E}_n[\tilde{1}_{;j}^2] + 4\mathbb{E}_n[\tilde{2}_{;j}^2] + 4\mathbb{E}_n[\tilde{3}_{;j}^2]: \quad (\text{V.1})$$

To bound  $\mathbb{E}_n[\tilde{2}_{;j}^2]$  use the mean value equation (O.2) and the lower bound on  $\bar{g}_j(Z)$  from Assumption 4.1

$$\begin{aligned} \mathbb{E}_n[\tilde{2}_{;j}^2] &= \mathbb{E}_n[\rho_j^2(X) D f Y - \bar{m}_j(Z) g^2 f b_j^{-1}(Z) - \bar{y}_j^{-1}(Z) g^2] \\ &\quad_{k;1} e^{-B_0} (1 + e^{C_0 k b_j - j k_1})^2 \mathbb{E}_n[\rho_j(X) D e^{-j Z} f Y - \bar{m}_j(Z) g^2 f \bar{g}_j(Z) - \bar{g}_j(Z) g^2] \end{aligned}$$

Applying (O.8), Lemma B.1, and  $s_k^{-k} < 1$  there is a constant  $\tilde{M}_1$  that does not depend on  $k$  such that in the event  $\bigcap_{m=1}^7 \Omega_{k,m}$  this is bounded

$$\tilde{M}_1 \frac{2}{k;1} s_k^{-2} \quad (\text{V.2})$$

To bound  $\mathbb{E}_n[\tilde{3}_{;j}^2]$  write  $b_j^{-1}(Z) - \bar{y}_j^{-1}(Z) = e^{-j Z} f e^{b_j^0 Z + j^0 Z} - 1 g$  and use the lower bound on  $\bar{g}_j(Z)$  from Assumption 4.1:

$$\begin{aligned} \mathbb{E}_n[\tilde{3}_{;j}^2] &= \mathbb{E}_n[\rho_j^2(X) D f \bar{m}_j(Z) - \bar{m}_j(Z) g^2 f b_j^{-1}(Z) - \bar{y}_j^{-1}(Z) g^2] \\ &\quad_{k;1} e^{-B_0} (1 + e^{C_0 k b_j - j k_1})^2 \mathbb{E}_n[\rho_j(X) e^{-j Z} f \bar{m}_j(Z) - \bar{m}_j(Z) g^2] \end{aligned}$$

Noting that this last term is exactly  $D^2_{;j}(b_j; \bar{y}_j; \bar{y}_j)$  and applying Lemma B.2, there is a constant  $\tilde{M}_2$  that does not depend on  $k$  such that on the event  $\bigcap_{m=1}^6 \Omega_{k,m}$  this is bounded

$$\tilde{M}_2 \frac{2}{k;1} s_k^{-2} \quad (\text{V.3})$$

Finally, to bound  $E_n[\tilde{\epsilon}_{1,j}^2]$  again use the lower bound on  $\bar{g}_j(Z)$  and decompose

$$\begin{aligned} E_n[\tilde{\epsilon}_{1,j}^2] &= E_n[p_j^2(X) f(m(z) - m(Z)g^2) 1_{D_j^-(Z)} g^2] \\ &\quad + E_n[p_j^2(X) (1 + e^{B_0})^2 f(m_j(Z) - m_j(Z)g^2)] \\ &\quad + E_n[p_j^2(X) (1 + e^{B_0})^2 C_0^2 k b_j^{-1} k^2] \end{aligned}$$

Again on the event  $\bigcap_{m=1}^6 \Omega_{k,m}$  apply Lemma B.2 this is bounded, for some constant  $\tilde{M}_3$  that does not depend on  $k$  by

$$\tilde{M}_3 \frac{2}{k;1} S_k^{-2} \quad (\text{V.4})$$

The result (B.9) follows by collecting (V.1)-(V.4).  $\square$

## B.2 Nonasymptotic Bounds for Residual Estimation

We now provide nonasymptotic bounds that are useful for showing consistent residual estimation. These bounds will be shown under the events in (B.2), (B.5), and (B.7) using the results in Lemmas B.1 and B.2.

**Lemma B.5** (Nonasymptotic Logistic Residual Bound). *Suppose that Assumption 4.1 and the conditions of Lemma B.1 hold. Then, in the event  $\Omega_{k;1} \setminus \Omega_{k;2}$  described on (B.2) there is a constant  $M_{;r}$  that does not depend on  $k$  such that:*

$$\max_{1 \leq j \leq k} E_n[(\hat{\vartheta}_{;j} - U_{;j})^2] \leq M_{;r} \frac{2}{k;1} S_k^{-2} \quad (\text{B.10})$$

*Proof.* Consider each  $j$  separately. By applying the mean value theorem (O.2) and Lemma B.1, we can write

$$\begin{aligned} (\hat{\vartheta}_{;j} - U_{;j})^2 &= p_j(X)^2 Df e^{b_j^0 Z} - e^{j^0 Z} g f e^{b_j^0 Z} - e^{j^0 Z} g \\ &\quad + E_n[p_j(X) Df e^{b_j^0 Z} - e^{j^0 Z} g e^{j^0 Z} u(b_j; j)^0 Z f_j^0 Z] - b_j^0 Z g \\ &\quad + E_n[e^{B_0 + M_0} Df e^{b_j^0 Z} - e^{j^0 Z} g f_j^0 Z] - b_j^0 Z g \end{aligned}$$

So that

$$\begin{aligned} E_n[(\hat{\vartheta}_{;j} - U_{;j})^2] &\leq e^{B_0 + M_0} \frac{2}{k;1} E_n[p_j(X) Df e^{b_j^0 Z} g f_j^0 Z - e^{j^0 Z} g] \\ &\quad + e^{B_0 + M_0} \frac{2}{k;1} S_k^{-2} \end{aligned}$$

$\square$

**Lemma B.6** (Nonasymptotic Linear Residual Bound). *Suppose that Assumptions 4.1 and 4.2 and the conditions of Lemma B.2 hold. Then, in the event  $\bigcap_{m=1}^6 \Omega_{k;m}$ , there is a constant  $M_{;r}$  that does not depend on  $k$  such that*

$$\max_{1 \leq j \leq k} E_n[(\hat{\vartheta}_{;j} - U_{;j})^2] \leq M_{;r} \frac{2}{k;1} S_k^{-2} \quad (\text{B.11})$$

*Proof.* Recall that  $\hat{\vartheta}_{;j} = p_j(X) D e^{b_j^0 Z} (Y - b_j^0 Z)$  and  $U_{;j} = p_j(X) D e^{j^0 Z} (Y - j^0 Z)$ . As an intermediary, define  $\hat{U}_{;j} = p_j(X) D e^{b_j^0 Z} (Y - j^0 Z)$ . We will show a bound on the empirical mean square error between

$\hat{U}_{\cdot j}$  and  $\dot{U}_{\cdot j}$  as well as on the empirical mean square error between  $\dot{U}_{\cdot j}$  and  $U_{\cdot j}$ . The bound in (B.11) will then follow from  $(a + b)^2 \leq 2a^2 + 2b^2$ .

First consider  $(\hat{U}_{\cdot j} - \dot{U}_{\cdot j})^2$ :

$$\begin{aligned} E_n[(\hat{U}_{\cdot j} - \dot{U}_{\cdot j})^2] &= E_n[\rho_j^2(X) D e^{2b_j^0 Z} (b_j^0 Z - \dot{b}_j^0 Z)^2] \\ &= E_n[\rho_j^2(X) D e^{2(b_j^0 Z - \dot{b}_j^0 Z)} (b_j^0 Z - \dot{b}_j^0 Z)^2] \\ &\leq e^{2B_0} e^{2M_0} \underbrace{E_n[\rho_j(X) D e^{b_j^0 Z} (b_j^0 Z - \dot{b}_j^0 Z)^2]}_{= D^{Z_{\cdot j}}(b_j^0; \dot{b}_j^0)} \\ &\leq e^{2M_0} e^{B_0} M_{1, k; 1} s_k^{-2} \end{aligned}$$

Where the last empirical expectation is bounded by Lemma B.2. Next, consider  $(\dot{U}_{\cdot j} - U_{\cdot j})^2$ :

$$\begin{aligned} E_n[(\dot{U}_{\cdot j} - U_{\cdot j})^2] &= E_n[\rho_j^2(X) D \tilde{r} e^{b_j^0 Z} e^{-b_j^0 Z} g^2(Y - \dot{b}_j^0 Z) g^2(Y - b_j^0 Z)] \\ &= E_n[\rho_j^2(X) D \tilde{r} e^{-b_j^0 Z} u(b_j^0)^{b_j^0 Z} (b_j^0 Z - \dot{b}_j^0 Z)^2 g^2(Y - \dot{b}_j^0 Z) g^2(Y - b_j^0 Z)] \\ &\leq 2e^{M_0} e^{B_0} C_{0, k; 1}^2 (M_1 s_k^{-k})^2 E_n[\rho_j(X) D e^{b_j^0 Z} (Y - \dot{b}_j^0 Z)^2] \end{aligned}$$

To proceed we assume that  $Z$  contains a constant. That is  $Z = (1; Z_2; \dots; Z_{d_z})$ . However, this is not necessary it just simplifies the proof a bit. We bound the final empirical expectation in the event  $\Omega_{k;5}$ . In this event we can bound

$$E_n[\rho_j(X) D e^{b_j^0 Z} (Y - \dot{b}_j^0 Z)^2] = (E_n - E)[\rho_j(X) D e^{b_j^0 Z} (Y - \dot{b}_j^0 Z)^2] + E[\rho_j(X) D e^{b_j^0 Z} (Y - m_j(X))^2] \leq e^{B_0} (D_0 + D_1):$$

Combining the above, and using the fact that  $s_k^{-k} < 1$  completes the result. □

### B.3 Probability Bounds for the First Stage

In this section we establish that each of the events in (B.2), (B.5), and (B.7) occurs under Assumptions 4.1 and 4.2 with probability approaching one.

