# Causal Inference for Nonlinear Outcome Models with Possibly Invalid Instrumental Variables

Sai Li[*] and Zijian Guo[†]

## Abstract

Instrumental variable methods are widely used for inferring the causal effect of an exposure on an outcome when the observed relationship is potentially affected by unmeasured confounders. Existing instrumental variable methods for nonlinear outcome models require stringent identifiability conditions. We develop a robust causal inference framework for nonlinear outcome models, which relaxes the conventional identifiability conditions. We adopt a flexible semi-parametric potential outcome model and propose new identifiability conditions for identifying the model parameters and causal effects. We devise a novel three-step inference procedure for the conditional average treatment effect and establish the asymptotic normality of the proposed point estimator. We construct confidence intervals for the causal effect by the bootstrap method. The proposed method is demonstrated in a large set of simulation studies and is applied to study the causal effects of lipid levels on whether the glucose level is normal or high over a mice dataset.

*Keywords:* unmeasured confounders; binary outcome; semi-parametric model; endogeneity; partial mean; Mendelian Randomization

[*]Department of Biostatistics, Epidemiology and Informatics, Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA 19104 (E-mail: *sai.li@pennmedicine.upenn.edu*).
[†]Department of Statistics, Rutgers University, Piscataway, NJ 08854 (E-mail: *zijguo@stat.rutgers.edu*).

# 1 Introduction

Inference for the causal effect is a fundamental task in many fields. For instance, in epidemiology and genetics, identifying causal risk factors for diseases and health-related conditions can deepen our understandings of etiology and biological processes. In many applications, the effect of an exposure on an outcome is possibly nonlinear. For example, binary outcome models are widely used for studying the health conditions and the occurrence of diseases (Davey Smith and Ebrahim 2003; Davey Smith and Hemani 2014). It is of importance to make accurate inference for causal effects in nonlinear outcome models.

The existence of unmeasured confounders is a major concern for inferring causal effects in observational studies. The instrumental variable (IV) approach is the state-of-the-art method for estimating the causal effects when the unmeasured confounders potentially affect the observed relationships (Wooldridge 2010). As illustrated in Figure 1, the success of IV-based methods requires the candidate IVs to satisfy three core conditions: conditioning on the observed covariates, (A1) the candidate IVs are associated with the exposure; (A2) the candidate IVs have no direct effects on the outcome; and (A3) the candidate IVs are independent with unmeasured confounders.
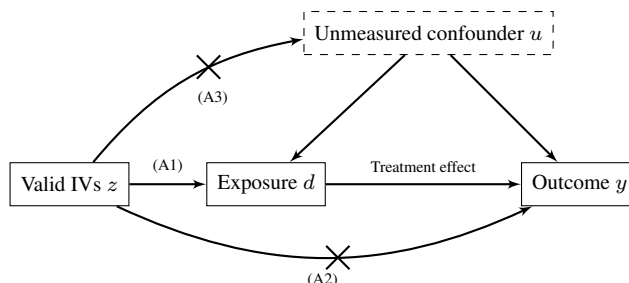


Figure 1: Illustration of IV assumptions (A1)-(A3).

The major challenge of applying IV-based methods is to identify IVs satisfying (A1)-(A3) simultaneously (Bowden et al. 2015; Kolesár et al. 2015; Kang et al. 2016). The assumptions (A2) and (A3) cannot even be tested in a data-dependent way. There is a pressing need to develop causal inference approaches when the candidate IVs are possibly invalid, say, violating assumptions (A2) or (A3) or both. There is a growing interest in using genetic variants as IVs, known as Mendelian Randomization (MR); see Voight et al. (2012) for an application example. Although genetic variants are subject to little environmental effects and are unlikely to have reverse causation

(Davey Smith and Ebrahim 2003; Lawlor et al. 2008), certain genetic variants are possibly invalid IVs due to the existence of pleiotropic effects (Davey Smith and Ebrahim 2003; Davey Smith and Hemani 2014), that is, one genetic variant can influence both the exposure and outcome simultaneously. In applications of MR, many outcome variables are dichotomous, e.g., the health conditions and disease status.

In the framework of linear outcome models, some recent progress has been made in inferring causal effects with possibly invalid IVs (Bowden et al. 2015; Kolesár et al. 2015; Bowden et al. 2016; Kang et al. 2016; Hartwig et al. 2017; Guo et al. 2018; Windmeijer et al. 2019). However, in consideration of binary and other nonlinear outcome models, existing methods (Blundell and Powell 2004; Rothe 2009) rely on the prior knowledge of a set of valid IVs. There is a lack of methods for inferring the causal effects in nonlinear outcome models with possibly invalid IVs.

## 1.1 Our results and contributions

The current paper focuses on inference for causal effects in nonlinear outcome models with unmeasured confounders. We propose a robust causal inference framework which covers a rich class of nonlinear outcome models and allows for possibly invalid IVs. Specifically, we propose a semi-parametric potential outcome model to capture the nonlinear effect, which includes logistic model, probit model, and multi-index models for continuous and binary outcome variables. The candidate IVs are allowed to be invalid and the invalid effects are modeled semi-parametrically, see equation (9). This generalizes the invalid IV framework for linear outcome models (Kang et al. 2016; Guo et al. 2018; Windmeijer et al. 2019), where the effect of invalid IVs is restricted to be additive and linear.

To identify the causal effect in semi-parametric outcome models, we introduce two identifiability conditions: dimension reduction (Condition 2.2) and majority rule (Condition 2.3). These identifiability conditions weaken the conventional conditions (summarized in Condition 2.1) for identifying the model parameters in semi-parametric outcome models (Blundell and Powell 2004; Rothe 2009). Specifically, the causal effect can be identified when a proportion of the candidate IVs are invalid and there is no knowledge on which candidate IVs are valid. We show that these two conditions are sufficient to identify the model parameters and conditional average treatment effect (CATE).

We propose a three-step inference procedure for CATE in $\underline{\text{S}}$emi-$\underline{\text{p}}$arametric $\underline{\text{out}}$come models with possibly invalid $\underline{\text{IV}}$s, termed as SpotIV. First, we estimate the reduced-form parameters based on semi-parametric dimension reduction methods. Second, we apply the median rule to estimate the model parameters by leveraging the fact that more than $50\%$ of candidate IVs are valid. Third, we develop a partial mean estimator to make inference for CATE. We establish the asymptotic normality of our proposed SpotIV estimator and construct confidence intervals for CATE by bootstrap. We demonstrate our proposed SpotIV method using a stock mice dataset and make inference for the casual effects of the lipid levels on whether the glucose level is normal or high.

We establish the asymptotic normality of our proposed SpotIV estimator of CATE, which can be viewed as a partial mean estimator. Our theoretical analysis generalizes the existing literature on partial means. The existing partial mean approaches (Newey 1994; Linton and Nielsen 1995) focus on the standard non-parametric regression settings with direct observations of the covariates. In contrast, the SpotIV estimator is a multi-index functional with indices estimated in a data-dependent way instead of directly observed. New techniques are proposed to handle the estimated indices and establish the asymptotic normality of the SpotIV estimator.

To sum up, the main contributions of this work are three-folded.

1. We introduce a robust causal inference framework for nonlinear outcome models allowing for possibly invalid IVs.

2. We propose new identification strategies of CATE in semi-parametric outcome models. To the authors' best knowledge, the SpotIV method is the first to make inference for causal effects in semi-parametric outcome models with possibly invalid IVs.

3. We develop new theoretical techniques to establish the asymptotic normality of the partial mean estimators with estimated indices.

## 1.2 Existing literature

Some recent progress has been made in inferring the causal effects with possibly invalid IVs under linear outcome models. With continuous outcome and exposure models, Bowden et al. (2015) and Kolesár et al. (2015) propose methods for causal effect estimation, which allow all candidate IVs to be invalid but assume orthogonality between the IV strengths and their invalid effects on the

outcome. Bowden et al. (2016), Kang et al. (2016), and Windmeijer et al. (2019) propose consistent estimators of causal effects assuming at least 50% of the IVs are valid. Hartwig et al. (2017) and Guo et al. (2018) consider linear outcome models under the assumption that the most common causal effect estimate is a consistent estimate of the true causal effect. Under this assumption, Guo et al. (2018) constructs confidence interval for the treatment effect and Windmeijer et al. (2019) further develops the inference procedure by refining the threshold levels of Guo et al. (2018). Verbanck et al. (2018) applies outlier detection methods to test horizontal pleiotropy. Spiller et al. (2019) proposes MRGxE, which assumes that the interaction effects of possibly invalid IVs and an environmental factor satisfy the valid IV assumptions (A1)-(A3). Tchetgen et al. (2019) introduces MR GENIUS which leverages a heteroscedastic covariance restriction. Bayesian approaches are also proposed to model invalid effects, to name a few, Thompson et al. (2017); Li (2017); Berzuini et al. (2020); Shapland et al. (2020). These methods are mainly developed for linear outcome models and cannot be extended to handle the inference problems in nonlinear outcome models.

There are two main streams of research on causal inference for nonlinear outcome models with unmeasured confounders. The first stream is based on parametric models, where the probit and logistic models are popular choices for modeling binary outcomes (Rivers and Vuong 1988; Vansteelandt et al. 2011). However, both models assume specific distributions of the unmeasured confounders, which limits their practical applications. The mixed-logistic model (Clarke and Windmeijer 2012), given in (32) of the current paper, is commonly used in observational studies. However, the IV-based two-stage method is biased for the mixed-logistic model (Cai et al. 2011). The main cause is that the odds ratio of the mixed-logistic model suffers from non-collapsibility. That is, the odds ratio depends on the distribution of unmeasured confounders and cannot be identified without distributional assumptions on the unmeasured confounders.

The second stream is based on semi-parametric models. Blundell and Powell (2004) and Rothe (2009) study causal inference for binary outcomes with double-index models assuming a known set of valid IVs and a valid control function. As mentioned, these assumptions can be impractical for applications such as MR. Moreover, the focus of Blundell and Powell (2004) and Rothe (2009) is on inference for model parameters, instead of causal estimands (e.g., CATE). In semi-parametric models, the model parameters are only identifiable up to certain linear transformations. The current

4

paper targets at inference for CATE, which can be uniquely identified, based on further innovations in methods and theory.

### 1.3 Organization of the rest of the paper

The rest of this paper is organized as follows. In Section 2, we introduce the model set-up and the identifiability conditions. In Section 3, we propose the strategies for identifying CATE. In Section 4, the SpotIV estimator is proposed to make inference for CATE. In Section 5, we provide theoretical guarantees for the proposed method. In Section 6, we investigate the empirical performance of the SpotIV estimator and compare it with the existing methods. In Section 7, our proposed method is applied to a dataset concerning the causal effects of high-density lipoproteins (HDL), low-density lipoproteins (LDL), and Triglycerides on the fasting glucose levels in a stock mice population. Section 8 concludes the paper.

## 2 Nonlinear Outcome Models with Possibly Invalid IVs

### 2.1 Models and causal estimands

For the $i$-th subject, $y_i \in \mathbb{R}$ denotes the observed outcome, $d_i \in \mathbb{R}$ denotes the exposure, $z_i \in \mathbb{R}^{p_z}$ denotes candidate IVs, and $x_i \in \mathbb{R}^{p_x}$ denotes baseline covariates. Define $p = p_z + p_x$ and we use $w_i = (z_i^\mathsf{T}, x_i^\mathsf{T})^\mathsf{T} \in \mathbb{R}^p$ to denote all measured covariates, including candidate IVs and baseline covariates. We assume that the data $\{y_i, d_i, w_i\}_{1 \le i \le n}$ are generated in *i.i.d.* fashions. Let $u_i$ denote the unmeasured confounder which can be associated with both exposure and outcome variables.

We define causal effects using the potential outcome framework (Neyman 1923; Rubin 1974). Let $y_i^{(d)} \in \mathbb{R}$ be the potential outcome if the $i$-th individual were to have exposure $d$. We consider the following nonlinear potential outcome model

$$\mathbb{E}[y_i^{(d)}|w_i = w, u_i = u] = q\left(d\beta + w^\mathsf{T}\kappa, u\right), \tag{1}$$

where $q : \mathbb{R}^2 \to \mathbb{R}$ is a (possibly unknown) link function, $\beta \in \mathbb{R}$ is the coefficient of the exposure, and $\kappa = (\kappa_z^\mathsf{T}, \kappa_x^\mathsf{T})^\mathsf{T} \in \mathbb{R}^p$ is the coefficient vector of the measured covariates. Model (1) includes a broad class of nonlinear potential outcome models, which can be used for both continuous and binary outcomes. The function $q$ can be either known or unknown. For binary outcomes,

if $q(a, b) = 1/(1 + \exp(-a - b))$, then the model (1) is logistic; if $q(a, b) = \mathbb{1}(a + b > 0)$ and $u_i$ is normal with mean zero, then the model (1) is the probit model.