**Lemma B.7** (Logistic Score Domination and Penalty Majorization). *Suppose Assumption 4.1 holds and that the penalty parameter  $\lambda_{\cdot j}$  is chosen as described in Section 4.1. Then, for  $n$  sufficiently large, the event  $\Omega_{k;1}$  holds with probability  $1 - o(1)$  where*

$$\Omega_{k;1} = C \max \left( \frac{4kn + 4k}{n^2}; \frac{\tilde{M}_{k;1} s_k; \tilde{c}_n^2 \ln^5(d_z n)}{n} \leq 1; \frac{\tilde{M}_{k;1}^4 \ln^7(d_z kn)}{n} \leq 1; \frac{1}{\ln^2(d_z kn)} \leq 1 \right) \quad (\text{B.12})$$

where  $C; \tilde{M}$  are absolute constants that do not depend on  $k$ . In particular so long as  $\lambda_{\cdot j} \rightarrow 0$  as  $n \rightarrow \infty$ , this shows that  $\Pr(\Omega_{k;1}) = 1 - o(1)$  under the rate conditions of Assumption 4.1.

Moreover, with probability at least  $1 - \frac{5k}{n} - \frac{4k}{n^2}$  there is a constant  $M_2$  that does not depend on  $k$  such that  $\Omega_{k;2}$  holds with

$$s_k^{-k} = \max\{M_2; M_4; M_5; M_6; M_7\} g_{k;1} \frac{\ln(d_z n)}{n} \quad (\text{B.13})$$

where  $M_4; M_5; M_6$  and  $M_7$  are all constants that also do not depend on  $k$  described in Lemma B.8 and Lemmas B.9–B.11. In particular, so long as  $k = n \rightarrow 0$ ,  $\Pr(\Omega_{k,2}) = 1 - o(1)$ .

*Proof.* Collecting the logistic nonasymptotic residual bound from Lemma B.5 and the probability bounds from Lemmas B.9–B.12 we find that, (eventually) with probability at least  $1 - \frac{4k}{n} - \frac{4k}{n^2}$ :

$$\max_{\substack{1 \leq j \leq k \\ 1 \leq l \leq d_z}} E_n[(\hat{\theta}_{j,l} - U_{j,l} Z_l)^2] \leq M_{j,l} C_0^2 \frac{k;1 S_k; \tilde{c}_n^2 \ln^3(d_z n)}{n}. \quad (\text{P.1})$$

where  $M_{j,l}$  is a constant that does not depend on  $k$ . Define the vectors

$$\begin{aligned} W_k &:= (U_{j,1} Z^0; \dots; U_{j,k} Z^0)^0 \in \mathbb{R}^{kd_z} \\ &:= (W_{k,1}^0; \dots; W_{k,k}^0)^0 \\ \mathcal{W}_k &:= (\hat{\theta}_{j,1} Z^0; \dots; \hat{\theta}_{j,k} Z^0)^0 \in \mathbb{R}^{kd_z} \\ &:= (\mathcal{W}_{k,1}^0; \dots; \mathcal{W}_{k,k}^0)^0. \end{aligned}$$

Notice by optimality of  $\hat{\theta}_{j,1}; \dots; \hat{\theta}_{j,k}$  that  $W_k$  is a mean zero vector. Under our assumptions the covariance matrix  $\Sigma_k = \frac{1}{n} \sum_{i=1}^n E[W_k W_k^0]$  exists and is finite. Define the sequences of constants

$$\begin{aligned} \tilde{c}_{j,n}^2 &:= M_{j,l} C_0^2 \frac{k;1 S_k; \tilde{c}_n^2 \ln^5(d_z n)}{n} \\ \tilde{c}_{j,n} &:= \frac{4k}{n} + \frac{4k}{n^2} \end{aligned}$$

Then, by (P.1) we have that with probability at least  $1 - \tilde{c}_{j,n}$

$$\Pr(k E_n[(\mathcal{W}_k - W_k)^2]_{k,l} > \frac{2}{n} \ln^2(d_z n)) \leq \tilde{c}_{j,n}. \quad (\text{P.2})$$

Let  $e_1; \dots; e_n$  be i.i.d normal random variables generated independently of the data as in Section 4.1. Define the scaled random variables and the multiplier bootstrap process

$$\begin{aligned} \mathfrak{S}_{j,n}^e &:= n^{-1/2} \sum_{i=1}^n e_i \mathcal{W}_{k,i} \\ &:= (\mathfrak{S}_{j,n}^{e^0}; \dots; \mathfrak{S}_{j,n}^{e^k})^0 \end{aligned}$$

and let  $\Pr_e$  denote the probability measure with respect to the  $e_j^0$ s conditional on the observed data. Assumption 4.1 implies that the conditions of (E.1) hold for  $Z = W_k$  with  $b$  replaced by  $c_U$  and  $B_n$  replaced by  $B_k = (C_0 C_U)^3 - 1$ . Further, via (P.2) the residual estimation requirement of with  $n$  and  $n$  replaced by  $\tilde{c}_{j,n}$  and  $\tilde{c}_{j,n}$ .

Let  $\mathfrak{q}_{j,k}(\cdot)$  be the  $\alpha$ -quantile of  $k \mathfrak{S}_{j,k}^{e^0}$  conditional on the data  $Z_i$  and the estimates  $\hat{Z}_i$ . Theorem E.4 then shows that there is a finite constant depending only on  $c_U$  such that

$$\max_{1 \leq j \leq k} \sup_{\alpha \in (0,1)} \Pr(k S_{j,k} - \mathfrak{q}_{j,k}(\alpha)) \leq C \max_{\tilde{c}_{j,n}; \tilde{c}_{j,n}} \left( \frac{B_k^4 \ln^7(k d_z n)}{n} \right)^{1-\alpha} \frac{1}{\ln^2(k d_z n)}.$$

This gives the first claim of Lemma B.7 by construction of  $\hat{Z}_j$ .

The second claim follows Lemma E.1. For this second claim we will consider the marginal convergence of each  $U_{\cdot j} Z$  as opposed to their joint convergence (the convergence of  $W_k$ ). First, notice that conditional on the data, the random vector  $E_n[e^{\theta_{\cdot j} Z}]$  is centered gaussian in  $\mathbb{R}^{d_z}$ . Lemma E.1 then shows that

$$\varphi_{\cdot j}(\cdot) \leq (2 + \frac{1}{2}) \frac{\sqrt{\ln(d_z)}}{n} \max_{1 \leq l \leq d_z} E_n[\theta_{\cdot j}^2 Z_l^2].$$

Furthermore, with probability at least  $1 - \frac{1}{n}$  we have that, for all  $j = 1, \dots, k$ :

$$\max_{1 \leq l \leq d_z} E_n[\theta_{\cdot j}^2 Z_l^2] \leq C_0^2 E_n[\theta_{\cdot j}^2] + 2C_0^2 (E_n[U_{\cdot j}^2] + E_n[(\theta_{\cdot j}^2 - U_{\cdot j})^2]) \leq 4C_0^2 \frac{2}{k} C_U^2 + \frac{2}{n} \ln^2(d_z n)$$

Under the rate conditions of Assumption 4.1,  $\frac{2}{n} \ln^2(d_z n)$  will eventually be smaller than 1 and so the claim in (B.13) holds with  $M_2 = 8C_0^2 C_U^2 - 1$ .  $\square$

**Lemma B.8** (Linear Score Domination and Penalty Majorization). *Suppose Assumptions 4.1 and 4.2 holds and that the penalty parameters  $\lambda_{\cdot j}$  and  $\gamma_{\cdot j}$  are chosen as described in Section 4.1. Then, for  $n$  sufficiently large, the event  $\Omega_{k,3}$  holds with probability  $1 - o(1)$  where:*

$$\frac{1}{n} = C \max \left( \frac{4kn + 4k}{n^2}; \frac{\tilde{M}_{k,1}^2 s_{k,1}^2 \tilde{c}_n^2 \ln^5(d_z n)}{n}; \frac{\tilde{M}_{k,1}^4 \ln^7(d_z kn)}{n}; \frac{1}{\ln^2(d_z kn)} \right) \quad (B.14)$$

where  $C, \tilde{M}$  are absolute constants that do not depend on  $k$ . In particular so long as  $\frac{1}{n} \rightarrow 0$  as  $n \rightarrow \infty$ , this shows that  $\Pr(\Omega_{k,3}) = 1 - o(1)$  under Assumption 4.1.

Moreover, with probability at least  $1 - \frac{5k}{n} - \frac{4k}{n^2}$  there is a constant  $M_4$  that does not depend on  $k$  such that  $\Omega_{k,4}$  holds with

$$\frac{1}{n} \leq \max \{M_2; M_4; M_5; M_6; M_7\} \frac{\sqrt{\ln(d_z n)}}{n} \quad (B.15)$$

where  $M_2; M_5; M_6$  and  $M_7$  are all constants that also do not depend on  $k$  described in Lemma B.7 and Lemmas B.9-B.11. In particular, so long as  $k/n \rightarrow 0$ ,  $\Pr(\Omega_{k,4}) = 1 - o(1)$ .

*Proof.* Apply the same steps as the proof of Lemma B.7 with

$$\begin{aligned} \frac{2}{n} &= M_{k,1} C_0^2 \frac{2}{k} s_{k,1}^2 \tilde{c}_n^2 \ln^5(d_z n) \\ \frac{1}{n} &= \frac{4}{n} + \frac{4}{n^2} \end{aligned}$$

$\square$

**Lemma B.9** (Probabilistic Bound on  $\Omega_{k,5}$ ). *Let  $\tilde{\Sigma}_{\cdot j}$  and  $\Sigma_{\cdot j} = E \tilde{\Sigma}_{\cdot j}$  be as in (B.4). Under Assumptions 4.1 and 4.2 if*

$$\frac{1}{n} \leq \frac{4}{k} (D_0^2 + D_0 D_1) C_0^2 \frac{h}{D_0^2 \log(d_z n) = n} + D_0 D_1 \frac{1}{\log(d_z n) = n}$$

Then  $\Pr(\Omega_{k,5}) \geq 1 - 2k^2/n$ . In particular, there is a constant  $M_5$  that does not depend on  $k$ , such that if  $\frac{1}{n} \leq \frac{1}{M_5 k}$  and  $k^2 \rightarrow 0$  as  $n \rightarrow \infty$  then under the conditions of Assumption 4.2,  $\Pr(\Omega_{k,5}) = 1 - o(1)$ .