We assume that $y_i^{(d)} \perp\!\!\!\perp d_i \mid (w_i^\mathsf{T}, u_i)$. This condition is mild as we can hypothetically identify the unmeasured variable $u_i$ such that $y_i^{(d)}$ and $d_i$ are conditionally independent. This is much weaker than the (strong) ignorability condition $y_i^{(d)} \perp\!\!\!\perp d_i \mid w_i$ (Rosenbaum and Rubin 1983). Under the condition $y_i^{(d)} \perp\!\!\!\perp d_i \mid (w_i^\mathsf{T}, u_i)$ and the consistency assumption (Imbens and Rubin 2015, e.g.), we can connect the conditional mean for the observed outcome $y_i$ and the potential outcome $y_i^{(d)}$ as

$$\mathbb{E}[y_i | d_i = d, w_i = w, u_i = u] = \mathbb{E}[y_i^{(d)} | d_i = d, w_i = w, u_i = u] = \mathbb{E}[y_i^{(d)} | w_i = w, u_i = u]. \quad (2)$$

As a result, the potential outcome model (1) leads to the following model for observed outcome $y_i$

$$\mathbb{E}[y_i | d_i = d, w_i = w, u_i = u] = q\left(d\beta + w^\mathsf{T}\kappa, u\right). \quad (3)$$

We focus on the continous exposure $d_i$ with linear conditional mean function

$$d_i = w_i^\mathsf{T}\gamma + v_i, \quad \mathbb{E}[v_i | w_i] = 0, \quad (4)$$

where $\gamma = (\gamma_z^\mathsf{T}, \gamma_x^\mathsf{T})^\mathsf{T}$ denotes the association between $w_i$ and $d_i$ and $v_i$ is the residual term. In observational studies, since the unmeasured confounder $u_i$ can be dependent with $v_i$, the exposure $d_i$ is associated with $u_i$ even after conditioning on the measured covariates $w_i$; see Figure 1.

The current paper studies the semi-parametric potential outcome model (1) and the exposure association model (4). The target causal estimand is CATE

$$\mathrm{CATE}(d, d' | w) := \mathbb{E}\left[y_i^{(d)} - y_i^{(d')} | w_i = w\right], \quad (5)$$

where $d \in \mathbb{R}$ and $d' \in \mathbb{R}$ denote two different exposure levels and $w \in \mathbb{R}^p$ denotes the specific value of measured covariates. The CATE can characterize the heterogeneity across subpopulations with different levels of measured covariates.

## 2.2 Review of the control function approach with valid IVs

While two-stage least squares based on valid IVs are popularly used for linear outcome models, the control function approach with valid IVs is widely adopted for causal inference when dealing with nonlinear outcome models (Blundell and Powell 2004; Rothe 2009; Petrin and Train 2010; Cai et al. 2011; Wooldridge 2015; Guo and Small 2016). The key idea of control functions is to treat the residual $v_i$ of the exposure model (4) as a proxy for the unmeasured confounder $u_i$ and to incorporate $v_i$ into the outcome model as an adjustment for the unmeasured confounder. The success of existing control function approaches relies on the following identifiability condition (Blundell and Powell 2004; Rothe 2009).

**Condition 2.1** (Valid IV and control function)**.** *The models for the candidate IVs $z_i \in \mathbb{R}^{p_z}$ satisfy $\|\gamma_z\|_2 \geq \tau_0 > 0$ in (4) and $\kappa_z = 0$ in (3), where $\tau_0$ is a positive constant. The conditional density $f_u(u_i|w_i, v_i)$ satisfies*

$$f_u(u_i|w_i, v_i) = f_u(u_i|v_i). \tag{6}$$

The condition $\|\gamma_z\|_2 \geq \tau_0 > 0$ assumes strong associations between the IVs and the exposure variable, which corresponds to the classical IV assumption (A1). The condition $\kappa_z = 0$ assumes that the IVs do not have direct effects on the outcome, which corresponds to (A2). Equation (6) assumes that conditioning on the control variable $v_i$, the unmeasured confounder $u_i$ is independent of the measured covariates $w_i$. This assumption can be viewed as a version of (A3) for nonlinear outcome models. In the special case of no baseline covariates $x_i$, condition (6) is equivalent to (A3) given that $v_i$ is independent of $z_i$. However, such a connection is not obvious in general. Condition 2.1 can be illustrated in Figure 1 by replacing (A3) with its nonlinear version (6).

Under Condition 2.1, the outcome model (3) can be written as

$$\mathbb{E}[y_i|d_i, w_i, v_i] = \int q(d_i\beta + w_i^\mathsf{T}\kappa, u_i) f_u(u_i|v_i) du_i = g_0\left(d_i\beta + x_i^\mathsf{T}\kappa_x, v_i\right), \tag{7}$$

where $g_0 : \mathbb{R}^2 \to \mathbb{R}$ is an unknown function. Inference for parameters $\beta$ and $\kappa_x$ in (7) has been studied in Blundell and Powell (2004) and Rothe (2009) under Condition 2.1.

Although Condition 2.1 is commonly adopted for the control function approach, it can be challenging to identify IVs satisfying Condition 2.1 in applications. As explained, the valid IV

assumptions (A2) and (A3) are likely to be violated when using genetic variants as IVs in the MR applications. Moreover, (6) is unlikely to hold when $u_i$ involves omitted variables, which may be associated with measured covariates $w_i$. As pointed out in Blundell and Powell (2004), a valid control function largely relies on including all the suspicious confounders into the model, which may be a strong assumption for practical applications. To make things worse, these identifiability assumptions, including both $\kappa_z = 0$ and (6), are untestable in a data-dependent way.

## 2.3   Identifiability conditions with possibly invalid IVs

To better accommodate for practical applications, we introduce new identifiability conditions, which weaken Condition 2.1.

**Condition 2.2** (Dimension reduction). *The conditional density $f_u(u_i|w_i, v_i)$ satisfies*

$$f_u(u_i|w_i, v_i) = f_u(u_i|w_i^\mathsf{T}\eta, v_i) \quad \text{for some } \eta \in \mathbb{R}^{p \times q}. \tag{8}$$

In contrast to (6), expression (8) allows the unmeasured confounder $u_i$ to depend on the measured covariates $w_i$ after conditioning on the control variable $v_i$. Condition 2.2 essentially requires a dimension reduction property of the conditional density $f_u(u_i|w_i, v_i)$. In particular, the dependence on $w_i$ is captured by the linear combinations $w_i^\mathsf{T}\eta \in \mathbb{R}^q$ conditioning on $v_i$. To better illustrate the main idea, we focus on the case of $q = 1$ and $\eta \in \mathbb{R}^p$ being a vector throughout the rest of the paper. Our framework and methods can be directly extended to the settings with some finite integer $1 \le q < p$. In view of (8), the conditional mean of the outcome can be written as

$$\mathbb{E}[y_i|d_i, w_i, v_i] = \int q(d_i\beta + w_i^\mathsf{T}\kappa, u_i) f_u(u_i|w_i^\mathsf{T}\eta, v_i) du_i = g^*\left(d_i\beta + w_i^\mathsf{T}\kappa, w_i^\mathsf{T}\eta, v_i\right). \tag{9}$$

In comparison to (7), the above model allows $\kappa_z \neq 0$ and has an additional additional index $w_i^\mathsf{T}\eta$, which is induced by the dependence of $u_i$ and $w_i^\mathsf{T}\eta$ as in (8).

Now we introduce another identifiability condition which states that a majority of the candidate IVs are valid. Let $\mathcal{S}$ be the set of relevant IVs, i.e., $\mathcal{S} = \{1 \le j \le p_z : \gamma_j \neq 0\}$ and $\mathcal{V}$ be the set of valid IVs, i.e.,

$$\mathcal{V} = \{j \in \mathcal{S} : (\kappa_z)_j = (\eta_z)_j = 0\}.$$

The set $\mathcal{S}$ contains all candidate IVs that are strongly associated with the exposure. The set $\mathcal{V}$ is a subset of $\mathcal{S}$ which contains all candidate IVs satisfying the classical IV assumptions $(\kappa_z)_j = 0$ and $(\eta_z)_j = 0$. For $j \in \mathcal{S} \cap \mathcal{V}^c$, the corresponding IV can have $(\kappa_z)_j \neq 0$ or $(\eta_z)_j \neq 0$ or both of them, i.e., these IVs violate the classical identifiability condition (Condition 2.1).

When the candidate IVs are possibly invalid, the main challenge of causal inference is that the set $\mathcal{V}$ is unknown a priori in data analysis. The following identifiability condition is needed for identifying the causal effect without any prior knowledge on the set of valid IVs $\mathcal{V}$.

**Condition 2.3** (Majority rule). *More than half of the relevant IVs are valid:* $|\mathcal{V}| > |\mathcal{S} \cap \mathcal{V}^c|$.

The majority rule assumes that more than half of the relevant IVs are valid but does not require prior knowledge of the set $\mathcal{V}$. The majority rule has been proposed in linear outcome models with invalid IVs (Bowden et al. 2016; Kang et al. 2016; Guo et al. 2018; Windmeijer et al. 2019).

To summarize, Conditions 2.2 and 2.3 are the new identifiability conditions to identify causal effects in the semi-parametric outcome model (1) with possibly invalid IVs. These two conditions (Figure 2) weaken Condition 2.1 and better accommodate for practical applications.
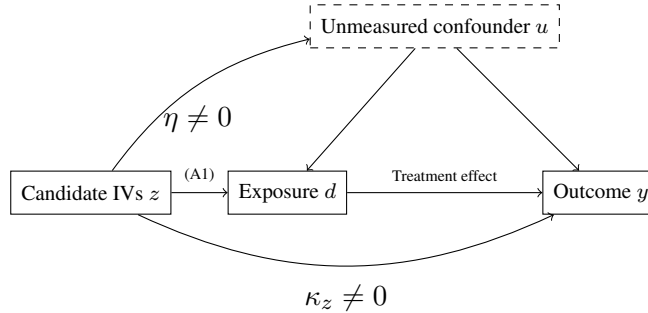


Figure 2: Illustration of the new identifiability conditions (Conditions 2.2 and 2.3) in the presence of unmeasured confounders.

## 3 Causal Effects Identification

In this section, we describe how to identify the CATE$(d, d'|w)$ defined in (5) for nonlinear outcome models under Conditions 2.2 and 2.3. We introduce another causal estimand, the average structural function (ASF),

$$\text{ASF}(d, w) = \int \mathbb{E}[y_i^{(d)}|w_i = w, v_i = v] f_v(v) dv, \tag{10}$$

where $f_v$ is the density of the residue $v_i$ defined in (4). For binary outcomes, the $\text{ASF}(d, w)$ represents the response probability for a given pair of $(d, w)$ (Newey 1994; Blundell and Powell 2004) and it is a policy relevant quantity in econometrics. The ASF is closely related to CATE in the sense that if $w_i$ and $v_i$ are independent, then

$$\text{CATE}(d, d'|w) = \text{ASF}(d, w) - \text{ASF}(d', w). \tag{11}$$

In the following, we present a three-step strategy for identifying ASF and CATE. The data-dependent algorithm is presented in Section 4.

### 3.1 Identification of the reduced-form parameters

The conditional mean function (9) can be re-written as

$$\mathbb{E}[y_i|d_i, w_i, v_i] = g^*((d_i, w_i^\mathsf{T})B^*, v_i) \quad \text{with} \quad B^* = \begin{pmatrix} \beta & 0 \\ \kappa & \eta \end{pmatrix} \in \mathbb{R}^{(p+1)\times 2}, \tag{12}$$

where $g^* : \mathbb{R}^3 \to \mathbb{R}$ is defined in (9). Due to the collinearity among $d_i, w_i$, and $v_i$, we cannot directly identify $B^*$ in the conditional mean model (12). We will deduce a reduced-form representation of (12) by combining it with (4). As $\mathbb{E}[y_i|w_i, v_i] = \mathbb{E}[y_i|d_i, w_i, v_i]$, we derive the reduced-form model

$$\mathbb{E}[y_i|w_i, v_i] = \mathbb{E}[y_i|w_i^\mathsf{T}\Theta^*, v_i] \quad \text{with} \quad \Theta^* = (\gamma, \mathrm{I}_p)B^* \in \mathbb{R}^{p\times 2}, \tag{13}$$

where $\mathrm{I}_p$ is the $p \times p$ identity matrix. Although $\Theta^*$ cannot not be uniquely identified in the above model, we can identify $\Theta^*$ up to a linear transformation; that is, we can identify some parameter $\Theta \in \mathbb{R}^{p\times M}$ such that

$$\mathbb{E}[y_i|w_i, v_i] = \mathbb{E}[y_i|w_i^\mathsf{T}\Theta, v_i] \quad \text{and} \quad \Theta = \Theta^*T \tag{14}$$

where $T \in \mathbb{R}^{2\times M}$ is a linear transformation matrix for a positive integer $M$. While $\Theta$ can have $M$ columns for any integer $M \geq 1$, it is implied by (13) that $M$ is at most two. In words, $w_i^\mathsf{T}\Theta$ is a sufficient summary of the mean dependence of $y_i$ on $w_i$ given $v_i$. In the semi-parametric literature, identifying some $\Theta$ satisfying (14) is closely related to the estimation of the central subspace or central mean space (Cook and Li 2002; Cook 2009). Our detailed implementation is described in

10

Section 4.1. In the rest of this section, we assume that there exists some reduced-form matrix $\Theta$ such that (14) holds and discuss how to identify the model parameters and the causal effects.

## 3.2 Identification of model parameters

The model parameter of interest is $B \in \mathbb{R}^{p \times M}$ such that

$$\Theta = (\gamma, \mathrm{I}_p)B, \tag{15}$$

where $B = B^*T$ with the same transformation $T$ in (14). The parameter $B$ is a linear transformation of original parameter $B^*$. Since $\Theta$ and $\gamma$ can be directly identified from the data, we can apply the majority rule (Condition 2.3) to identify the matrix $B$ based on (15). Specifically, for $1 \leq m \leq M$, define $b_m = \mathrm{Median}(\{\Theta_{j,m}/\gamma_j\}_{j \in \mathcal{S}})$, where $\mathcal{S}$ denotes the set of relevant IV. We identify $B$ as

$$B = \begin{pmatrix} b_1 & \ldots & b_M \\ \Theta_{.,1} - b_1\gamma & \ldots & \Theta_{.,M} - b_M\gamma \end{pmatrix} \tag{16}$$

for some $\Theta$ satisfying (14). Here $\Theta_{.,j}$ denotes the $j$-th column of $\Theta$. The rationale for $B$ in (16) is the same as the application of majority rule in linear outcomes models: each candidate IV can produce an estimate of the causal effect $\beta$ based on the ratio of the reduced-form parameter and the IV strength $\gamma$; the median of these ratios will be $\beta$ if more than half of the relevant IVs are assumed to be valid. The definition of $B$ in (16) generalizes this idea to semi-parametric outcome models.

The following proposition shows that $(d_i, w_i^\intercal)B$ and $v_i$ are a sufficient summary of the conditional mean of $y_i$ given $d_i, w_i$ and $v_i$.

**Proposition 3.1.** *Under Conditions 2.2 and 2.3, the parameter $B$ defined in (16) satisfies (15) and* $\mathbb{E}[y_i|d_i, w_i, v_i] = \mathbb{E}[y_i|(d_i, w_i^\intercal)B, v_i]$.

With $B$ in (16), we define the conditional mean function $g : \mathbb{R}^{M+1} \to \mathbb{R}$ as

$$g((d_i, w_i^\intercal)B, v_i) = \mathbb{E}[y_i|(d_i, w_i^\intercal)B, v_i]. \tag{17}$$

11

As a remark, the conditional mean function $g$ implicitly depends on $B$ but $g((d_i, w_i^\mathsf{T})B, v_i) = \mathbb{E}[y_i|d_i, w_i, v_i]$ is invariant to $B$.

**Remark 3.1.** Some other conditions for identifying $B$ can be used to replace the majority rule in Proposition 3.1. First, a version of the orthogonal condition considered in Bowden et al. (2015) and Kolesár et al. (2015) is sufficient for identifying $B$ in the current framework. Specifically, if both $\kappa$ and $\eta$ are orthogonal to $\gamma$, then the correlation between $\Theta_{.,m}$ and $\gamma$ is $b_m$ for $m = 1, \ldots, M$. Second, the plurality rule considered in Guo et al. (2018) can be used to identify the parameter $B$. Although the plurality rule is a relaxation of the majority rule, the implementation of the plurality rule depends on the limiting distribution of the estimated parameters, which is computationally expensive in the semi-parametric scenario.

### 3.3 Identification of causal estimands

In the following proposition, we demonstrate how to identify ASF and CATE based on the parameter $B$ defined in (16) and the function $g$ defined in (17).

**Proposition 3.2.** *Under Conditions 2.2 and 2.3, it holds that*

$$\mathbb{E}\left[y_i^{(d)}|w_i = w, v_i = v\right] = g\left((d, w^\mathsf{T})B, v\right) \tag{18}$$

*where $B$ is defined in (16) and $g$ is defined in (17).*

Proposition 3.2 implies that the conditional mean of the potential outcome can be identified via the identification of the model parameter $B$ and the nonparametric function $g$. As $B$ can be identified as in (16), $g(\cdot)$ can be identified using the conditional mean of the observed outcome. Hence, the ASF$(d, w)$ defined in (10) can be identified by taking an integration of $g((d, w^\mathsf{T})B, v_i)$ with respect to the density of $v_i$. The CATE can be identified via its relationship with ASF function as in (11).

## 4 Methodology: SpotIV

In this section we formally introduce the SpotIV method, which implements the three-step identification strategies derived in Section 3 in a data-dependent way. We illustrate the procedure for

binary outcome models in Sections 4.1 to 4.3 and discuss its generalization to continuous nonlinear outcome models in Section 4.4.

## 4.1 Step 1: estimation of the reduced-form parameters

We estimate the reduced-form parameter $\Theta$ satisfying (14) based on the semi-parametric dimension reduction methods. Various approaches have been proposed for semi-parametric dimension reduction; see, for example, Li (1991); Xia et al. (2002); Ma and Zhu (2012). Notice that the linear space spanned by $\Theta$ defined in (14) is different from the broadly studied mean dimension-reduction space or central subspace (Cook 2009) as the index $v_i$ is given. Our specific procedure is derived from the sliced-inverse regression approach (SIR) (Li 1991).

Let $\Sigma = \text{Cov}((w_i^\intercal, v_i)^\intercal) \in \mathbb{R}^{(p+1) \times (p+1)}$ denote the covariance matrix of $(w_i^\intercal, v_i)^\intercal$ and $\alpha(y_i) = \mathbb{E}[\Sigma^{-1/2}(w_i^\intercal, v_i)^\intercal | y_i] \in \mathbb{R}^{p+1}$ denote the inverse regression function. For the covariance matrix $\Omega = \text{Cov}(\alpha(y_i)) \in \mathbb{R}^{(p+1) \times (p+1)}$, we use $M_\Omega = \text{rank}(\Omega)$ to denote its rank and $\Phi \in \mathbb{R}^{(p+1) \times M_\Omega}$ to denote the matrix of eigenvectors corresponding to non-zero eigenvalues. We first introduce an estimation procedure of $\Phi$ by assuming a known rank $M_\Omega$. A consistent estimate of $M_\Omega$ will be provided in (22). We fit the first-stage model (4) based on least squares,

$$\hat\gamma = (W^\intercal W)^{-1} W^\intercal d \quad \text{and} \quad \hat v = d - W\hat\gamma. \tag{19}$$

Define $\widehat\Sigma = \frac{1}{n} \sum_{i=1}^{n} (w_i^\intercal, \hat v_i)^\intercal (w_i^\intercal, \hat v_i)$. For $k = 0, 1$, we estimate $\alpha(k)$ by

$$\hat\alpha(k) = \frac{1}{\sum_{i=1}^{n} \mathbb{1}(y_i = k)} \sum_{i=1}^{n} \mathbb{1}(y_i = k) \widehat\Sigma^{-1/2}(w_i^\intercal, \hat v_i)^\intercal$$

and estimate $\Omega$ by $\widehat\Omega = \widehat{\mathbb{P}}(y_i = 1)\widehat{\mathbb{P}}(y_i = 0)\{\hat\alpha(1) - \hat\alpha(0)\}\{\hat\alpha(1) - \hat\alpha(0)\}^\intercal$, where $\widehat{\mathbb{P}}(y_i = 1) = \sum_{i=1}^{n} \mathbb{1}(y_i = 1)/n$ and $\widehat{\mathbb{P}}(y_i = 0) = 1 - \widehat{\mathbb{P}}(y_i = 1)$. Let $\hat\lambda_1 \geq \cdots \geq \hat\lambda_{p+1}$ denote the eigenvalues of $\widehat\Omega$ and $\widehat\Phi \in \mathbb{R}^{(p+1) \times M_\Omega}$ denotes the matrix of the eigenvectors of $\widehat\Omega$ corresponding to $\hat\lambda_1, \ldots, \hat\lambda_{M_\Omega}$.

Now we introduce an estimate of $\Theta$ using the matrix $\widehat\Phi$. Define

$$(i^*, j^*) = \underset{1 \leq i, j \leq M_\Omega}{\arg\min} \, i + j \tag{20}$$

$$\text{subject to} \ |\text{cor}(\widehat\Phi_{1:p,i}, \widehat\Phi_{1:p,j})| \leq 1 - \sqrt{\frac{\log n}{n}},$$

where $\widehat{\Phi}_{1:p,j}$ denotes the first $p$ elements of $\widehat{\Phi}_{\cdot,j} \in \mathbb{R}^{p+1}$ and $\mathrm{cor}(a,b) = \langle a,b \rangle / (\|a\|_2 \|b\|_2)$ if $a \neq 0$ and $b \neq 0$ and $\mathrm{cor}(a,b) = 0$ otherwise. If all vectors $\{\widehat{\Phi}_{1:p,i}\}_{1 \leq i \leq M_\Omega}$ are collinear, there is no solution to (20) with a high probability. Taking this into consideration, we construct the estimator of $\Theta$ as

$$\widehat{\Theta} = \begin{cases} (\widehat{\Phi}_{1:p,i^*}, \widehat{\Phi}_{1:p,j^*}) & \text{if (20) has a solution,} \\ \widehat{\Phi}_{1:p,1} & \text{otherwise.} \end{cases} \tag{21}$$

We now provide explanations for (20) and (21). Let $\Phi_{1:p,\cdot} \in \mathbb{R}^{p \times M_\Omega}$ denote the sub-matrix containing the first $p$ rows of $\Phi$. We can show that a valid $\Theta$ satisfying (14) is a basis of the column space $\Phi_{1:p,\cdot}$. The columns of $\widehat{\Theta}$ in (21) estimate a minimum set of basis for the column space of $\Phi_{1:p,\cdot}$. Since (13) implies $M = \mathrm{rank}(\Theta) \leq 2$, the column rank of $\Phi_{1:p,\cdot}$ is at most two. If (20) has a solution, then the column space of $\Phi_{1:p,\cdot}$ is two-dimensional with high probability and hence $\widehat{\Theta}$ in (21) takes two linearly independent columns of $\widehat{\Phi}_{1:p,\cdot}$; if (20) does not have a solution, then the column space of $\Phi_{1:p,\cdot}$ is one-dimensional with high probability and $\widehat{\Theta}$ takes the first column of $\widehat{\Phi}_{1:p,\cdot}$. Indicated by the definition (21), $\widehat{M} = \mathrm{rank}(\widehat{\Theta})$ is either one or two.

To determine $M_\Omega$, a BIC-type procedure in Zhu et al. (2006) can be applied. Specifically, the dimension $M_\Omega$ can be estimated as,

$$\widehat{M}_\Omega = \underset{1 \leq m \leq 3}{\arg\max}\, C(m) \text{ with } C(m) = \frac{n}{2} \sum_{i=m+1}^{p} \{\log(\hat{\lambda}_i + 1) - \hat{\lambda}_i\} \mathbb{1}(\hat{\lambda}_i > 0) - \frac{C_n \cdot m(2p - m + 1)}{2}, \tag{22}$$

where $C_n = n^{c_0}$ (with $0 < c_0 < 1$) is a penalty constant and $m(2p - m + 1)/2$ is the degree of freedom. The true dimension $M_\Omega$ is at most three because the dimension of $\Theta$ in (14) is at most two. The consistency of $\widehat{M}_\Omega$ follows from Theorem 2 in Zhu et al. (2006) under mild conditions. For a better illustration of this approach, we assume $M_\Omega$ is known in the following.

**Remark 4.1.** Other dimension reduction methods can be used to estimate $\Theta$. We adopt the SIR approach mainly for its computational efficiency. The computational cost of the SIR estimate $\widehat{\Phi}$ is relatively low in comparison to the semi-parametric ordinary least square estimator (Ichimura 1993) and semi-parametric maximum likelihood estimator for binary outcomes (Klein and Spady 1993). The aforementioned two methods are based on kernel approximations of $g(\cdot)$ and the optimization is not convex in general, which requires much more computational power than SIR.

## 4.2 Step 2: estimation of the model parameter $B$

We proceed to estimate the model parameter $B$ defined in (16). To apply the majority rule, we first select the set of relevant IVs by

$$\widehat{\mathcal{S}} = \left\{ 1 \leq j \leq p_z : |\widehat{\gamma}_j| \geq \widehat{\sigma}_v \sqrt{2\{\widehat{\Sigma}^{-1}\}_{j,j} \log n / n} \right\}, \tag{23}$$

where $\widehat{\sigma}_v^2 = \sum_{i=1}^n (d_i - w_i \widehat{\gamma})^2 / n$ and $\widehat{\Sigma}$ is defined after (19). The term $\log n$ is the adjustment for the multiplicity of the selection procedure. Under mild conditions, $\widehat{\mathcal{S}}$ is shown to be a consistent estimate of $\mathcal{S}$. As a remark, such a thresholding has been proposed in Guo et al. (2018) and a possibly finer threshold can be found in Windmeijer et al. (2019). With $\widehat{\gamma}$ and $\widehat{\Theta}$ defined in (19) and (21), respectively, we provide an estimator of $B$ by leveraging the majority rule detailed in (16). Specifically, for $m = 1, \ldots, \widehat{M}$, we define $\widehat{b}_m = \text{Median}\left(\left\{\widehat{\Theta}_{j,m}/\widehat{\gamma}_j\right\}_{j \in \widehat{\mathcal{S}}}\right)$ and

$$\widehat{B} = \begin{pmatrix} \widehat{b}_1 & \cdots & \widehat{b}_{\widehat{M}} \\ \widehat{\Theta}_{\cdot,1} - \widehat{b}_1 \widehat{\gamma} & \cdots & \widehat{\Theta}_{\cdot,\widehat{M}} - \widehat{b}_{\widehat{M}} \widehat{\gamma} \end{pmatrix}. \tag{24}$$

where $\widehat{\Theta}_{\cdot,m}$ denotes the $m$-th column of $\widehat{\Theta}$.