*Proof.* We show that this happens with probability  $1 - 2^{-2}$  for each  $j = 1; \dots; k$ . For any  $l; h = 1; \dots; d_z$ , the variable

$$p_j(X) e^{j^{\circ} Z D f Y} \bar{m}_j(Z) g^2 Z_l Z_h$$

is the product of  $p_k(X) e^{j^{\circ} Z Z_l Z_h}$ , which is bounded in absolute value by  $C_{0;7}^2 e^{B_0}$ , and  $D f Y \bar{m}_j(Z) g$ , which is uniformly sub-gaussian conditional on  $Z$ . By Lemma E.7 we have:

$$\mathbb{E} \left[ j(\tilde{\Sigma}_{\cdot j})_{lh} - (\tilde{\Sigma}_{\cdot j})_{lh} \right]^k \leq \frac{k!}{2} (2 C_{0;7}^2 e^{B_0} D_0^2)^k \leq 2 C_{0;7}^2 e^{B_0} D_0 D_1^2; \quad k = 2; 3; \dots$$

Apply the above and Lemma E.5 with  $t = \log(d_z^2) = n$  to obtain

$$\Pr \left[ j(\tilde{\Sigma}_{\cdot j})_{lh} - (\tilde{\Sigma}_{\cdot j})_{lh} > 2 e^{B_0} C_{0;7}^2 D_0^2 t + 2 e^{B_0} C_{0;7}^2 D_0 D_1 \sqrt{2t} \right] \leq 2^{-2} = d_z^2;$$

A union bound completes the argument.  $\square$

**Lemma B.10** (Probabilistic Bound on  $\Omega_{k;6}$ ). *Let  $\tilde{\Sigma}_{\cdot j}$  and  $\Sigma_{\cdot j} = \mathbb{E} \tilde{\Sigma}_{\cdot j}$  be as in (B.4). Under Assumptions 4.1 and 4.2 if*

$$C_{0;7}^2 \sqrt{2(e^{B_0} + 1) C_0 \log(d_z^2) = n};$$

then  $\Pr(\Omega_{k;6}) \geq 1 - 2k^{-2}$ . In particular, there is a constant  $M_6$  that does not depend on  $k$ , such that if  $C_{0;7}^2 M_6 \sqrt{\log(d_z^2) = n}$  and  $k^{-2} \rightarrow 0$  as  $n \rightarrow \infty$  then under the conditions of Assumption 4.2,  $\Pr(\Omega_{k;6}) = 1 - o(1)$ .

*Proof.* Consider each  $j$  separately. For any  $l; h = 1; \dots; d_z$ , note  $j(\tilde{\Sigma}_{\cdot j})_{lh} = j p_j(X) D e^{j^{\circ} Z Z_l Z_h} C_{0;7}^2 e^{B_0}$  so that  $(\tilde{\Sigma}_{\cdot j})_{lh} - (\Sigma_{\cdot j})_{lh}$  is mean zero and bounded in absolute values by  $2 C_{0;7}^2 e^{B_0}$ . Applying Lemma E.3 with  $C_{0;7}^2 e^{B_0} \sqrt{\log(d_z^2) = n}$  yields:

$$\Pr \left[ j(\tilde{\Sigma}_{\cdot j})_{lh} - (\Sigma_{\cdot j})_{lh} > 2 e^{B_0} C_{0;7}^2 \sqrt{\log(d_z^2) = n} \right] \leq 2^{-2} = d_z^2;$$

A union bound completes the argument.  $\square$

**Lemma B.11** (Probabilistic Bound on  $\Omega_{k;7}$ ). *Let  $\tilde{\Sigma}_{\cdot j}^1$  and  $\Sigma_{\cdot j}^1 = \mathbb{E} \tilde{\Sigma}_{\cdot j}^1$  be as in (B.7). Under Assumptions 4.1 and 4.2 if*

$$C_{0;7}^2 \sqrt{4(D_0^2 + D_1^2)^{1-2} e^{B_0} C_0^2 \log(d_z^2) = n};$$

then  $\Pr(\Omega_{k;7}) \geq 1 - 2k^{-2}$ . In particular, there is a constant  $M_7$  that does not depend on  $k$  such that if  $C_{0;7}^2 M_7 \sqrt{\log(d_z^2) = n}$  and  $k^{-2} \rightarrow 0$  as  $n \rightarrow \infty$  then, under the conditions of Assumptions 4.1 and 4.2,  $\Pr(\Omega_{k;7}) = 1 - o(1)$ .

*Proof.* We deal with each  $j$  term separately. The variables  $p_j(X) e^{j^{\circ} Z j Y} \bar{m}_j(Z) j Z_l Z_h$  are uniformly sub-gaussian conditional on  $Z$  because  $j p_j(X) e^{j^{\circ} Z Z_l Z_h} C_{0;7}^2 e^{B_0} C_0^2$  and  $D j Y \bar{m}_j(Z) j$  is uniformly sub-gaussian. Applying Lemma E.4 for  $C_{0;7}^2 e^{B_0} C_{0;7}^2 C_0^2 \sqrt{8(D_0^2 + D_1^2)^2 \log(d_z^2) = n}$  yields

$$\Pr \left[ j(\tilde{\Sigma}_{\cdot j}^1)_{lh} - (\Sigma_{\cdot j}^1)_{lh} > 2 e^{B_0} C_{0;7}^2 \sqrt{\log(d_z^2) = n} \right] \leq 2^{-2} = d_z^2;$$

A union bound completes the argument.  $\square$

#### B.4 Probability Bounds for Residual Estimation

For showing consistent residual estimation, we employ the following two lemmas.

**Lemma B.12** (Deterministic Logistic Score Domination). *Under Assumption 4.1 let*

$$\bar{c}_k = \frac{C_0}{k} \frac{e^{B_0} + 1}{2} \frac{1}{\ln(d_Z) = n}.$$

Then if for all  $j = 1, \dots, k$  we let  $\bar{c}_j = \bar{c}_k$ ,  $\Pr(\Omega_{k,1} \setminus \Omega_{k,2}) \leq 2k \bar{c}_k$ . In particular, there is a constant  $M_8^p$  that does not depend on  $k$  such that if  $\bar{c}_k \leq \frac{M_8^p}{k}$ ,  $\Pr(\Omega_{k,1} \setminus \Omega_{k,2}) \leq 2k n^p$ .

*Proof.* Let us recall that

$$k S_{j,k} = \max_{l=1, \dots, d_Z} \int E_n[\rho_j(X) f(D) e^{jZ} + (1-D)gZ_l] j.$$

Notice for each  $1 \leq l \leq d_Z$ ,  $S_{j,l} = \rho_j(X) f(D) e^{jZ} + (1-D)gZ_l$  is bounded in absolute value by  $C_0 \frac{e^{B_0} + 1}{k}$  and is mean zero by optimality of  $\bar{c}_j$ . For  $\bar{c}_k \leq \frac{C_0}{k} \frac{e^{B_0} + 1}{\ln(d_Z) = n}$  apply Lemma E.3 to see the result.  $\square$

**Lemma B.13** (Deterministic Linear Score Domination). *Under Assumptions 4.1 and 4.2 let*

$$\bar{c}_k = \frac{C_0}{k} \frac{e^{B_0} + 1}{8(D_0^2 + D_1^2)} \frac{1}{\ln(d_Z) = n}.$$

Then if for all  $j = 1, \dots, k$  we let  $\bar{c}_j = \bar{c}_k$ ,  $\Pr(\Omega_{k,3} \setminus \Omega_{k,4}) \leq 2k \bar{c}_k$ . In particular, there is a constant  $M_9^p$  that does not depend on  $k$  such that if  $\bar{c}_k \leq \frac{M_9^p}{k}$ ,  $\Pr(\Omega_{k,3} \setminus \Omega_{k,4}) \leq 2k n^p$ .

*Proof.* Notice  $S_{j,l} = \rho_j(X) f(D) e^{jZ} f_Y(\mathbf{m}_j(Z)) gZ_l$  for  $l = 1, \dots, p$ . By optimality of  $\bar{c}_j$ ,  $S_{j,l}$  is mean zero. Under Assumption 4.1,  $j S_{j,l} \leq e^{B_0} C_0 j f_Y(\mathbf{m}_j(Z)) g$  so by Assumption 4.2 the variables  $S_{j,l}$  are uniformly sub-gaussian conditional on  $Z$  in the following sense:

$$\max_{l=1, \dots, p} \tilde{D}_0^2 E[\exp(S_{j,l}^2 = \tilde{D}_0^2) - 1] \leq \tilde{D}_1^2$$

for  $\tilde{D}_0 = \frac{C_0}{k} D_0 e^{B_0}$  and  $\tilde{D}_1 = \frac{C_0}{k} D_1 e^{B_0}$ . Apply Lemma E.4 for  $\bar{c}_k$  defined above in the statement of Lemma B.13 and union bound to obtain the result.  $\square$

## C Additional Second Stage Results

**Theorem C.1** (Integrated Rate of Convergence). *Assume that Condition 1 and Assumption 5.1 hold. In addition suppose that  $\frac{2}{k} \log k = n^{-1} \rightarrow 0$  and  $c_k \rightarrow 0$ . Then if either the propensity score or outcome regression model are correctly specified:*

$$k \mathbb{E} \|g_k - g_0\|_{L_2} = (E[(g(x) - g_0(x))^2])^{1/2} \cdot \frac{1}{p} \frac{1}{k=n} + c_k \quad (\text{C.1})$$

*Proof.* We begin with a matrix law of large numbers from Rudelson (1999), which is used to show  $\mathbb{E} \|g_k - g_0\|_{L_2} \leq \frac{1}{p} \frac{1}{k=n} + c_k$ .