## 4.3 Step 3: inference for causal effects

We propose inference procedures for $\text{ASF}(d, w)$ defined in (10) and for $\text{CATE}(d, d'|w)$ defined in (5). In view of Proposition 3.2, after identifying the parameter matrix $B$, we further estimate function $g(\cdot)$ defined in (17). With $\widehat{B}$ defined in (24), we estimate $g$ by a kernel estimator $\widehat{g}$. Let $s_i = ((d, w^\intercal)B, v_i)^\intercal$ denote the true index at the given level $(d, w^\intercal)$. Denote the estimated indices as $\widehat{s}_i = ((d, w^\intercal)\widehat{B}, \widehat{v}_i)^\intercal$ and $\widehat{t}_i = ((d_i, w_i^\intercal)\widehat{B}, \widehat{v}_i)^\intercal$, for $1 \leq i \leq n$. Define the kernel $K_H(a, b)$ for $a, b \in \mathbb{R}^{\widehat{M}+1}$ as $K_H(a, b) = \prod_{l=1}^{\widehat{M}+1} \frac{1}{h_l} k\left(\frac{a_l - b_l}{h_l}\right)$ where $h_l$ is the bandwidth for the $l$-th argument and $k(x) = \mathbf{1}(|x| \leq 1/2)$. To focus on the main result, we take $K_H$ in the form of product kernel and $k(x)$ as the box kernel and set $h_l = h$ for $1 \leq l \leq \widehat{M} + 1$. We estimate $\{g(s_i)\}_{1 \leq i \leq n}$ by the kernel estimator

$$\widehat{g}(\widehat{s}_i) = \frac{\frac{1}{n} \sum_{j=1}^n y_j K_H(\widehat{s}_i, \widehat{t}_j)}{\frac{1}{n} \sum_{j=1}^n K_H(\widehat{s}_i, \widehat{t}_j)} \quad \text{for} \quad 1 \leq i \leq n$$

and estimate $\mathrm{ASF}(d, w) = \int g(s_i) f_v(v_i) dv_i$ by a sample average with respect to $\widehat{v}_i$ (or equivalently $\widehat{s}_i$),

$$\widehat{\mathrm{ASF}}(d, w) = \frac{1}{n} \sum_{i=1}^{n} \widehat{g}(\widehat{s}_i). \tag{25}$$

Estimating $\mathrm{ASF}(d', w)$ analogously, we estimate $\mathrm{CATE}(d, d'|w)$ as

$$\widehat{\mathrm{CATE}}(d, d'|w) = \widehat{\mathrm{ASF}}(d, w) - \widehat{\mathrm{ASF}}(d', w). \tag{26}$$

In Section 5.2, we establish the asymptotic normality of $\widehat{\mathrm{CATE}}(d, d'|w)$ under regularity conditions. By approximating its variance by bootstrap, we construct the confidence interval for $\mathrm{CATE}(d, d'|w)$ as

$$\left( \widehat{\mathrm{CATE}}(d, d'|w) - z_{1-\alpha/2}\hat{\sigma}^*, \quad \widehat{\mathrm{CATE}}(d, d'|w) + z_{1-\alpha/2}\hat{\sigma}^* \right), \tag{27}$$

where $z_{1-\alpha/2}$ is the $1 - \alpha/2$ quantile of standard normal and $\hat{\sigma}^*$ is the standard deviation estimated by $N$ bootstrap samples.

### 4.4 Extensions to continuous nonlinear outcome models

The SpotIV procedure for binary outcomes detailed in Section 4.1 to 4.3 can be extended to deal with continuous nonlinear outcome models. The main change is to use a different estimator of $\Omega = \mathrm{Cov}(\alpha(y_i)) \in \mathbb{R}^{(p+1) \times (p+1)}$. Specifically, $\Omega$ can be estimated based on SIR (Li 1991) or kernel-based method (Zhu and Fang 1996). With such an estimate of $\Omega$, we can apply the same procedure in Sections 4.1 to 4.3 and make inference for CATE for continuous outcome models. We examine the numerical performance of our proposal for continuous nonlinear outcome models in Section 6.

## 5 Theoretical Justifications

In this section we provide theoretical justifications of our proposed method for binary outcome models. In Section 5.1, we present the estimation accuracy of the model parameter matrix $\widehat{B}$. In Section 5.2, we establish the asymptotic normality of the proposed SpotIV estimator under proper conditions.

## 5.1 Estimation accuracy of model parameter matrix

We introduce the required regularity conditions in the following and start with the moment conditions on the observed data.

**Condition 5.1.** (Moment conditions) *The observed data* $(y_i, d_i, w_i^\intercal)^\intercal$, $i = 1, \ldots, n$, *are i.i.d. generated with* $\mathbb{E}[v_i|w_i] = 0$ *and* $\mathbb{E}[(w_i^\intercal, v_i)^\intercal(w_i^\intercal, v_i)]$ *being positive definite. Moreover,* $\{w_{i,j}\}_{1 \leq j \leq p}$ *and* $v_i$ *are sub-Gaussian random variables.*

Next, we introduce the regularity conditions for the SIR method. Let $\mathbb{P}_{\mathbb{S}}(w_i^\intercal, v_i)^\intercal$ denote the projection of $(w_i^\intercal, v_i)^\intercal \in \mathbb{R}^{p+1}$ onto a linear subspace $\mathbb{S}$ of $\mathbb{R}^{p+1}$. Let $\mathbb{C}$ denote the intersection of all the subspaces $\mathbb{S}$ such that $\mathbb{P}(y_i = 1|w_i, v_i) = \mathbb{P}(y_i = 1|\mathbb{P}_{\mathbb{S}}(w_i^\intercal, v_i)^\intercal)$. The linear subspace $\mathbb{C}$ is indeed the central subspace for the distribution of $y_i$ conditioning on $w_i$ and $v_i$ (Cook 2009).

**Condition 5.2** (Regularity conditions for SIR). *The linear subspace* $\mathbb{C}$ *uniquely exists. It holds that* $\mathbb{E}[w_i|\mathbb{P}_{\mathbb{C}}(w_i^\intercal, v_i)^\intercal]$ *is linear in* $\mathbb{P}_{\mathbb{C}}(w_i^\intercal, v_i)^\intercal$. *The nonzero eigenvalues of* $\Omega = Cov(\alpha(y_i))$ *are simple, where* $\alpha(y_i) = \mathbb{E}[\Sigma^{-1/2}(w_i^\intercal, v_i)^\intercal|y_i] \in \mathbb{R}^{p+1}$ *denotes the inverse regression function.*

Existence and uniqueness of $\mathbb{C}$ can be guaranteed under mild conditions (Cook 2009). The condition that $\mathbb{E}[w_i|\mathbb{P}_{\mathbb{C}}(w_i^\intercal, v_i)^\intercal]$ is linear in $\mathbb{P}_{\mathbb{C}}(w_i^\intercal, v_i)^\intercal$ is known as the linearity assumption and is standard for SIR methods (Li 1991; Cook and Lee 1999; Chiaromonte et al. 2002). A sufficient condition for the linearity assumption is that, $w_i$ is normal and is independent of $v_i$. The simple nonzero eigenvalues of $\Omega$ guarantee the uniqueness of the matrix $\Phi$ as the true parameters. Similar assumptions have been imposed in Zhu and Fang (1996). The next lemma establishes the convergence rate of $\widehat{B} - B$.

**Lemma 5.1.** *Assume Conditions 2.2, 2.3, 5.1, and 5.2 hold. Then*

$$\mathbb{P}\left(\|\widehat{B} - B\|_2 \geq c_1\sqrt{t/n}\right) \leq \exp(-c_2 t) + \mathbb{P}(E_1^c), \tag{28}$$

*where* $\mathbb{P}(E_1^c) \to 0$ *as* $n \to \infty$ *and* $c_1, c_2 > 0$ *are positive constants.*

As shown in Lemma 5.1, the proposed $\widehat{B}$ converges at rate $n^{-1/2}$. The true parameter $B$ and the event $E_1$ are given in the proof of Lemma 5.1 in the supplementary materials. Intuitively speaking, the high probability event $E_1$ is the intersection of the events $\{\widehat{M} = M\}$, $\{\widehat{S} = S\}$, and

{The median $\hat{b}_m$ are evaluated at valid IVs}. As a remark, the result in Lemma 5.1 still holds if the estimator $\widehat{\Theta}$ is replaced with other $\sqrt{n}$-consistent estimators of $\Theta$.

## 5.2 Asymptotic normality

In the following, we establish the asymptotic normality of the proposed SpotIV estimator and shall focus on the case $M = 2$. We introduce the following assumptions on the density function $f_t$ of $t_i = ((d_i, w_i^\intercal)B, v_i)^\intercal \in \mathbb{R}^3$ and the unknown function $g$ defined in (17) at $s_i = ((d, w^\intercal)B, v_i)^\intercal \in \mathbb{R}^3$. We define

$$\mathcal{N}_h(s) = \left\{ t \in \mathbb{R}^3 : \|t - s\|_\infty \leq h \right\}, \tag{29}$$

where $\|\cdot\|_\infty$ denotes the vector maximum norm.

**Condition 5.3** (Smoothness conditions). *(a) The density function $f_t$ of $t_i = ((d_i, w_i^\intercal)B, v_i)^\intercal$ has a convex support $\mathcal{T} \subset \mathbb{R}^3$ and satisfies $c_0 \leq f_t(s_i) \leq C_0$ for all $1 \leq i \leq n$, $\int_{t \in \mathcal{T}^{\mathrm{int}}} f_t(t)dt = 1$ and $\max_{1 \leq i \leq n} \sup_{t \in \mathcal{N}_h(s_i) \cap \mathcal{T}} \|\nabla f_t(t)\|_\infty \leq C$, where $\mathcal{T}^{\mathrm{int}}$ is the interior of $\mathcal{T}$, $\mathcal{N}_h(s)$ is defined in (29), $\nabla f_t$ is the gradient of $f_t$ and $C_0 > c_0 > 0$ and $C > 0$ are positive constants. The density $f_v$ of $v_i$ is bounded and has a convex support $\mathcal{T}_v$.*

*(b) The function $g$ defined in (17) is twicely differentiable. For any $1 \leq i \leq n$, $g(s_i)$ is bounded away from zero and one. The function $g$ satisfies $\max_{1 \leq i \leq n} \sup_{t \in \mathcal{N}_h(s_i) \cap \mathcal{T}} \|\nabla g(t)\|_2 \leq C$ and $\max_{1 \leq i \leq n} \sup_{t \in \mathcal{N}_h(s_i) \cap \mathcal{T}} \lambda_{\max}(\triangle g(t)) \leq C$, where $\mathcal{N}_h(s)$ is defined in (29), $\|\nabla g(t)\|_2$ and $\lambda_{\max}(\triangle g(t))$ respectively denote the $\ell_2$ norm of the gradient vector and the largest eigenvalue of the hessian matrix of $g$ evaluated at $t$ and $C > 0$ is a positive constant.*

*(c) For any $v \in \mathcal{T}_v$, then the evaluation point $(d, w^\intercal)^\intercal$ satisfies $((d, w^\intercal)B + \Delta^\intercal, v)^\intercal \in \mathcal{T}$ for any $\Delta \in \mathbb{R}^2$ satisfying $\|\Delta\|_\infty \leq h$.*

Condition 5.3(a) and 5.3(b) are mainly imposed for the regularities of the density function $f_t$, $f_v$, and the conditional mean function $g$ at $s_i = ((d, w^\intercal)B, v_i)^\intercal$ or its neighborhood $\mathcal{N}_h(s_i)$. Here the randomness of $s_i$ only depends on $v_i$ for the pre-specified evaluation point $(d, w^\intercal)^\intercal$. Condition 5.3(c) essentially assumes that the evaluation point $(d, w^\intercal)$ is not at the tail of the joint distribution of $(d_i, w_i^\intercal)$. These conditions are mild and will be verified in the supplementary materials, see Propositions A.3, A.4, and A.5. Specifically, when $M = 2$, there is a one-to-one correspondence

18

between $t_i$ and $t_i^* = ((d_i, w_i^\mathsf{T})B^*, v_i)$, where $B^*$ denotes the parameter matrix defined in (12). We will verify Condition 5.3 (a) under the regularity conditions on the density function $f_{t^*}$ of $t_i^*$. 5.3 (b) will be implied by the regularity conditions on the potential outcome model $q(\cdot)$ defined in (1). If $q(\cdot)$ is continuous, it suffices to require that $q(\cdot)$ has bounded second derivatives and the conditional density $f_u(u_i|w_i^\mathsf{T}\eta, v_i)$ belongs to a location-scale family with smooth mean and variance functions. If $q(\cdot)$ is an indicator function, then $g$ becomes the conditional density of $u_i$ given $w_i^\mathsf{T}\eta$ and $v_i$ and it suffices to require this conditional density function to satisfy Condition 5.3 (b). Examples of $q$ functions satisfying Condition 5.3 (b) include the logistic or probit models with uniformly bounded $v_i$.

The following theorem establishes the asymptotic normality of the proposed ASF estimator.

**Theorem 5.1.** *Suppose that, $M = 2$, Condition 5.3 holds, and the bandwidth satisfies $h = n^{-\mu}$ for $0 < \mu < 1/4$. For any estimator $\widehat{B}$ satisfying (51), with probability larger than $1 - n^{-c} - \mathbb{P}(E_1^c)$,*

$$\left| \widehat{\mathrm{ASF}}(d, w) - \mathrm{ASF}(d, w) \right| \leq C \left( \frac{1}{\sqrt{nh^2}} + h^2 \right). \tag{30}$$

*where $\mathbb{P}(E_1^c) \to 0$ as $n \to \infty$ and $c > 0$ and $C > 0$ are some positive constant. Taking $h = n^{-\mu}$ for $0 < \mu < 1/6$, we have*

$$\frac{n}{\sqrt{\mathrm{V}}} \left( \widehat{\mathrm{ASF}}(d, w) - \mathrm{ASF}(d, w) \right) \xrightarrow{d} N(0, 1) \quad with \quad \mathrm{V} = \sqrt{\sum_{j=1}^{n} a_j^2 g(t_j)(1 - g(t_j))}$$

*where $a_j = \frac{1}{n} \sum_{i=1}^{n} \frac{K_H(s_i, t_j)}{\frac{1}{n}\sum_{j=1}^{n} K_H(s_i, t_j)}$ for $1 \leq j \leq n$ and $\xrightarrow{d}$ denotes the convergence in distribution. The asymptotic standard error satisfies*

$$\mathbb{P}\left( c_0/\sqrt{nh^2} \leq \sqrt{\mathrm{V}}/n \leq C_0/\sqrt{nh^2} \right) \geq 1 - n^{-c}$$

*for some positive constants $C_0 \geq c_0 > 0$ and $c > 0$.*

A few remarks are in order for this main theorem. Firstly, the rate in (30) is the same as the optimal rate of estimating a twicely differentiable function in two dimensions (Tsybakov 2008). Though the unknown target function $\mathrm{ASF}(d, w)$ can be viewed as a two-dimension function on linear combinations of $d$ and $w$, it cannot be directly estimated using the classical nonparametric

methods. In contrast, we have to first estimate the unknown function $g$ in three dimensions and then further estimate the target $\mathrm{ASF}(d, w)$. After a careful analysis, we establish that, even though $\mathrm{ASF}(d, w)$ involves estimating the three-dimension function $g$, the final convergence rate can be reduced to the same rate as estimating two-dimensional twicely differentiable smooth functions.

Secondly, beyond Condition 5.3, the above theorem requires a suitable bandwidth condition $h = n^{-\mu}$ with $0 < \mu < 1/6$ for establishing the asymptotic normality, which is standard in nonparametric regression in two dimensions (Wasserman 2006). This bandwidth condition essentially requires the variance component to dominate its bias, that is, $(nh^2)^{-1/2} \gg h^2$. Thirdly, we can establish asymptotic normality for a large class of initial estimators $\widehat{B}$ as long as they satisfy (51). By Lemma 5.1, our proposed estimator $\widehat{B}$ belongs to this class of initial estimators with a high probability.

Lastly, we shall emphasize the technical novelties of establishing Theorem 5.1. The proposed estimator of $\mathrm{ASF}(d, w)$ can be viewed as integrating the three-dimension function $g$. The main step in the proof is to show that the error or asymptotic variance of estimating $\mathrm{ASF}(d, w)$ is the same as estimating two-dimension twicely differentiable functions. This type of results has been established in Newey (1994) and Linton and Nielsen (1995) under the name "partial mean". However, our proof is distinguished from the standard partial mean problem in the sense that we do not have access to direct observations of $s_i$ and $t_i$ but only have their estimators $\widehat{s}_i$ and $\widehat{t}_i$ for $1 \leq i \leq n$. Due to the dependence between the estimators $\{\widehat{s}_i, \widehat{t}_i\}_{1 \leq i \leq n}$ and the errors $y_i - g((d_i, w_i)B, v_i)$, it is challenging to adopt the standard partial mean techniques and establish asymptotic normality. We have developed new techniques to decouple the dependence between $\{\widehat{s}_i, \widehat{t}_i\}_{1 \leq i \leq n}$ and the errors. The techniques depend on introducing "enlarged support kernels" to control the errors between $K_H(\widehat{s}_i, \widehat{t}_i)$ and $K_H(s_i, t_i)$. These techniques are of independent interest for other related problems in handling partial means with estimated indexes.

We now provide theoretical guarantees for $\widehat{\mathrm{CATE}}(d, d'|w)$ defined in (26). Similar to the definition of $s_i$, we define $r_i = ((d', w^\intercal)B, v_i)$ as the corresponding multiple indices by fixing $(d_i, w_i^\intercal)$ at the given level $(d', w^\intercal)$. The following corollary establishes the asymptotic normality of the proposed estimator $\widehat{\mathrm{CATE}}(d, d'|w)$.

**Corollary 5.1.** *Suppose that Condition 5.3 holds for both $\{s_i\}_{1 \leq i \leq n}$ and replacing $\{s_i\}_{1 \leq i \leq n}$ and $w$ by $\{r_i\}_{1 \leq i \leq n}$ and $w'$, respectively. Suppose that, $M = 2$, $v_i$ is independent of $w_i$, the bandwidth*

*satisfies $h = n^{-\mu}$ for $0 < \mu < 1/6$, and $|d - d'| \cdot \max\{|B_{11}|, |B_{21}|\} \geq h$. For any estimator $\widehat{B}$ satisfying* (51)*, then*

$$\frac{n}{\sqrt{V_{\text{CATE}}}} \left( \widehat{\text{CATE}}(d, d'|w) - \text{CATE}(d, d'|w) \right) \xrightarrow{d} N(0, 1) \text{ with } V_{\text{CATE}} = \sqrt{\sum_{j=1}^{n} c_j^2 g(t_j)(1 - g(t_j))}$$

*where $c_j = \frac{1}{n} \sum_{i=1}^{n} \left( \frac{K_H(s_i, t_j)}{\frac{1}{n} \sum_{j=1}^{n} K_H(r_i, t_j)} - \frac{K_H(r_i, t_j)}{\frac{1}{n} \sum_{j=1}^{n} K_H(r_i, t_j)} \right)$, for $1 \leq j \leq n$. The asymptotic standard error satisfies*

$$\mathbf{P} \left( c_0/\sqrt{nh^2} \leq \sqrt{V_{\text{CATE}}}/n \leq C_0/\sqrt{nh^2} \right) \geq 1 - n^{-c} \tag{31}$$

*for some positive constants $C_0 \geq c_0 > 0$ and $c > 0$.*

Corollary 5.1 is closely related to Theorem 5.1. The asymptotic normality of $\text{ASF}(d', w)$ can be established with a similar argument to Theorem 5.1 with replacing $s_i$ by $r_i$. When $v_i$ is independent of the measured covariates $w_i$, we apply (11) to compute CATE by taking the difference of $\widehat{\text{ASF}}(d, w)$ and $\widehat{\text{ASF}}(d', w)$. An extra step is to show that the asymptotic normal component of $\widehat{\text{ASF}}(d, w) - \widehat{\text{ASF}}(d', w)$ dominates its bias component. To ensure this, an extra assumption on the difference between $d$ and $d'$, $|d - d'| \cdot \max\{|B_{11}|, |B_{21}|\} \geq h$, is needed to guarantee the lower bound for $\sqrt{V_{\text{CATE}}}/n$ in (31).

## 6 Numerical Studies

In this section, we assess the empirical performance of the proposed method for both binary and continuous outcome models. We detail our optimization method as follows. Following Zhu et al. (2006), we select $\widehat{M}_{\Omega}$ according to (22) with $C_n = c^{-1} \log n$ where $c$ is the number of observations in each slice. We estimate $\widehat{\Phi}$ using the SIR method in the R package `np` (Hayfield and Racine 2008) and then obtain $\widehat{\Theta} \in \mathbb{R}^{p \times \widehat{M}}$ via (21). Next, we estimate $\mathcal{S}$ by $\widehat{\mathcal{S}}$ defined in (23) and estimate $B$ by $\widehat{B}$ defined in (24). Finally, we estimate CATE as in (26) with the bandwidth selected by 5-fold cross validation. To construct confidence intervals for CATE, we use the standard deviation of $N = 50$ bootstrap realizations of CATEs to estimate its standard error. The R code for our proposal is available at `https://github.com/saili0103/SpotIV`.

We consider four simulation scenarios in the following and plot their corresponding $\text{ASF}(d, w)$ (as a function of $d$) in Figure 3 with $p = 7$ and $w = (0, \ldots, 0, 0.1)^{\mathsf{T}} \in \mathbb{R}^7$. The first two sce-

narios correspond to binary outcome models and the last two scenarios correspond to continuous nonlinear outcome models. The ASF and hence the CATE functions are all nonlinear across these scenarios.
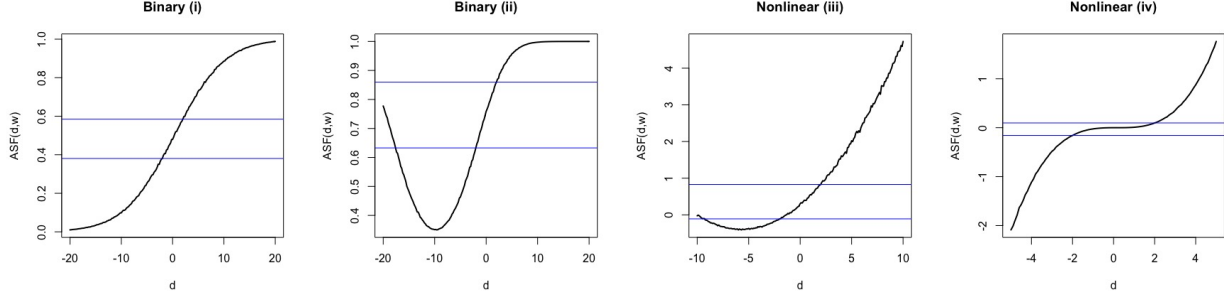


Figure 3: The curves correspond to the functions $\mathrm{ASF}(d, w)$ in the four scenarios considered in this section. The blue lines give the true values for $d = -2$ and $d = 2$ in each scenario.

## 6.1 Binary outcome models

The exposure $d_i$ is generated as $d_i = z_i^\intercal \gamma + v_i$, where $\gamma = c_\gamma \cdot (1, 1, 1, -1, -1, -1, -1)^\intercal$ and $v_i$ are *i.i.d.* normal with mean zero and variance $\sigma_v^2 = 1$. We vary the strength of the IV, $c_\gamma \in \{0.4, 0.6, 0.8\}$, and consider the setting with no measured covariates $x_i$, i.e., $w_i = z_i$. We generate two distributions of the $z_i$: (1) $\{z_i\}_{1 \leq i \leq n}$ are *i.i.d.* $N(0, \mathrm{I}_p)$; (2) $\{z_i\}_{1 \leq i \leq n}$ are *i.i.d.* uniformly distributed in $[-1.73, 1.73]$. We generate the outcome models as follows.

(i) We generate $y_i$, $1 \leq i \leq n$, via the logistic model

$$\mathbb{P}(y_i = 1 \mid d_i, w_i, u_i) = \mathrm{logit}\,(d_i\beta + w_i^\intercal \kappa + u_i). \tag{32}$$

with $\beta = 0.25$, $\kappa = \eta = (0, 0, 0, 0, 0, 0.4, -0.4)^\intercal$ and $\mathrm{logit}(x) = 1/(1 + \exp(-x))$. We generate the unmeasured confounder $u_i$ as

$$u_i = 0.25v_i + w_i^\intercal \eta + \xi_i, \ \xi_i \sim N(0, (w_i^\intercal \eta)^2). \tag{33}$$

The model (32) is known as the mixed-logistic model. After integrating out $u_i$ conditioning on $v_i, w_i$, the conditional distribution $y_i$ given $d_i, w_i$ is in general not logistic.

22

(ii) We generate $y_i$, $1 \le i \le n$, via

$$\mathbb{P}(y_i = 1 | d_i, w_i, u_i) = \text{logit}\left(d_i\beta + w_i^\intercal\kappa + u_i + (d_i\beta + w_i^\intercal\kappa + u_i)^2/3\right),$$

with $\beta = 0.25$ and $\kappa = \eta = (0, 0, 0, 0, 0, 0.4, -0.4)^\intercal$. We generate the unmeasured confounder $u_i$ as

$$u_i = \exp(0.25v_i + w_i^\intercal\eta) + \xi_i, \ \xi_i \sim U[-1, 1]. \tag{34}$$

In both configurations, conditioning on $w_i$, the unmeasured confounder $u_i$ is correlated with $v_i$ and $d_i$ and the majority rule is satisfied: the first five IVs are valid and the last two are invalid. We construct 95% confidence intervals for CATE$(d, d'|w)$. We compare the proposed SpotIV estimator with two state-of-the-art methods. The first one is the semi-parametric MLE with valid control function and valid IVs (Rothe 2009), shorthanded as Valid-CF. While the Valid-CF is not derived for the invalid setting, the main purpose of this comparison is to understand how invalid IVs affect the accuracy of the causal inference approaches by assuming valid IVs. We also compare SpotIV with a method called Logit-Median, which is detailed in Section C in the supplementary material. This method models the conditional outcome model as a logistic function, which can be a mis-specified model after integrating out the unmeasured confounder $u_i$. The same majority rule as the proposed SpotIV method is implemented to estimate the model parameters. The purpose of making this comparison is to understand the effect of the mis-specified outcome model. Detailed implementation of Valid-CF and Logit-Median are described in Section C in the supplement.