**Lemma C.1** (Rudelson’s LLN for Matrices). *Let  $Q_1; \dots; Q_n$  be a sequence of independent, symmetric, non-negative  $k \times k$  matrix valued random variables with  $k \geq 2$  such that  $Q = E[E_n Q_i]$  and  $kQ_i k \leq M$  a.s. Then for  $\mathfrak{Q} = E_n[Q_i]$ ,*

$$\Delta := E k \mathfrak{Q} - Q k \leq \frac{M \log k}{n} + \sqrt{\frac{M k Q k \log k}{n}}.$$

*In particular if  $Q_i = p_i p_i^\top$  with  $k p_i k \leq c_k$  almost surely, then*

$$\Delta := E k \mathfrak{Q} - Q k \leq \frac{c_k \log k}{n} + \sqrt{\frac{c_k^2 k Q k \log k}{n}}.$$

Now, to prove Theorem C.1 we have that:

$$\begin{aligned} k \mathfrak{g}_k - g_0 k_{L,2} &= k p^k(x)^\top \mathfrak{b}^k - p^k(x)^\top k_{L,2} + k p^k(x)^\top k_{L,2} - g k_{L,2} \\ &= k p^k(x)^\top \mathfrak{b}^k - p^k(x)^\top k_{L,2} + c_k \end{aligned}$$

where under the normalization  $Q = I_k$  we have that

$$k p^\top \mathfrak{b} - p^\top k_{L,2} = k \mathfrak{b} - k$$

Further,

$$\begin{aligned} k \mathfrak{b}^k - k k &= k \mathfrak{Q}^{-1} E[p^k(x) (p - \bar{Y})] k + k \mathfrak{Q}^{-1} E_n[p^k(x) (k + r_k)] k \\ &= k \mathfrak{Q}^{-1} E[p^k(x) (p - \bar{Y})] k + k \mathfrak{Q}^{-1} E_n[p^k(x) - k] k + k \mathfrak{Q}^{-1} E_n[p^k(x) r_k] k \end{aligned}$$

By the matrix LLN (Lemma C.1) we have that since  $\frac{2}{k} \log k = n^{-1} \rightarrow 0$ ,  $k \mathfrak{Q} - Q k \rightarrow 0$ . This means that with probability approaching one all eigenvalues of  $\mathfrak{Q}$  are bounded away from zero, in particular they are larger than  $1/2$ . So w.p.a 1

$$k E[p^k(x) (p - \bar{Y})] k + k E_n[p^k(x) - k] k + k E_n[p^k(x) r_k] k$$

Under Condition 1 the first term is  $O_p(\frac{1}{k=n})$ . By equation (A.48) in Belloni et al. (2015) the third term is bounded in probability by  $c_k$ . For the second term apply the third condition in Assumption 5.1 to see

$$E k E_n[p^k(x) - k] k^2 = E \sum_{j=1}^k p_j(x)^2 = n^{-2} E_n[p^k(x) p^k(x)^\top = n] \cdot E[p^k(x) p^k(x)^\top = n] = k = n.$$

This gives  $k E_n[p^k(x) - k] k = o_p(\frac{1}{k=n})$  and thus shows (C.1). □

The following lemma is a building block for asymptotic pointwise normality. It establishes conditions under which the coefficient estimator  $\mathfrak{b}$  is asymptotically linear in the sense of Bickel et al. (1993).

**Lemma C.2** (Pointwise Linearization). *Suppose that Condition 1 and Assumption 5.1, hold. In addition assume that  $\frac{2}{k} \log k = n^{-1} \rightarrow 0$ . Then for any  $\delta \geq S^k - 1$ ,*

$$P_{\bar{n}}^{-\delta} ( \mathfrak{b}^k - k ) = O_p[G_n[p^k(x) (k + r_k)] + R_{1n}(\delta)] \tag{C.2}$$



is bounded in probability uniformly over all  $x \in X$ . It will form an important building block in uniform consistency and strong approximation results presented in Theorems 5.2 and 5.3.

**Lemma C.3** (Uniform Linearization). *Suppose that Condition 1 and Assumption 5.1-5.3 hold. Then if either the propensity score model or our outcome regression model is correctly specified:*

$$\rho_{\bar{n}}(x)^\theta(b^k - \eta) = (x)^\theta G_n[\rho^k(x) - (\eta + r_k)] + R_{1n}(x) \quad (\text{C.6})$$

where  $R_{1n}(x)$  describes the design error and satisfies

$$R_{1n}(x) \cdot \rho \leq \frac{\sqrt{\frac{2 \log k}{n}}}{n^{1-m}} \mathbb{P} \overline{\log k} + \frac{\rho}{k} \cdot c_k := \bar{R}_{1n} \quad (\text{C.7})$$

uniformly over  $x \in X$ . Moreover,

$$\rho_{\bar{n}}(x)^\theta(b^k - \eta) = (x)^\theta G_n[\rho^k(x) - \eta] + R_{1n}(x) + R_{2n}(x) \quad (\text{C.8})$$

where  $R_{2n}(x)$  describes the sampling error and satisfies, uniformly over  $x \in X$ :

$$R_{2n}(x) \cdot \rho \leq \mathbb{P} \overline{\log k} \cdot c_k := \bar{R}_{2n} \quad (\text{C.9})$$

*Proof.* As in the proof of Lemma C.2, we decompose

$$\begin{aligned} \rho_{\bar{n}}(x)^\theta(b^k - \eta) &= \rho_{\bar{n}}(x)^\theta \mathbb{E}_n[\rho^k(x) - (\eta - \bar{Y})] \\ &\quad + (x)^\theta G_n[\rho^k(x) - (\eta + r_k)] \\ &\quad + (x)^\theta [\mathbb{E}^{-1} - \mathbb{E}_n] G_n[\rho^k(x) - (\eta + r_k)]; \end{aligned} \quad (\text{C.10})$$

Using Condition 1, the matrix LLN (Lemma C.1), and bounded eigenvalues of the design matrix, we have that:

$$\sup_{x \in X} \rho_{\bar{n}}(x)^\theta \mathbb{E}_n[\rho^k(x) - (\eta - \bar{Y})] = o_p(1):$$

Since this is  $o_p(1)$ , we can simply include this term in  $R_{1n}(x)$ . Now derive bounds on  $R_{1n}(x)$  and  $R_{2n}(x)$ .

**Step 1:** Conditional on the data let

$$T := \{t = (t_1, \dots, t_n) \in \mathbb{R}^n : t_i = (x)^\theta (\mathbb{E}^{-1} - \mathbb{E}_n) \rho^k(x) - \eta; x \in X\}:$$

Define the norm  $\|\cdot\|_{k, n, 2}$  on  $\mathbb{R}^n$  by  $\|t\|_{k, n, 2}^2 = n^{-1} \sum_{i=1}^n t_i^2$ . For an  $\epsilon > 0$  an  $\epsilon$ -net of the normed space  $(T; \|\cdot\|_{k, n, 2})$  is a subset  $T_\epsilon$  of  $T$  such that for every  $t \in T$  there is a point  $t_\epsilon \in T_\epsilon$  such that  $\|t - t_\epsilon\|_{k, n, 2} < \epsilon$ . The covering number  $N(T; \|\cdot\|_{k, n, 2}; \epsilon)$  of  $T$  is the infimum of the cardinality of  $\epsilon$ -nets of  $T$ .

Let  $\epsilon_1, \dots, \epsilon_n$  be independent Rademacher random variables that are independent of the data. Let  $\epsilon = (\epsilon_1, \dots, \epsilon_n)$ . Let  $\mathbb{E}[\cdot]$  denote the expectation with respect to the distribution of  $\epsilon$ . By Dudley's inequality

(Dudley, 1967),

$$E \int_0^1 \sup_{x \in X} (x)^{\theta} [\mathbb{Q}^{-1} I] G_n [ip^k(x) - k] \cdot \int_0^1 \mathbb{Q} \frac{1}{\log N(T; k, k_{n,2}; \cdot) d''} d'';$$

where

$$\begin{aligned} &:= 2 \sup_{t \in T} kt k_{n,2} \\ &= 2 \sup_{x \in X} E_n [(x)^{\theta} (\mathbb{Q}^{-1} I) p^k(x) - k]^2]^{1/2} \\ &= 2 \max_{1 \leq i \leq n} j_{-k;ij} k \mathbb{Q}^{-1} I k k \mathbb{Q} k^{1-2}; \end{aligned}$$

by (A.2). Now, for any  $x \in X$ ,

$$\begin{aligned} &E_n [(x)^{\theta} (\mathbb{Q}^{-1} I) p^k(x) - k - (\bar{x})^{\theta} (\mathbb{Q}^{-1} I) p^k(x) - k]^2]^{1/2} \\ &= \max_{1 \leq i \leq n} j_{-k;ij} k (x) - (\bar{x}) k k \mathbb{Q}^{-1} I k k \mathbb{Q} k^{1-2} \\ &= \frac{1}{k} \max_{1 \leq i \leq n} j_{-k;ij} k \mathbb{Q}^{-1} I k k \mathbb{Q} k^{1-2} k_X - \bar{x} k \end{aligned}$$

So, for some  $C > 0$ ,

$$N(T; k, k_{n,2}; \cdot) \leq \frac{C \frac{1}{k} \max_{1 \leq i \leq n} j_{-k;ij} k \mathbb{Q}^{-1} I k k \mathbb{Q} k^{1-2} \int d_x}{\cdot};$$

This gives us that

$$\int_0^1 \mathbb{Q} \frac{1}{\log(N(T; k, k_{n,2}; \cdot) d'')} d'' \leq \max_{1 \leq i \leq n} j_{-k;ij} k \mathbb{Q}^{-1} I k k \mathbb{Q} k^{1-2} \int_0^1 \mathbb{Q} \frac{1}{d_x \log(C \frac{1}{k} \cdot) d''} d'';$$