All simulation results are calculated over 500 replications. In Table 1, we report the inference results for CATE$(-2, 2|w)$ in binary outcome model (i). The proposed SpotIV method has the empirical coverage close to the nominal level for both Gaussian and uniform $w_i$. The estimation errors get smaller when the IVs become stronger or when the sample size becomes larger. In contrast, the Valid-CF method, assuming all IVs to be valid, has larger estimation errors, mainly due to the bias of using invalid IVs. The empirical coverage of the Valid-CF is lower than 95% in most settings.

In Table 2, we report the inference results CATE$(-2, 2|w)$ in binary outcome model (ii). The pattern is similar to that in Table 1 for binary outcome model (i). The Valid-CF approach has a larger bias and lower coverage when IVs become stronger. This is because when IVs are stronger,

| | | $N(0, \mathrm{I}_p)$ | | | | | | $U[-1.73, 1.73]$ | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | SpotIV | | | Valid-CF | | | SpotIV | | | Valid-CF | | |
| $n$ | $c_\gamma$ | MAE | COV | SE | MAE | COV | SE | MAE | COV | SE | MAE | COV | SE |
| 500 | 0.4 | 0.094 | 0.962 | 0.14 | 0.121 | 0.877 | 0.13 | 0.098 | 0.968 | 0.14 | 0.100 | 0.906 | 0.13 |
| 500 | 0.6 | 0.064 | 0.942 | 0.10 | 0.081 | 0.883 | 0.11 | 0.064 | 0.962 | 0.10 | 0.090 | 0.920 | 0.11 |
| 500 | 0.8 | 0.055 | 0.950 | 0.09 | 0.075 | 0.917 | 0.10 | 0.050 | 0.960 | 0.09 | 0.084 | 0.920 | 0.10 |
| 1000 | 0.4 | 0.067 | 0.960 | 0.10 | 0.088 | 0.892 | 0.11 | 0.065 | 0.956 | 0.10 | 0.089 | 0.906 | 0.11 |
| 1000 | 0.6 | 0.048 | 0.980 | 0.07 | 0.064 | 0.922 | 0.08 | 0.041 | 0.960 | 0.07 | 0.062 | 0.893 | 0.08 |
| 1000 | 0.8 | 0.038 | 0.946 | 0.06 | 0.060 | 0.920 | 0.08 | 0.040 | 0.956 | 0.06 | 0.059 | 0.903 | 0.08 |
| 2000 | 0.4 | 0.051 | 0.960 | 0.07 | 0.072 | 0.874 | 0.09 | 0.050 | 0.946 | 0.08 | 0.075 | 0.870 | 0.09 |
| 2000 | 0.6 | 0.032 | 0.932 | 0.05 | 0.043 | 0.916 | 0.06 | 0.033 | 0.954 | 0.05 | 0.049 | 0.912 | 0.06 |
| 2000 | 0.8 | 0.028 | 0.970 | 0.05 | 0.046 | 0.870 | 0.06 | 0.034 | 0.954 | 0.05 | 0.047 | 0.903 | 0.06 |

Table 1: Inference for CATE$(-2, 2|w)$ in the binary outcome model (i). The columns indexed with "MAE", "COV" and "SE" report the median absolute errors of $\widehat{\mathrm{CATE}}(-2, 2|w)$, the empirical coverages of the confidence intervals and the average of estimated standard errors of the point estimators, respectively. The columns indexed with "SpotIV" and "Valid-CF" correspond to the proposed method and the method assuming valid IVs, respectively.

the variance of the estimator is smaller and the bias is relatively more significant. The empirical coverage of Logit-Median (Table 5 in the supplement) also gets lower with a larger sample size and a stronger IV. This demonstrates the bias caused by the model mis-specification.

## 6.2 General nonlinear outcome models

We consider two nonlinear continuous outcome models.

(iii) We generate $y_i$, $i = 1, \ldots, n$ via $y_i = d_i\beta + z_i^{\mathsf{T}}\kappa + u_i + (d_i\beta + z_i^{\mathsf{T}}\kappa + u_i)^2/3$, where $u_i$ is generated via (33).

(iv) We generate $y_i$, $i = 1, \ldots, n$ via $y_i = u_i(d_i\beta + z_i^{\mathsf{T}}\kappa)^3$, where $u_i$ is generated via (34). This is an example of double-index format of (1).

The true parameters in (iii) and (iv) are set to be the same as in Section 6.1.

We compare the SpotIV estimator with the two-stage hard-thresholding (TSHT) method (Guo et al. 2018), which is proposed to deal with possibly invalid IVs in linear outcome models. The purpose of this comparison is to understand the effect of mis-specifying a nonlinear model as linear. The proposed SpotIV method has coverage probabilities close to 95% in model (iii) and model (iv) (Table 3 and Table 4). In comparison, the TSHT does not guarantee the 95% coverage

| | | $N(0, \mathrm{I}_p)$ | | | | | | $U[-1.73, 1.73]$ | | | | | |
| | | SpotIV | | | Valid-CF | | | SpotIV | | | Valid-CF | | |
| $n$ | $c_\gamma$ | MAE | COV | SE | MAE | COV | SE | MAE | COV | SE | MAE | COV | SE |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 500 | 0.4 | 0.085 | 0.940 | 0.13 | 0.105 | 0.867 | 0.11 | 0.091 | 0.940 | 0.13 | 0.089 | 0.880 | 0.10 |
| 500 | 0.6 | 0.061 | 0.930 | 0.09 | 0.077 | 0.873 | 0.09 | 0.063 | 0.940 | 0.09 | 0.081 | 0.882 | 0.09 |
| 500 | 0.8 | 0.050 | 0.960 | 0.08 | 0.073 | 0.863 | 0.08 | 0.052 | 0.920 | 0.06 | 0.068 | 0.884 | 0.08 |
| 1000 | 0.4 | 0.060 | 0.962 | 0.09 | 0.074 | 0.893 | 0.09 | 0.064 | 0.949 | 0.09 | 0.071 | 0.854 | 0.08 |
| 1000 | 0.6 | 0.046 | 0.946 | 0.07 | 0.069 | 0.843 | 0.07 | 0.052 | 0.929 | 0.07 | 0.071 | 0.854 | 0.07 |
| 1000 | 0.8 | 0.039 | 0.944 | 0.06 | 0.062 | 0.763 | 0.06 | 0.043 | 0.940 | 0.06 | 0.065 | 0.800 | 0.06 |
| 2000 | 0.4 | 0.049 | 0.952 | 0.07 | 0.066 | 0.843 | 0.07 | 0.047 | 0.954 | 0.07 | 0.062 | 0.833 | 0.07 |
| 2000 | 0.6 | 0.034 | 0.946 | 0.05 | 0.061 | 0.800 | 0.06 | 0.035 | 0.931 | 0.05 | 0.061 | 0.786 | 0.05 |
| 2000 | 0.8 | 0.027 | 0.938 | 0.04 | 0.057 | 0.720 | 0.04 | 0.032 | 0.934 | 0.04 | 0.065 | 0.674 | 0.05 |

Table 2: Inference for CATE$(-2, 2|w)$ in the binary outcome model (ii). The columns indexed with "MAE", "COV" and "SE" report the median absolute errors of $\widehat{\mathrm{CATE}}(-2, 2|w)$, the empirical coverages of the confidence intervals and the average of estimated standard errors of the point estimators, respectively. The columns indexed with "SpotIV" and "Valid-CF" correspond to the proposed method and the method assuming valid IVs, respectively.

and has larger estimation errors, mainly due to the fact that the TSHT method is developed for linear outcome models.

# 7 Applications to Mendelian Randomization

We apply the proposed SpotIV method to make inference for the effects of the lipid levels on the glucose level in a stock mice population. The dataset is available at https://wp.cs.ucl.ac.uk/outbredmice/heterogeneous-stock-mice/. It consists of 1,814 subjects, where for each subject, 10,346 polymorphic genetic markers, certain phenotypes, and baseline covariates are available. After removing observations with missing values, the remaining sample size is 1,269.

Fasting glucose level is an important indicator of type-2 diabetes and rodent models have been broadly used to study the risk factors of diabetes for adults (Islam and du Loots 2009; King 2012). According to Fajardo et al. (2014), we dichotomize the fasting glucose level at 11.1 (unit: mmol/L) and consider $\leq 11.1$ as normal and $> 11.1$ as high (pre-diabetic and diabetic). The proportion of high fasting glucose level is approximately 25.1%. We study the causal effects of three lipid levels (HDL, LDL, and Triglycerides) on whether the fasting glucose level is normal or high for this stock mice population. We include "gender" and "age" as baseline covariates. The polymorphic markers and covariates are standardized before analysis.

|  |  | $N(0, \mathrm{I}_p)$ | | | | | | $U[-1.73, 1.73]$ | | | | | |
|  |  | SpotIV | | | TSHT | | | SpotIV | | | TSHT | | |
| $n$ | $c_\gamma$ | MAE | COV | SE | MAE | COV | SE | MAE | COV | SE | MAE | COV | SE |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 500 | 0.4 | 0.249 | 0.970 | 0.46 | 1.146 | 0.310 | 0.35 | 0.261 | 0.962 | 0.44 | 1.13 | 0.276 | 0.34 |
| 500 | 0.6 | 0.201 | 0.960 | 0.35 | 0.315 | 0.634 | 0.24 | 0.202 | 0.954 | 0.35 | 0.258 | 0.666 | 0.24 |
| 500 | 0.8 | 0.161 | 0.948 | 0.31 | 0.290 | 0.594 | 0.18 | 0.176 | 0.962 | 0.31 | 0.269 | 0.610 | 0.18 |
| 1000 | 0.4 | 0.174 | 0.960 | 0.30 | 0.200 | 0.916 | 0.26 | 0.183 | 0.970 | 0.29 | 0.190 | 0.916 | 0.25 |
| 1000 | 0.6 | 0.125 | 0.974 | 0.23 | 0.127 | 0.902 | 0.18 | 0.135 | 0.938 | 0.22 | 0.136 | 0.896 | 0.17 |
| 1000 | 0.8 | 0.128 | 0.958 | 0.20 | 0.128 | 0.842 | 0.14 | 0.125 | 0.943 | 0.20 | 0.108 | 0.856 | 0.13 |
| 2000 | 0.4 | 0.124 | 0.942 | 0.21 | 0.146 | 0.886 | 0.18 | 0.126 | 0.931 | 0.20 | 0.129 | 0.894 | 0.18 |
| 2000 | 0.6 | 0.090 | 0.969 | 0.16 | 0.120 | 0.840 | 0.12 | 0.113 | 0.914 | 0.16 | 0.114 | 0.826 | 0.12 |
| 2000 | 0.8 | 0.078 | 0.946 | 0.13 | 0.111 | 0.756 | 0.10 | 0.100 | 0.920 | 0.14 | 0.114 | 0.770 | 0.09 |

Table 3: Inference for CATE$(-2, 2|w)$ in continuous outcome model (iii). The columns indexed with "MAE", "COV" and "SE" report the median absolute errors of $\widehat{\mathrm{CATE}}(-2, 2|w)$, the empirical coverages of the confidence intervals and the average of estimated standard errors of the point estimators, respectively. The columns indexed with "SpotIV" and "TSHT" correspond to the proposed method and the method proposed in (Guo et al. 2018), respectively.

|  |  | $N(0, \mathrm{I}_p)$ | | | | | | $U[-1.73, 1.73]$ | | | | | |
|  |  | SpotIV | | | TSHT | | | SpotIV | | | TSHT | | |
| $n$ | $c_\gamma$ | MAE | COV | SE | MAE | COV | SE | MAE | COV | SE | MAE | COV | SE |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 500 | 0.4 | 0.075 | 0.988 | 0.14 | 1.887 | 0.370 | 0.63 | 0.067 | 0.974 | 0.13 | 2.463 | 0.164 | 0.39 |
| 500 | 0.6 | 0.059 | 0.982 | 0.12 | 2.263 | 0.157 | 0.51 | 0.061 | 0.956 | 0.11 | 2.552 | 0.030 | 0.33 |
| 500 | 0.8 | 0.059 | 0.980 | 0.12 | 2.562 | 0.044 | 0.44 | 0.060 | 0.952 | 0.11 | 2.918 | 0 | 0.30 |
| 1000 | 0.4 | 0.063 | 0.954 | 0.11 | 1.749 | 0.156 | 0.48 | 0.046 | 0.974 | 0.10 | 1.758 | 0.006 | 0.30 |
| 1000 | 0.6 | 0.048 | 0.970 | 0.09 | 2.053 | 0.106 | 0.39 | 0.045 | 0.974 | 0.09 | 2.131 | 0 | 0.24 |
| 1000 | 0.8 | 0.045 | 0.966 | 0.09 | 2.531 | 0.010 | 0.37 | 0.052 | 0.976 | 0.09 | 2.675 | 0 | 0.22 |
| 2000 | 0.4 | 0.042 | 0.974 | 0.08 | 1.804 | 0.020 | 0.37 | 0.044 | 0.974 | 0.08 | 1.743 | 0 | 0.22 |
| 2000 | 0.6 | 0.036 | 0.980 | 0.07 | 2.122 | 0.014 | 0.30 | 0.040 | 0.972 | 0.07 | 2.039 | 0 | 0.18 |
| 2000 | 0.8 | 0.035 | 0.980 | 0.07 | 2.613 | 0 | 0.28 | 0.038 | 0.974 | 0.07 | 2.479 | 0 | 0.17 |

Table 4: Inference of CATE$(-2, 2|w)$ in continuous outcome model (iv). The columns indexed with "MAE", "COV" and "SE" report the median absolute errors of $\widehat{\mathrm{CATE}}(-2, 2|w)$, the empirical coverages of the confidence intervals and the average of estimated standard errors of the point estimators, respectively. The columns indexed with "SpotIV" and "TSHT" correspond to the proposed method and the method proposed in (Guo et al. 2018), respectively.