By Assumption 5.3 we have that  $E[\max_{1 \leq i \leq n} j_{-k;ij} j | X] \leq P n^{1-m}$  where  $X = (x_1, \dots, x_n)$ . In addition  $\frac{2m-(m-2)}{k} \log k = n - 1$  for  $m > 2$  gives that  $\frac{2}{k} \log k = n - 1 \rightarrow 0$ . So,  $k \mathbb{Q}^{-1} I k \leq P (\frac{2}{k} \log k = n)^{1-2}$  and  $k \mathbb{Q}^{-1} k \leq P$ . Combining this all with  $\log \frac{1}{k} \leq \log k$  implies

$$E \int_0^1 \sup_{x \in X} (x)^{\theta} [\mathbb{Q}^{-1} I] G_n [p^k(x) - k] \cdot \int_0^1 \mathbb{Q} \frac{1}{\log(N(T; k, k_{n,2}; \cdot) d'')} d'' \leq 2 E \int_0^1 \sup_{x \in X} (x)^{\theta} [\mathbb{Q}^{-1} I] G_n [ip^k(x) - k] \cdot \int_0^1 \mathbb{Q} \frac{1}{\log(N(T; k, k_{n,2}; \cdot) d'')} d'' \leq P n^{1-m} \frac{\frac{2}{k} \log^2 k}{n}$$

where the first line is due to symmetrization inequality. This gives us

$$\sup_{x \in X} (x)^{\theta} [\mathbb{Q}^{-1} I] G_n [p^k(x) - k] \leq P n^{1-m} \frac{\frac{2}{k} \log^2 k}{n} \tag{C.11}$$

**Step 2:** Now simply report the results on approximation error from Belloni et al. (2015). Since the

approximation error is the same for all signals  $Y(-_k; \bar{m}_k)$ , there is no Hadamard product to deal with.

$$\sup_{x \in X} (x)^\theta \left[ \frac{1}{n} \int G_n[p^k(x)r_k] \cdot P \right]^{\frac{2 \log k}{n}} \cdot \frac{P_{-k} c_k}{k} \quad (\text{C.12})$$

$$\sup_{x \in X} (x)^\theta G_n[p^k(x)r_k] \cdot P \cdot \frac{P_{-k} c_k}{\log k} \quad (\text{C.13})$$

Looking at (C.10) and combining (C.11)-(C.12) gives the bound on  $R_{1n}(\cdot)$  while (C.13) gives the bound on  $R_{2n}(\cdot)$ .  $\square$

**Theorem C.2** (Validity of Gaussian Bootstrap). *Suppose that the assumptions of Theorem 5.3 hold with  $a_n = \log n$  and the assumptions of Theorem 5.4 hold with  $a_n = O(n^{-b})$  for some  $b > 0$ . In addition, suppose that there exists a sequence  $\frac{\theta}{n}$  obeying  $1 \leq \frac{\theta}{n} \leq k p^k(x) k$  uniformly for all  $x \in X$  such that  $k p^k(x) - p^k(x) k = \frac{\theta}{n} L_n k x - x k$ , where  $\log L_n \leq \log n$ . Let  $N_k^b$  be a bootstrap draw from  $N(0; I_k)$  and  $P^\tau$  be the distribution conditional on the observed data  $(Y_i; D_i; Z_i)_{i=1}^n$ . Then the following approximation holds uniformly in  $\cdot^{-1}(X)$ :*

$$\frac{p^k(x)^\theta \Omega^{1=2}}{\Omega^{1=2} p^k(x)} N_k^b =^d \frac{p^k(x)^\theta \Omega^{1=2}}{k \Omega^{1=2} p^k(x) k} + o_{P^\tau}(\log^{-1} N) \quad (\text{C.14})$$

*Proof.* See Theorem 3.4 in [Semenova and Chernozhukov \(2021\)](#).  $\square$

## D Alternative CV-Type Method for Penalty Parameter Selection

In this section we consider a procedure for penalty parameter selection where we use the pilot penalty parameters described in (4.2) directly, after choosing constants  $c_{-j}$  and  $c_{+j}$  from a (finite) set via cross validation. For each  $j$  we will assume that

$$c_{-j}; c_{+j} \geq \Lambda_n \in [c_n; \bar{c}_n] \quad (\text{D.1})$$

where  $j \Lambda_n j$  can be fairly large (on the order of  $n^2 = k$ ).

### D.1 Theory Overview

Let  $M_5; M_6; M_7; M_8^2; M_9^2$  be constants that do not depend on  $k$  as in Lemmas B.9–B.13. Whenever

$$c_n \frac{\ln^3(d_z n)}{n} \leq \max_{k; \tau} M_5; M_6; M_7; M_8^2; M_9^2 \leq \frac{\ln(d_z n)}{n}, \quad (\text{D.2})$$

we will have that, under Assumption 4.1(i)-(iv) and Assumption 4.2(i)-(iii) the event  $\bigcap_{k=1}^T \Omega_{k; \tau}$  occurs with probability at least  $1 - 10k = n^2$  for the  $2k$  pilot penalty parameters chosen with any values  $c_{-j}; c_{+j} \geq \Lambda_n$  and

$$c_{-k} := \bar{c}_n \frac{\ln^3(d_z n)}{n}.$$

In this event, apply Lemmas B.1 and B.2 to obtain the following finite sample bounds for the parameter estimates

$$\begin{aligned} \max_{1 \leq j \leq k} D_{\cdot j}^z(b_j; \bar{\cdot}_j) & \leq M_0 \frac{S_k \bar{c}_n^2 \ln^3(d_z n)}{n} \quad \text{and} \quad \max_{1 \leq j \leq k} \|b_j - \bar{\cdot}_j\|_{k_1} \leq M_0 S_k \bar{c}_n \frac{\ln^3(d_z n)}{n} \\ \max_{1 \leq j \leq k} D_{\cdot j}^z(b_j; \bar{\cdot}_j; \bar{\cdot}_j) & \leq M_1 \frac{S_k \bar{c}_n^2 \ln^3(d_z n)}{n} \quad \text{and} \quad \max_{1 \leq j \leq k} \|b_j - \bar{\cdot}_j\|_{k_1} \leq M_1 S_k \bar{c}_n \frac{\ln^3(d_z n)}{n} \end{aligned}$$

and Lemma B.3 to obtain the following finite sample bound for the weighted means:

$$\max_{1 \leq j \leq k} |E_n[\rho_j(X)(Y(b_j; \bar{m}_j) - Y(\bar{\cdot}_j; \bar{m}_j))]| \leq M_2 \frac{\bar{c}_n^2 S_k \ln^3(d_z n)}{n} \quad (\text{D.3})$$

$$\max_{1 \leq j \leq k} |E_n[\rho_j^2(X)(Y(b_j; \bar{m}_j) - Y(\bar{\cdot}_j; \bar{m}_j))^2]| \leq M_3 \frac{\bar{c}_{k,1}^2 \bar{c}_n^2 S_k^2 \ln^3(d_z n)}{n} \quad (\text{D.4})$$

Combining (D.2) and (D.3) we can see that Condition 1 can be obtained under Assumption 4.1(i)-(iv), Assumption 4.2(i)-(iii), and the following modified sparsity bounds

$$\frac{k_j \Lambda_n^j}{n^2} \neq 0; \quad \frac{\bar{c}_n^{-1} k_{k,1}}{\ln(d_z n)} \neq 0 \quad \text{and} \quad \frac{\bar{c}_n^2 S_k k^{1=2} \ln^3(d_z n)}{n} \neq 0; \quad (\text{D.5})$$

Similarly combining (D.2) and (D.4), Condition 2 can additionally be obtained by strengthening the rates in (D.5) to include

$$\frac{\bar{c}_{k,1}^2 \bar{c}_n^2 S_k k^2 \ln^3(d_z n)}{n^{(m-1)=m}} \neq 0 \quad (\text{D.6})$$

for  $m > 2$  as in Assumption 5.3. These rates are comparable and in certain cases may be more palatable than those presented in the main text, Assumptions 4.1 and 4.2. They come at the cost of slower rates of convergence for the weighted means as seen by comparing eqs. (D.3)–(D.4) to eqs. (4.9) and (4.12).

## D.2 Practical Implementation

In practice, the constants  $M_5; M_6; M_7; M_8^2; M_9^2$  from Lemmas B.9–B.13 are roughly on the order of  $kZk_1$ . We therefore recommend setting

$$\begin{aligned} \bar{c}_n &= \frac{1}{2 \log^{1=2}(d_z n)} \max_{1 \leq i \leq n} \|k p^k(X_i)\|_{k_1} \max_{1 \leq i \leq n} \|k Z_i\|_{k_1} \\ \bar{c}_n &= \frac{3 \log^{1=2}(d_z n)}{2} \max_{1 \leq i \leq n} \|k p^k(X_i)\|_{k_1} \max_{1 \leq i \leq n} \|k Z_i\|_{k_1} \end{aligned}$$

and letting  $\Lambda_n$  be a set of points evenly spaced between  $\bar{c}_n$  and  $\bar{c}_n$ . The cross validation procedure then can be implemented in the following steps.

1. Split the sample into  $K$  folds.
2. Consider a single pair of values for  $c; \bar{c}$  and designate a fold to hold out.
3. Estimate nuisance model parameters using  $\hat{\cdot}_j^{\text{pilot}}$  and  $\hat{\cdot}_j^{\text{pilot}}$  on the remaining folds.

4. Evaluate the resulting models on held out fold using non-penalized loss functions.
5. Repeat  $K$  times and record average loss over all folds.
6. Choose values of  $c_{\cdot j}$  and  $c_{\cdot j}$  with the lowest average loss.

In practice we find this procedure works well with small  $K$ , around  $K = 5$  and with  $j\Lambda_{nj}$  on the order of about 10-20.

## E High Dimensional Probability Results

### E.1 High Dimensional Central Limit and Bootstrap Theorems

**Lemma E.1** (Gaussian Quantile Bound). *Let  $Y = (Y_1, \dots, Y_p)$  be centered Gaussian in  $\mathbb{R}^p$  with  $\max_{1 \leq j \leq p} \mathbb{E}[Y_j^2] \leq 2$  and  $\rho \geq \frac{2}{\ln(p)}$ . Let  $q^Y(1 - \frac{\rho}{2})$  denote the  $(1 - \frac{\rho}{2})$ -quantile of  $\|Y\|_1$  for  $\rho \in (0, 1)$ . Then  $q^Y(1 - \frac{\rho}{2}) \leq (2 + \frac{\rho}{2}) \frac{1}{\rho} \ln(p)$ .*

*Proof.* See [Chetverikov and Sørensen \(2021\)](#), Lemma D.2. □

Now let  $Z_1, \dots, Z_n$  be independent, mean zero random variables in  $\mathbb{R}^p$ , and denote their scaled average and variance by

$$S_n := \frac{1}{n} \sum_{i=1}^n Z_i \text{ and } \Sigma := \frac{1}{n} \sum_{i=1}^n \mathbb{E}[Z_i Z_i^T];$$

For  $\mathbb{R}^p$  values random variables  $U$  and  $V$ , define the distributional measure of distance

$$(U; V) := \sup_{A \in \mathcal{A}_p} \Pr(U \in A) - \Pr(V \in A)$$

where  $\mathcal{A}_p$  denotes the collection of all hyperrectangles in  $\mathbb{R}^p$ . For any symmetric positive matrix  $M \in \mathbb{R}^{p \times p}$ , write  $N_M := N(\mathbf{0}; M)$ .

**Theorem E.1** (High-Dimensional CLT). *If, for some finite constants  $b > 0$  and  $B_n \leq 1$ ,*

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E}[Z_{ij}^2] \leq b; \frac{1}{n} \sum_{i=1}^n \mathbb{E}[|Z_{ij}|^{2+k}] \leq B_n^k \text{ and } \mathbb{E}[\max_{1 \leq j \leq p} Z_{ij}^4] \leq B_n^4; \tag{E.1}$$

for all  $i \in \{1, \dots, ng\}$  and  $j \in \{1, \dots, pg\}$  and  $k \in \{1, 2, 3, 4\}$ , then there exists a finite constant  $C_b$ , depending only on  $b$ , such that:

$$(S_n; N) \leq C_b \frac{B_n^4 \ln^7(pn)}{n}^{1/6};$$

*Proof.* See [Chernozhukov et al. \(2017\)](#), Proposition 2.1. □

Let  $\hat{Z}_i$  be an estimator of  $Z_i$  and let  $e_1, \dots, e_n$  be i.i.d  $N(0, 1)$  and independent of both the  $Z_i$ 's and  $\hat{Z}_i$ 's. Define  $\hat{S}_n^e := \frac{1}{n} \sum_{i=1}^n e_i \hat{Z}_i$  and let  $\Pr_e$  denote the conditional probability measure computed with respect to the  $e_i$ 's for fixed  $Z_i$ 's and  $\hat{Z}_i$ 's. Also abbreviate

$$\tilde{(\hat{S}_n^e; N)} := \sup_{A \in \mathcal{A}_p} \Pr_e[\hat{S}_n^e \in A] - \Pr(N \in A);$$

**Theorem E.2** (Multiplier Bootstrap for Many Approximate Means). *Let (E.1) hold for some finite constants  $b > 0$  and  $B_n \geq 1$ , and let  $f_{n, \mathcal{G}_N}$  and  $f_{n, \mathcal{G}_N}$  be sequences in  $\mathbb{R}_{++}$  converging to zero such that*

$$\Pr^{\otimes} \max_{j=1, \dots, p} \frac{1}{n} \sum_{i=1}^n (\mathbb{Z}_{ij} - Z_{ij})^2 > \frac{2}{\ln^2(\rho n)} A \quad (E.2)$$

Then, there exists a finite constant  $C_b$  depending only on  $b$  such that with probability at least  $1 - \frac{1}{\ln^2(\rho n)}$ ,

$$\mathbb{P}(\mathbb{S}_n^e; N) \leq C_b \max_{j=1, \dots, p} \frac{B_n \ln^6(\rho n)}{n} \quad ; \quad ;$$

*Proof.* See Belloni et al. (2018), Theorem 2.2 or Chetverikov and Sørensen (2021) Theorem D.2.  $\square$

We now consider a partition of  $Z$  and  $\mathbb{Z}$  into  $k$  subvectors.

$$Z := (Z_1^0; \dots; Z_k^0)^0 \in \mathbb{R}^{d_1; \dots; d_k} \quad \text{and} \quad \mathbb{Z} := (\mathbb{Z}_1^0; \dots; \mathbb{Z}_k^0)^0 \in \mathbb{R}^{d_1; \dots; d_k}$$

where  $\sum_{j=1}^k d_j = p$ . Given such a partition, for any symmetric, positive definite  $M \in \mathbb{R}^{p \times p}$  let  $N_{M,j}$  denote the marginal distribution of the subvector of  $N_M$  corresponding to the indices of partition  $j$ . In other words,  $N_{M,1}$  would denote the marginal distribution of the first  $d_1$  elements of an  $\mathbb{R}^p$  vector with distribution  $N_M$ ,  $N_{M,2}$  would denote the marginal distribution of the next  $d_2$  elements and so on. For each  $j = 1; \dots; k$  define  $q_{M,j}^N: \mathbb{R} \rightarrow \mathbb{R}$  as the (extended) quantile function of  $kN_{M,j}k_1$ ,

$$q_{M,j}^N(\cdot) := \inf \{ t \in \mathbb{R} : \Pr(kN_{M,j}k_1 \leq t) \geq \cdot \}$$

Define  $q_{M,j}^N(\cdot) = +\infty$  if  $\cdot < 0$  and  $-\infty$  if  $\cdot > 1$  so that  $q_{M,j}^N$  is always monotone (strictly) increasing.

**Lemma E.2.** *Let  $M \in \mathbb{R}^{p \times p}$  be symmetric positive definite, let  $U$  be a random variable in  $\mathbb{R}^p$ . Partition  $U$  into  $k$  subvectors,  $U = (U_1^0; \dots; U_k^0)^0 \in \mathbb{R}^{d_1; \dots; d_k}$  where  $d_1 + \dots + d_k = p$ . For each  $j = 1; \dots; k$  let  $q_j$  denote the quantile function of  $kU_jk_1$ . Then for any  $j = 1; \dots; k$ ,*

$$q_{M,j}^N(\cdot - 2(U; N_M)) = q_j(\cdot) = q_{M,j}^N(\cdot + (U; N_M)) \quad \text{for all } \cdot \in (0; 1):$$

*Proof.* Proof is a slight modification of that of Lemma D.3 in Chetverikov and Sørensen (2021). Main idea is to add and subtract a  $kN_{M,j}k_1$  term and use the fact that the approximation is achieved over all hyperrectangles. We show the bound holds for each  $j = 1; \dots; k$ . Without loss of generality, consider  $U_1$ . Let  $N_{M,1}$  denote the marginal distribution of the first  $d_1$  elements of a  $\mathbb{R}^p$  vector with distribution  $N_M$ .

$$\begin{aligned} \Pr(kU_1k_1 \leq t) &= \Pr(kN_{M,1}k_1 \leq t) + \Pr(kU_1k_1 \leq t) - \Pr(kN_{M,1}k_1 \leq t) \\ &= \Pr(kN_{M,1}k_1 \leq t) + \Pr(U \in [t; t]^{d_1} \in \mathbb{R}^{d_1}) - \Pr(N_M \in [t; t]^{d_1} \in \mathbb{R}^{d_1}) \\ &= \Pr(kN_{M,1}k_1 \leq t) + (U; N_M) \end{aligned}$$

for any  $t \in \mathbb{R}$ . A similar construction will give that

$$\Pr(kU_1k_1 \leq t) = \Pr(kN_{M,1}k_1 \leq t) + (U; N_M):$$



Substituting  $t = q_{M,1}^N(2(U; N_M))$  into the upper bound on  $\Pr(kU_1k_1 \leq t)$  gives the lower bound statement, while  $t = q_{M,1}^N(2(U; N_M))$  and using the lower bound on  $\Pr(kU_1k_1 \leq t)$  gives the upper bound statement.  $\square$

As with  $Z$  partition  $S_n$  and  $\mathfrak{S}_n^e$  into

$$S_n = (S_{n,1}^0; \dots; S_{n,k}^0)^0 \in \mathbb{R}^{d_1; \dots; d_k} \quad \text{and} \quad \mathfrak{S}_n^e = (\mathfrak{S}_{n,1}^{e0}; \dots; \mathfrak{S}_{n,k}^{e0})^0 \in \mathbb{R}^{d_1; \dots; d_k};$$

For each  $j = 1; \dots; k$  define  $q_{n,j}(\cdot)$  as the  $\cdot$ -quantile of  $kS_{n,j}k_1$

$$q_{n,j}(\cdot) := \inf\{t \in \mathbb{R} : \Pr(kS_{n,j}k_1 \leq t) \geq \cdot\} \text{ for } \cdot \in (0;1);$$

Let  $\mathfrak{q}_{n,j}(\cdot)$  be the  $\cdot$ -quantile of  $k\mathfrak{S}_{n,j}^e k_1$ , computed conditionally on  $X_i$  and  $\mathfrak{X}_i$ 's,

$$\mathfrak{q}_{n,j}(\cdot) := \inf\{t \in \mathbb{R} : \Pr_e(k\mathfrak{S}_{n,j}^e k_1 \leq t) \geq \cdot\} \text{ for } \cdot \in (0;1);$$

**Theorem E.3** (Quantile Comparasion). *If (E.1) holds for some finite constants  $b > 0$  and  $B_n \geq 1$ , and*

$$n := 2C_b \frac{B_n^4 \ln^7(pn)}{n} \stackrel{!}{=} 1-6$$

*denotes the upper bound in Theorem E.1 multiplied by two, then for all  $j = 1; \dots; k$*

$$q_{n,j}^N(1 - \frac{1}{n}) \leq q_{n,j}(1 - \frac{1}{n}) \leq q_{n,j}^N(1 + \frac{1}{n}) \text{ for all } \cdot \in (0;1);$$

*If, in addition, (E.2) holds for some sequences  $f_n, g_n$  and  $\tilde{f}_n, \tilde{g}_n$  converging to zero, and*

$$\frac{1}{n} \leq 2C_b \max\left\{ \frac{8}{n}, \frac{B_n^4 \ln^6(pn)}{n} \right\} \stackrel{!}{=} 1-6 \stackrel{9}{=}$$

*denotes the upper bound in Theorem E.2 multiplied by two, then with probability at least  $1 - \frac{1}{n} = 1 - \ln^2(pn)$  we have for all  $j = 1; \dots; k$ ,*

$$q_{n,j}^N(1 - \frac{1}{n}) \leq \mathfrak{q}_{n,j}(1 - \frac{1}{n}) \leq q_{n,j}^N(1 + \frac{1}{n}) \text{ for all } \cdot \in (0;1);$$

*Proof.* From Lemma E.2 with  $U = S_n$  we obtain

$$q_{n,j}^N(1 - 2(S_n; N)) \leq q_{n,j}(1 - \frac{1}{n}) \leq q_{n,j}^N(1 + (S_n; N));$$

The first chain of inequalities then follows from  $2(S_n; N) \leq \frac{1}{n}$  by Theorem E.1.