## 7.1 Construction of factor IVs

There are two main challenges of directly using all polymorphic markers as candidate instruments: a large number of polymorphic markers and the high correlation among some polymorphic markers (Bush and Moore 2012). To address these challenges, we propose a two-step procedure to construct the candidate IVs. Taking the HDL exposure as an example. In the first step, we select polymorphic markers which have "not-too-small" marginal associations with HDL. Specifically, for a given SNP, we regress the HDL level on this SNP and two measured covariates and select all the polymorphic markers with corresponding $p$-value $< 10^{-3}$. For HDL, we select 2514 polymorphic markers and form a matrix $Z^o$ with columns corresponding to the selected 2514 polymorphic markers. In the second step, we use the leading principal components of $Z^o$ as factor IVs by running the PCA analysis. This idea is closely related using factor models for the IV-exposure relationship (Bai and Ng 2010), which has demonstrated the benefits of strengthening the IVs when having many candidate IVs at hand. Let $Z^o = UDV^\mathsf{T}$ be the singular value decomposition of $Z^o$, where $D$ is a diagonal matrix containing singular values of $Z^o$. Since some columns of $Z^o$ are highly correlated, the singular values can decay to zero fast. We select the top $J^*$ principal components such that at least $90\%$ of the variance is maintained, that is,

$$\widehat{\mathcal{I}}(0.9) = \{1 \leq j \leq J^*\}, \ \ \text{where } J^* = \min\left\{1 \leq J \leq 2514 : \sum_{j=1}^{J} D_{j,j}^2 / \sum_{j=1}^{2514} D_{j,j}^2 \geq 0.9\right\}.$$

We then construct IVs based on the selected principal components as $Z = Z^o V_{,\widehat{\mathcal{I}}(0.9)}$, where $V$ is the right orthogonal matrix defined via the SVD of $Z^o$. For HDL, the number of principal components selected is 24. A plot of the cumulative proportion of explained variance is given in Section C.2 of the supplementary material. For LDL and Triglycerides exposures, we perform the same pre-processing steps to construct the candidate IVs and obtain 18 and 14 candidate IVs, respectively.

## 7.2 CATE of lipids

We study the CATE of three different lipid levels (HDL, LDL, and Triglycerides) on the highness of fasting glucose levels. We apply the proposed SpotIV method and include the Valid-CF method as a comparison. The exposures are standardized in the analysis. In Figure 4, we report estimated

CATE$(d, 0|w_F)$ and CATE$(d, 0|w_M)$, where $w_F$ and $w_M$ are the sample averages of the measured covariates for female and male mice, respectively. We consider $d' = 0$ and $d$ ranges from the 20% quantile to the 80% quantile of the standardized exposure.

For the HDL and LDL exposures, both methods give estimates of CATE close to zero at different levels of $d$. This indicates null CATEs of HDL and LDL on the fasting glucose levels. The proposed SpotIV method produces wider confidence intervals because the adjustment to possibly invalid IVs introduces more uncertainty. For Triglycerides, both methods show an increasing pattern of CATE with a larger $d$. This indicates that increased Triglycerides level can cause increased glucose levels at given levels of baseline covariates. One can see that the slope of the estimated CATE functions is larger with SpotIV than with Valid-CF.



Figure 4: The constructed 95% CIs for CATE$(d, 0|w_M)$ and CATE$(d, 0|w_F)$ with HDL, LDL, and Triglycerides exposures at different levels of $d$. The first and third columns report the results given by SpotIV and Valid-CF for CATE$(d, 0|w_M)$, respectively. The second and fourth columns report the results given by SpotIV and Valid-CF for CATE$(d, 0|w_F)$, respectively.

Because the number of candidate IVs are relatively large in this application, the uncertainty in the estimated causal effect is relatively high. To reduce the uncertainty in the estimated causal effect, we also consider the causal estimand

$$\text{CATE}(d, d'|x) = \int \mathbb{E}[y_i^{(d)} - y_i^{(d')}|z_i = \tilde{z}, x_i = x, v_i = \tilde{v}]f_{z,v}(\tilde{z}, \tilde{v}|x_i = x)d(\tilde{z}, \tilde{v}), \quad (35)$$

where $f_{z,v}$ denotes the joint density of candidate IVs and the control variable conditioning on the baseline covariates. That is, the effects of candidate IVs are marginalized out by conditioning on the baseline covariates (age and gender). In Figure 6 in the supplement, we report estimated $\text{CATE}(d, 0|x_F)$ and $\text{CATE}(d, 0|x_M)$, where $x_F$ and $x_M$ are the sample averages of the baseline covariates for female and male mice, respectively. The results are similar to the results in Figure 4 but with narrower confidence intervals.

## 8 Conclusion and Discussion

This work develops a robust causal inference framework for nonlinear outcome models in the presence of unmeasured confounders. In the semi-parametric potential outcome model, we propose new identifiability conditions to identify CATE, which weaken the classical identifiability conditions and better accommodate for the practical applications. The focus of the current work is on the inference of $\text{CATE}(d, d'|w)$ while other causal estimands of interest include the average treatment effect and $\text{CATE}(d, d'|x)$ defined in (35), which are left for future research.

## Acknowledgement

## SUPPLEMENTARY MATERIAL

Supplement to "Causal Inference for Nonlinear Outcome Models with Possibly Invalid Instrumental Variables". In the Supplementary Materials, we provide the proofs of all the theoretical results and more results on simulations and data applications.

# References

Bai, J. and S. Ng (2010). Instrumental variable estimation in a data rich environment. *Econometric Theory*, 1577–1606.

Bennett, G. (1962). Probability inequalities for the sum of independent random variables. *Journal of the American Statistical Association 57*(297), 33–45.

Berzuini, C., H. Guo, S. Burgess, and L. Bernardinelli (2020). A bayesian approach to mendelian randomization with multiple pleiotropic variants. *Biostatistics 21*(1), 86–101.

Blundell, R. W. and J. L. Powell (2004). Endogeneity in semiparametric binary response models. *The Review of Economic Studies 71*(3), 655–679.

Bowden, J., G. Davey Smith, and S. Burgess (2015). Mendelian randomization with invalid instruments: effect estimation and bias detection through egger regression. *International journal of epidemiology 44*(2), 512–525.

Bowden, J., G. Davey Smith, P. C. Haycock, and S. Burgess (2016). Consistent estimation in mendelian randomization with some invalid instruments using a weighted median estimator. *Genetic epidemiology 40*(4), 304–314.

Bush, W. S. and J. H. Moore (2012). Genome-wide association studies. *PLoS computational biology 8*(12).

Cai, B., D. S. Small, and T. R. T. Have (2011). Two-stage instrumental variable methods for estimating the causal odds ratio: Analysis of bias. *Statistics in medicine 30*(15), 1809–1824.

Chiaromonte, F., R. D. Cook, and B. Li (2002). Sufficient dimensions reduction in regressions with categorical predictors. *The Annals of Statistics 30*(2), 475–497.

Clarke, P. S. and F. Windmeijer (2012). Instrumental variable estimators for binary outcomes. *Journal of the American Statistical Association 107*(500), 1638–1652.

Cook, R. D. (2009). *Regression graphics: Ideas for studying regressions through graphics*, Volume 482. John Wiley & Sons.

Cook, R. D. and H. Lee (1999). Dimension reduction in binary response regression. *Journal of the American Statistical Association 94*(448), 1187–1200.

Cook, R. D. and B. Li (2002). Dimension reduction for conditional mean in regression. *The Annals of Statistics 30*(2), 455–474.

Davey Smith, G. and S. Ebrahim (2003). Mendelian randomization: can genetic epidemiology contribute to understanding environmental determinants of disease? *International journal of epidemiology 32*(1), 1–22.

Davey Smith, G. and G. Hemani (2014). Mendelian randomization: genetic anchors for causal inference in epidemiological studies. *Human molecular genetics 23*(R1), R89–R98.

Fajardo, R. J., L. Karim, V. I. Calley, and M. L. Bouxsein (2014). A review of rodent models of type 2 diabetic skeletal fragility. *Journal of Bone and Mineral Research 29*(5), 1025–1040.

Guo, Z., H. Kang, T. T. Cai, and D. S. Small (2018). Confidence intervals for causal effects with invalid instruments by using two-stage hard thresholding with voting. *Journal of the Royal Statistical Society: Series B (Statistical Methodology) 80*(4), 793–815.

Guo, Z. and D. S. Small (2016). Control function instrumental variable estimation of nonlinear causal effect models. *The Journal of Machine Learning Research 17*(1), 3448–3482.

Hartwig, F. P., G. Davey Smith, and J. Bowden (2017). Robust inference in summary data mendelian randomization via the zero modal pleiotropy assumption. *International journal of epidemiology 46*(6), 1985–1998.

Hayfield, T. and J. S. Racine (2008). Nonparametric econometrics: The np package. *Journal of Statistical Software 27*(5).

Ichimura, H. (1993). Semiparametric least squares (sls) and weighted sls estimation of single-index models. *Journal of Econometrics 58*(1-2), 71–120.

Imbens, G. W. and D. B. Rubin (2015). *Causal inference in statistics, social, and biomedical sciences*. Cambridge University Press.

Islam, M. S. and T. du Loots (2009). Experimental rodent models of type 2 diabetes: a review. *Methods and findings in experimental and clinical pharmacology 31*(4), 249–261.

Kang, H., A. Zhang, T. T. Cai, and D. S. Small (2016). Instrumental variables estimation with some invalid instruments and its application to mendelian randomization. *Journal of the American Statistical Association 111*(513), 132–144.

King, A. J. (2012). The use of animal models in diabetes research. *British journal of pharmacology 166*(3), 877–894.

Klein, R. W. and R. H. Spady (1993). An efficient semiparametric estimator for binary response models. *Econometrica: Journal of the Econometric Society*, 387–421.

Kolesár, M., R. Chetty, J. Friedman, E. Glaeser, and G. W. Imbens (2015). Identification and inference with many invalid instruments. *Journal of Business & Economic Statistics 33*(4), 474–484.

Lawlor, D. A., R. M. Harbord, J. A. Sterne, et al. (2008). Mendelian randomization: using genes as instruments for making causal inferences in epidemiology. *Statistics in medicine 27*(8), 1133–1163.

Li, K.-C. (1991). Sliced inverse regression for dimension reduction. *Journal of the American Statistical Association 86*(414), 316–327.

Li, S. (2017). Mendelian randomization when many instruments are invalid: hierarchical empirical bayes estimation. *arXiv preprint arXiv:1706.01389*.

Linton, O. and J. P. Nielsen (1995). A kernel method of estimating structured nonparametric regression based on marginal integration. *Biometrika*, 93–100.

Ma, Y. and L. Zhu (2012). A semiparametric approach to dimension reduction. *Journal of the American Statistical Association 107*(497), 168–179.

Newey, W. K. (1994). Kernel estimation of partial means and a general variance estimator. *Econometric Theory 10*(2), 1–21.

Neyman, J. S. (1923). On the application of probability theory to agricultural experiments. essay on principles. *Annals of Agricultural Sciences 10*, 1–51.

Petrin, A. and K. Train (2010). A control function approach to endogeneity in consumer choice models. *Journal of marketing research 47*(1), 3–13.

Rivers, D. and Q. H. Vuong (1988). Limited information estimators and exogeneity tests for simultaneous probit models. *Journal of econometrics 39*(3), 347–366.

Rosenbaum, P. R. and D. B. Rubin (1983). The central role of the propensity score in observational studies for causal effects. *Biometrika 70*(1), 41–55.

Rothe, C. (2009). Semiparametric estimation of binary response models with endogenous regressors. *Journal of Econometrics 153*(1), 51–64.

Rubin, D. B. (1974). Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of educational Psychology 66*(5), 688.

Shapland, C. Y., Q. Zhao, and J. Bowden (2020). Profile-likelihood bayesian model averaging for two-sample summary data mendelian randomization in the presence of horizontal pleiotropy.

*bioRxiv*.

Spiller, W., D. Slichter, J. Bowden, and G. Davey Smith (2019). Detecting and correcting for bias in mendelian randomization analyses using gene-by-environment interactions. *International journal of epidemiology 48*(3), 702–712.

Tchetgen, E. J. T., B. Sun, and S. Walter (2019). The genius approach to robust mendelian randomization inference. *arXiv preprint arXiv:1709.07779*.