For the second claim, apply Lemma E.2 with  $U = \mathfrak{S}_n^e$  and condition on the  $Z_i$ 's and  $\mathfrak{Z}_i$ 's obtain

$$q_{n,j}^N(1 - 2(\mathfrak{S}_n^e; N)) \leq \mathfrak{q}_{n,j}(1 - \frac{1}{n}) \leq q_{n,j}^N(1 + (\mathfrak{S}_n^e; N));$$

The second chain of inequalities then follows on the event  $2(\mathfrak{S}_n^e; N) \leq \frac{1}{n}$ , which by Theorem E.2 happens

with probability at least  $1 - \frac{1}{n} = 1 - \ln^{-2}(pn)$ . □

**Theorem E.4** (Multiplier Bootstrap Consistency). *Let (E.1) and (E.2) hold for some constants  $b > 0$  and  $B_n \rightarrow 1$  and some sequences  $f_{n, \mathcal{G}_N}$  and  $f_{n, \mathcal{G}_N}$  in  $\mathbb{R}_{++}$  converging to zero. Then, there exists a finite constant  $C_b$ , depending only on  $b$ , such that*

$$\max_{1 \leq j \leq k} \sup_{z \in (0,1)} \Pr(kS_{n,j}k_1 \geq \varphi_{n,j}(1-z)) \leq C_b \max_{1 \leq j \leq k} \left( \frac{B_n^4 \ln^7(pn)}{n} + \frac{1}{\ln^2(pn)} \right)$$

*Proof.* By Theorem E.1 and Theorem E.3,

$$\begin{aligned} \Pr(kS_{n,j}k_1 \geq \varphi_{n,j}(1-z)) &\leq \Pr(kS_{n,j}k_1 \geq \varphi_{n,j}^N(1-z) + \frac{1}{n}) + \frac{1}{n} \\ &\leq \Pr(kN_{n,j}k_1 \geq \varphi_{n,j}^N(1-z) + \frac{1}{n}) + \frac{1}{n} \\ &\leq \frac{1}{n} + \frac{1}{n} + \frac{1}{n} + \frac{1}{\ln^2(pn)} \end{aligned}$$

Where the second inequality is making use of the same rectangle argument as before. A parallel argument shows that

$$\Pr(kS_{n,j}k_1 \leq \varphi_{n,j}(1-z)) \leq \frac{1}{n} + \frac{1}{n} + \frac{1}{n} + \frac{1}{\ln^2(pn)}$$

Combining these two inequalities gives the result. □

### E.2 Concentration and Tail Bounds

We make use of the following concentration and tail bounds. Lemmas E.3–E.7 can be found in [Bühlmann and van de Geer \(2011\)](#). The proof of Lemma E.8 is trivial but provided here.

**Lemma E.3.** *Let  $(Y_1, \dots, Y_n)$  be independent random variables such that  $E[Y_i] = 0$  for  $i = 1, \dots, n$  and  $\max_{i=1, \dots, n} |Y_{ij}| \leq c_0$  for some constant  $c_0$ . Then, for any  $t > 0$ ,*

$$\Pr\left(\frac{1}{n} \sum_{i=1}^n Y_i > t\right) \leq 2 \exp\left(-\frac{nt^2}{2c_0^2}\right)$$

**Lemma E.4.** *Let  $(Y_1, \dots, Y_n)$  be independent random variables such that  $E[Y_i] = 0$  for  $i = 1, \dots, n$  and  $(Y_1, \dots, Y_n)$  are uniformly sub-gaussian:  $\max_{1 \leq i \leq n} c_1^2 E[\exp(Y_i^2/c_1^2) - 1] \leq c_2^2$  for some constants  $(c_1, c_2)$ . Then for any  $t > 0$ ,*

$$\Pr\left(\frac{1}{n} \sum_{i=1}^n Y_i > t\right) \leq 2 \exp\left(-\frac{nt^2}{8(c_1^2 + c_2^2)}\right)$$

**Lemma E.5.** *Let  $(Y_1, \dots, Y_n)$  be independent variables such that  $E[Y_i] = 0$  for  $i = 1, \dots, n$  and*

$$\frac{1}{n} \sum_{i=1}^n E[jY_{ij}^k] \leq \frac{k!}{2} c_3^k c_4^2, \quad k = 2, 3, \dots$$

for some constants  $(c_3; c_4)$ . Then, for any  $t > 0$ ,

$$\Pr \left\{ \frac{1}{n} \sum_{i=1}^n Y_i > c_3 t + c_4 \frac{\rho_{-}}{2t} \right\} \leq 2 \exp(-nt):$$

**Lemma E.6.** Suppose that  $Y$  is sub-gaussian:  $c_1^2 E[\exp(Y^2=c_1^2) - 1] \leq c_2^2$  for some constants  $(c_1; c_2)$ . Then

$$E[|Y|^k] \leq \Gamma \left( \frac{k}{2} + 1 \right) (c_1^2 + c_2^2) c_1^{k-2}; \quad k = 2; 3; \dots;$$

**Lemma E.7.** Suppose that  $X$  is bounded,  $|X| \leq c_0$ , and  $Y$  is sub-gaussian,  $c_2^2 E[\exp(Y^2=c_1^2) - 1] \leq c_2^2$  for some constants  $(c_1; c_2)$ . Then  $Z = XY^2$  satisfies

$$E[|Z|^k] \leq E[|Z|^k] \leq \frac{k!}{2} c_3^k c_4^2; \quad k = 2; 3; \dots;$$

for  $c_3 = 2c_0c_1^2$  and  $c_4 = 2c_0c_1c_2$ .

**Lemma E.8.** Suppose that  $Y$  is sub-gaussian in the following sense, there exist positive constants  $c_0; c_1 > 0$  such that  $c_0^2 E[\exp(Y^2=c_0^2) - 1] \leq c_1^2$ . Then

$$E[|Y|] \leq c_1^2 = c_0 + c_0.$$

*Proof.* Using the fact that  $e^{x^2} > |x|$  gives

$$\begin{aligned} c_0^2 E[\exp(Y^2=c_0^2) - 1] \leq c_1^2 &\Rightarrow E[\exp(Y^2=c_0^2)] \leq c_1^2 + 1 \\ &\Rightarrow E[|Y=c_0|] \leq c_1^2 + 1 \\ &\Rightarrow E[|Y|] \leq c_1^2 = c_0 + c_0 \end{aligned}$$

□

## F Additional Details on Empirical Application

As mentioned in the setup, to avoid outlier contamination we drop the top 3% and bottom 3% of birthweights by maternal age. We also drop ages for which there are fewer than 10 smoker or non smoker observations. The result is a dataset with 4107 (of an initial 4602) observations on the outcome variable, birthweight. In addition to the 21 control variables ( $Z$ ) available in the dataset, we further generate an additional 11 interaction/higher order variables that we believe may be useful in controlling for confounding. Table F.1 provides a final summary of the data after our cleaning process. The generated variables represent the bottom 11 variables in Table F.1.<sup>1</sup>

In conducting analysis, we found it quite helpful to the stability of the final model assisted estimator to do some light trimming of the estimated propensity score and outcome regression models. In particular we trim the estimated propensity score(s) to be between 0.01 and 0.99 and trim the estimated mean regression models so that they take a value no more than roughly 12.5% higher or lower than the maximum or minimum value of  $Y$  observed in the data.

<sup>1</sup>This table is generated using the wonderful stargazer package in R (Hlavac, 2022).

Table F.1: Summary of Data used in Empirical Exercise

Statistic	N	Mean	St. Dev.	Min	Max
bweight	4,107	3,384.354	447.616	1,544	4,668
mmarried	4,107	0.708	0.455	0	1
mhispanic	4,107	0.034	0.181	0	1
fhispanic	4,107	0.038	0.192	0	1
foreign	4,107	0.054	0.226	0	1
alcohol	4,107	0.031	0.174	0	1
deadkids	4,107	0.252	0.434	0	1
mage	4,107	26.125	5.025	16	36
medu	4,107	12.703	2.470	0	17
fage	4,107	27.000	9.022	0	60
fedu	4,107	12.324	3.624	0	17
nprenatal	4,107	10.822	3.613	0	40
monthslb	4,107	21.938	30.255	0	207
order	4,107	1.858	1.056	0	12
msmoke	4,107	0.390	0.890	0	3
mbsmoke	4,107	0.183	0.386	0	1
mrace	4,107	0.847	0.360	0	1
frace	4,107	0.822	0.382	0	1
prenatal	4,107	1.204	0.507	0	3
birthmonth	4,107	6.556	3.352	1	12
lbweight	4,107	0.025	0.155	0	1
fbaby	4,107	0.443	0.497	0	1
prenatal1	4,107	0.803	0.398	0	1
mbsmoke *alcohol	4,107	0.017	0.128	0	1
medu *fedu	4,107	161.518	64.291	0	289
mage *fage	4,107	730.422	328.522	0	2,088
msmoke^2	4,107	0.944	2.422	0	9
msmoke *alcohol	4,107	0.037	0.302	0	3
mage^2	4,107	707.741	262.383	256	1,296
mage *mmarried	4,107	19.570	13.090	0	36
mage *medu	4,107	336.911	108.588	0	612
mage *fedu	4,107	328.405	128.438	0	612
monthslb^2	4,107	1,396.431	3,509.883	0	42,849
msmoke *monthslb^2	4,107	750.703	4,407.398	0	112,908

Because the control variables are all of different magnitudes, it is common to do some normalization before estimating the  $\ell_1$ -regularized propensity score and outcome regression models so that all variables are “punished” equally by the penalty. We normalize our data by scaling each variable to take on values between zero and one.

In addition to the results presented in Section 8 we present some additional specifications below. Figure F.1 presents results from using a local constant regression with 3 knots in the first stage while Figure F.2. Both show that the results of analysis generally hold up in a variety of specifications, namely that the effect of maternal smoking on birthweight is negative and increasing in magnitude with age. The model-assisted estimates typically produce values that are more in line with previous work, though the shape of CATE estimates that use standard loss functions in the first stage are more stable to second stage basis.

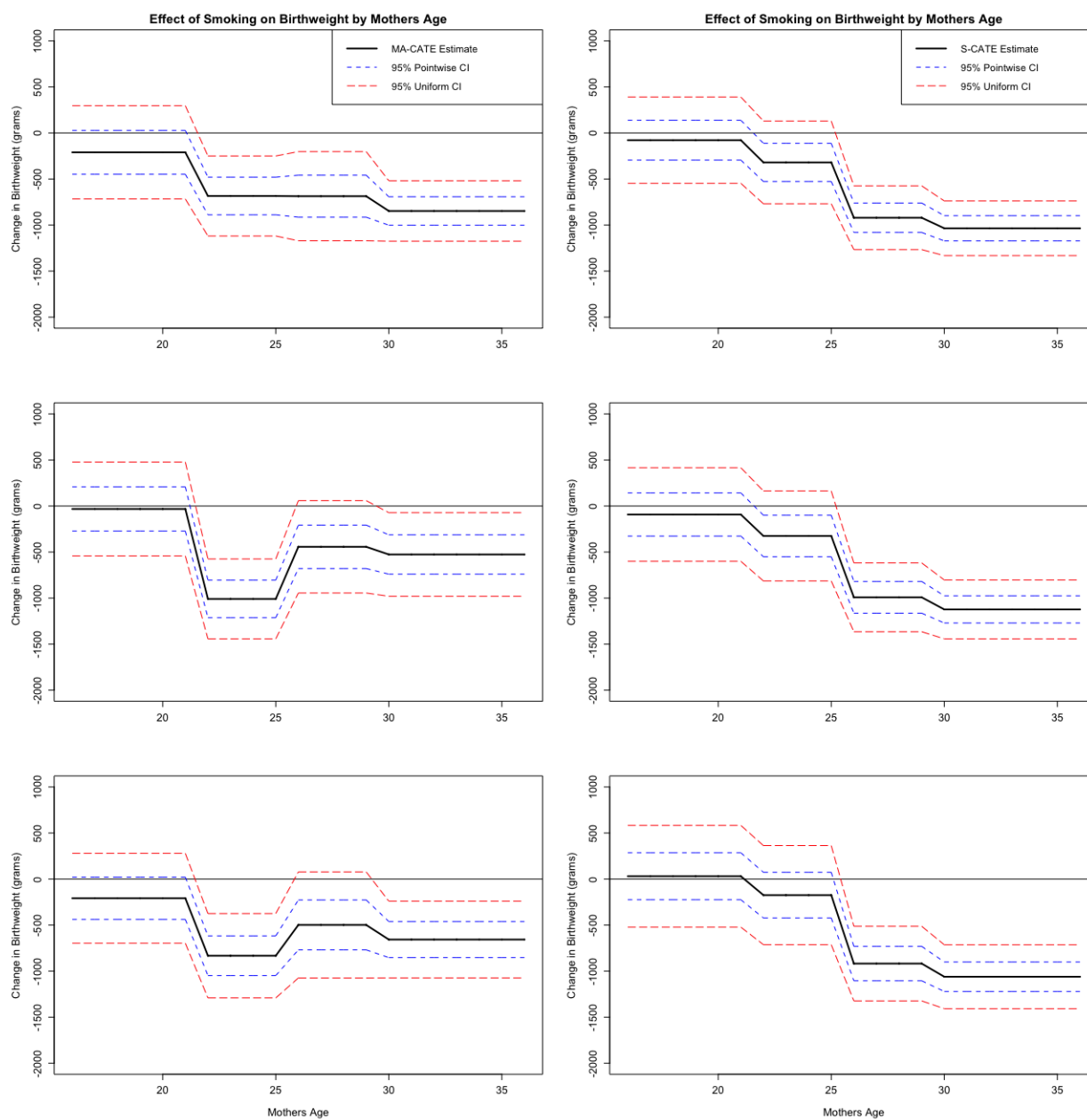


Figure F.1: CATE of maternal smoking estimated using model assisted estimating equations (left) and standard MLE/OLS estimating equations (right). Top row uses the 99<sup>th</sup> quantile of the bootstrap distribution to select the penalty parameters, second row uses 95<sup>th</sup> quantile, and final row uses the 90<sup>th</sup> quantile. Second stage is computed using a local constant with 3 knots. 95% pointwise confidence intervals are displayed in blue short dashes and 95% uniform confidence bands are displayed in long red dashes.

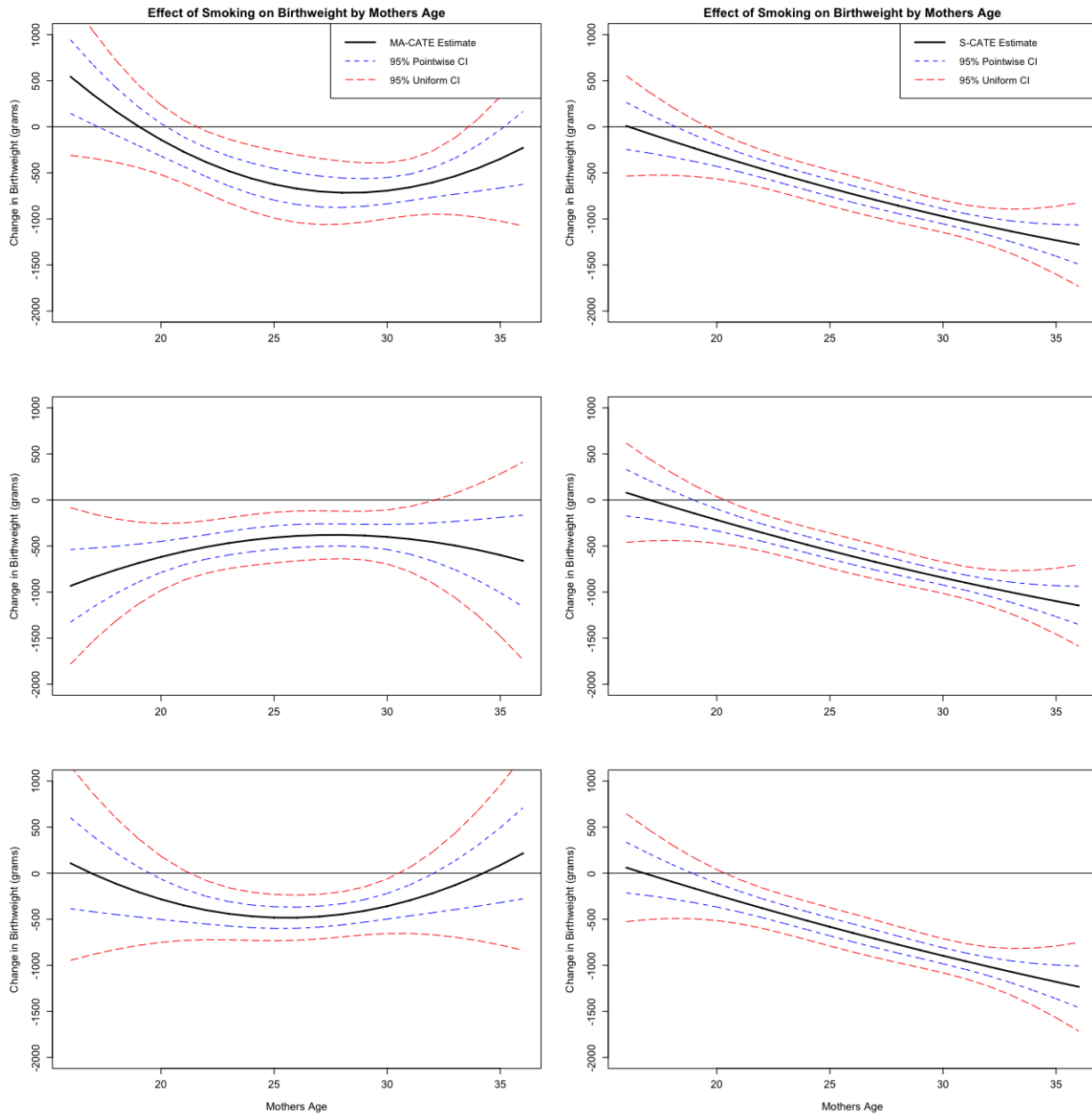


Figure F.2: CATE of maternal smoking estimated using model assisted estimating equations (left) and standard MLE/OLS estimating equations (right). Top row uses the 99<sup>th</sup> quantile of the bootstrap distribution to select the penalty parameters, second row uses 95<sup>th</sup> quantile, and final row uses the 90<sup>th</sup> quantile. Second stage is computed using second degree b-splines with 2 knots. 95% pointwise confidence intervals are displayed in blue short dashes and 95% uniform confidence bands are displayed in long red dashes.