Thompson, J. R., C. Minelli, J. Bowden, F. M. Del Greco, D. Gill, E. M. Jones, C. Y. Shapland, and N. A. Sheehan (2017). Mendelian randomization incorporating uncertainty about pleiotropy. *Statistics in Medicine 36*(29), 4627–4645.

Tsybakov, A. B. (2008). *Introduction to nonparametric estimation*. Springer Science & Business Media.

Vansteelandt, S., J. Bowden, M. Babanezhad, and E. Goetghebeur (2011). On instrumental variables estimation of causal odds ratios. *Statistical Science 26*(3), 403–422.

Verbanck, M., C.-y. Chen, B. Neale, and R. Do (2018). Detection of widespread horizontal pleiotropy in causal relationships inferred from mendelian randomization between complex traits and diseases. *Nature genetics 50*(5), 693–698.

Vershynin, R. (2010). Introduction to the non-asymptotic analysis of random matrices. *arXiv preprint arXiv:1011.3027*.

Voight, B. F., G. M. Peloso, M. Orho-Melander, et al. (2012). Plasma hdl cholesterol and risk of myocardial infarction: a mendelian randomisation study. *The Lancet 380*(9841), 572–580.

Wasserman, L. (2006). *All of nonparametric statistics*. Springer Science & Business Media.

Windmeijer, F., H. Farbmacher, N. Davies, and G. Davey Smith (2019). On the use of the lasso for instrumental variables estimation with some invalid instruments. *Journal of the American Statistical Association 114*(527), 1339–1350.

Windmeijer, F., X. Liang, F. P. Hartwig, and J. Bowden (2019). The confidence interval method for selecting valid instrumental variables. Technical report, Department of Economics, University of Bristol, UK.

Wooldridge, J. M. (2010). *Econometric analysis of cross section and panel data*. MIT press.

Wooldridge, J. M. (2015). Control function methods in applied econometrics. *Journal of Human Resources 50*(2), 420–445.

Xia, Y., H. Tong, W. K. Li, and L.-X. Zhu (2002). An adaptive estimation of dimension reduction space. *Journal of the Royal Statistical Society: Series B (Statistical Methodology) 64*(3), 363–410.

Zhu, L., B. Miao, and H. Peng (2006). On sliced inverse regression with high-dimensional covariates. *Journal of the American Statistical Association 101*(474), 630–643.

Zhu, L.-X. and K.-T. Fang (1996). Asymptotics for kernel estimate of sliced inverse regression. *The Annals of Statistics 24*(3), 1053–1068.

# A  Proofs

In this section we provide proofs for the theoretical results stated in the main paper and postpone the proofs of technical lemmas to Section B. We present the proofs for Propositions 3.1 and 3.2 in Sections A.1 and A.2, respectively. In Section A.3, we provide the proof for Lemma 5.1. In Section A.4, we provide sufficient conditions to verify Condition 5.3. We prove Theorem 5.1 and Corollary 5.1 in Sections A.5 and A.6, respectively.

In following proofs, $c_1, c_2, \ldots$ and $C_1, C_2, \ldots$ are positive constants which can be different at different places. For a matrix $A$, let $\dim(A)$ denote the column rank of $A$. For a sequence of random variables $X_n$, we use $X_n \xrightarrow{p} X$ and $X_n \xrightarrow{d} X$ to represent that $X_n$ converges to $X$ in probability and in distribution, respectively. For two positive sequences $a_n$ and $b_n$, $a_n \lesssim b_n$ means that $\exists C > 0$ such that $a_n \leq Cb_n$ for all $n$; $a_n \asymp b_n$ if $a_n \lesssim b_n$ and $b_n \lesssim a_n$, and $a_n \ll b_n$ if $\limsup_{n\to\infty} a_n/b_n = 0$.

## A.1  Proposition 3.1

By the definition of $\Theta$,

$$\Theta = \Theta^* T = \left( (\beta\gamma + \kappa)T_{1,1} + \eta T_{2,1} \quad \ldots \quad (\beta\gamma + \kappa)T_{1,M} + \eta T_{2,M} \right). \tag{36}$$

Hence, for $m = 1, \ldots, M$,

$$b_m = \text{Median}\left( \{\frac{\Theta_{j,m}}{\gamma_j}\}_{j\in\mathcal{S}} \right) = \text{Median}\left( \beta T_{1,m} + \{\frac{\kappa_j T_{1,m} + \eta_j T_{2,m}}{\gamma_j}\}_{j\in\mathcal{S}} \right).$$

34

Under the majority rule, for $m = 1, \ldots, M$,

$$\text{Median}\left(\{\frac{\kappa_j T_{1,m} + \eta_j T_{2,m}}{\gamma_j}\}_{j \in \mathcal{S}}\right) = 0.$$

Hence, $b_m = \beta T_m$ and

$$\Theta_{.,m} - b_m \gamma = \kappa T_{1,m} + \eta T_{2,m}.$$

As a result,

$$B = \begin{pmatrix} \beta T_{1,1} & \cdots & \beta T_{1,M} \\ \kappa T_{1,1} + \eta T_{2,1} & \cdots & \kappa T_{1,M} + \eta T_{2,M} \end{pmatrix} = B^* T. \tag{37}$$

### A.2 Proof of Proposition 3.2

Next, we show that

$$\mathbb{E}[y_i | d_i, w_i, v_i] = \mathbb{E}[y_i | (d_i, w_i^\mathsf{T}) B, v_i]$$

for $B = B^* T$. As $d_i$ is a function of $w_i$ and $v_i$, it holds that

$$\mathbb{E}[y_i | d_i, w_i, v_i] = \mathbb{E}[y_i | w_i, v_i] = \mathbb{E}[y_i | w_i^\mathsf{T} \Theta, v_i].$$

Since $\Theta^* = (\gamma, I_p) B^*$, $\Theta = (\gamma, I_p) B$. Therefore,

$$\mathbb{E}[y_i | d_i, w_i, v_i] = \mathbb{E}[y_i | w_i^\mathsf{T} \Theta, v_i] = \mathbb{E}[y_i | w_i^\mathsf{T} (\gamma, I_p) B, v_i] = \mathbb{E}[y_i | (w_i^\mathsf{T} \gamma, w_i^\mathsf{T}) B, v_i]$$

$$= \mathbb{E}[y_i | (w_i^\mathsf{T} \gamma + v_i, w_i^\mathsf{T}) B, v_i] = \mathbb{E}[y_i | (d_i, w_i^\mathsf{T}) B, v_i].$$

Therefore,

$$\mathbb{E}[y_i | d_i = d, w_i = w, v_i = v] = \int q(d\beta + w^\mathsf{T}\kappa, u_i) f_u(u_i | w^\mathsf{T}\eta, v) du_i$$

$$= \mathbb{E}[y_i | (d_i, w_i^\mathsf{T}) B = (d, w^\mathsf{T}) B, v_i = v] = g((d, w^\mathsf{T}) B, v).$$

Based on (1) and (8), it is not hard to see that for any $d_0 \in \mathbb{R}$,

$$
\begin{aligned}
\mathbb{E}\left[y_i^{(d_0)}|w_i = w, v_i = v\right] &= \mathbb{E}\left[\mathbb{E}[y_i^{(d_0)}|w_i = w, v_i = v, u_i]|w_i = w, v_i = v\right] \\
&= \int q(d_0\beta + w^\mathsf{T}\kappa, u_i)f_u(u_i|w^\mathsf{T}\eta, v)du_i \\
&= g\left((d_0, w^\mathsf{T})B, v\right).
\end{aligned}
\tag{38}
$$

## A.3 Proof of Lemma 5.1

**Proposition A.1.** *The parameter matrix $\Phi_{1:p,\cdot}$ satisfies (14).*

*Proof of Proposition A.1.* We first show that $\mathbb{E}[(w_i^\mathsf{T}, v_i)|\mathbb{P}_\mathbb{C}(w_i, v_i)]$ is linear in $\mathbb{P}_\mathbb{C}(w_i, v_i)$. Notice that if $v_i \in \mathbb{P}_\mathbb{C}(w_i, v_i)$,

$$
\mathbb{E}[v_i|\mathbb{P}_\mathbb{C}(w_i, v_i)] = v_i
$$

and if $v_i \notin \mathbb{P}_\mathbb{C}(w_i, v_i)$

$$
\mathbb{E}[v_i|\mathbb{P}_\mathbb{C}(w_i, 0)] = 0,
$$

where the last step is due to $\mathbb{E}[v_i|w_i] = 0$. Together with Condition 5.2, we arrive at $\mathbb{E}[(w_i^\mathsf{T}, v_i)|\mathbb{P}_\mathbb{C}(w_i, v_i)]$ is linear in $\mathbb{P}_\mathbb{C}(w_i, v_i)$. That is, the linearity assumption holds for all the covariates. By Proposition 2.1 in Chiaromonte et al. (2002), we know that the space spanned by the columns of $\Omega$ is the central subspace of $y_i|w_i, v_i$. Since the columns of $\Phi$ are eigenvectors of $\Omega$ corresponding to nonzero eigenvalues, $\Phi = (\phi_1, \ldots, \phi_{M_\Omega})$ spans the central subspace of $y_i|w_i, v_i$, which is $\mathbb{C}$. That is,

$$
\mathbb{E}[y_i|w_i, v_i] = \mathbb{E}[y_i|(w_i^\mathsf{T}, v_i)^\mathsf{T}\Phi].
$$

Let $e = (\mathbf{0}_p^\mathsf{T}, 1)^\mathsf{T}$. Let $P_e$ be the projection onto the linear space of $e$ and $P_e^\perp = I_p - P_e$. By some simple algebra,

$$
\begin{aligned}
\Phi &= P_e\Phi + P_e^\perp\Phi \\
&= \begin{pmatrix} 0 & \cdots & 0 \\ (\phi_1)_{p+1} & \cdots & (\phi_{M_\Omega})_{p+1} \end{pmatrix} + \begin{pmatrix} \phi_{1:p,1} & \cdots & \phi_{1:p,M_\Omega} \\ 0 & \cdots & 0 \end{pmatrix}.
\end{aligned}
\tag{39}
$$

Put it in another way,

$$(\phi_1, \ldots, \phi_{M_\Omega}) \subseteq \text{Span}(\Phi_{1:p,}) \oplus \text{Span}(e).$$

Hence,

$$\mathbb{E}[y_i|w_i, v_i] = \mathbb{E}[y_i|w_i^\mathsf{T}\Phi_{1:p,}, v_i].$$

On the other hand, by the definition of central subspace, we know that

$$\Phi \subseteq \text{span} \begin{pmatrix} \Theta^* & 0 \\ 0 & 1 \end{pmatrix}.$$

In view of (39), we know that

$$\text{Span}(\Phi_{1:p,.}) \subseteq \text{Span}(\Theta^*).$$

Hence, the dimension of $\text{Span}(\Phi_{1:p,.})$ is no larger than 2 and

$$\Phi_{1:p,.} = \Theta^* T$$

for some linear transformation $T$.

We first define the probabilistic limit of $\widehat{\Theta}$. Let

$$\Theta = \begin{cases} (\Phi_{1:p,i^*}, \Phi_{1:p,j^*}) & \text{if } rank(\Phi_{1:p,.}) = 2 \\ \Phi_{1:p,1} & \text{otherwise,} \end{cases} \tag{40}$$

where

$$(i^*, j^*) = \underset{1 \le i,j \le M_\Omega}{\arg\min} \; i + j$$

$$\text{subject to } |\text{cor}(\Phi_{1:p,i}, \Phi_{1:p,j})| < 1.$$

Notice that $\Theta$ in (40) is uniquely defined.

**Proposition A.2** (Convergence rate of $\widehat{\Theta}$). *Assume that Conditions 5.1 and 5.2 hold and $0 < \mathbb{P}(y_i = 1) < 1$. Then for some positive constants $c_1$ and $c_2$,*

$$\mathbb{P}\left(\|\widehat{\Theta} - \Theta\|_2 \geq c_1\sqrt{t/n}\right) \leq \exp(-c_2 t) + \mathbb{P}(E_0^c), \tag{41}$$

*$E_0$ is defined in (45) and where $\mathbb{P}(E_0) \to 1$.*

*Proof of Proposition A.2.* Notice that

$$\Omega = \Sigma^{-1/2}Cov(\alpha(y_i))\Sigma^{-1/2} = \Sigma^{-1/2}\mathbb{E}[\alpha(y_i)\alpha(y_i)^\intercal]\Sigma^{-1/2}$$

as $\mathbb{E}[\alpha(y_i)] = \mathbb{E}[(w_i^\intercal, v_i)] = 0$. The following decomposition holds

$$\|\widehat{\Omega} - \Omega\|_2 \leq 2\|\Sigma^{-1/2} - \widehat{\Sigma}^{-1/2}\|_2\|cov(\alpha(y_i))\Sigma^{-1/2}\|_2$$

$$+ \|\Sigma^{-1/2}\|_2^2\|cov(\alpha(y_i)) - \frac{1}{n}\sum_{i=1}^{n}\hat{\alpha}(y_i)\hat{\alpha}(y_i)^\intercal\|_2 + r_n, \tag{42}$$

where $r_n$ is of smaller order than the first two terms.

For the first term,

$$\|\Sigma^{-1/2} - \widehat{\Sigma}^{-1/2}\|_2 \leq \|\Sigma - \widehat{\Sigma}\|_2\|\Sigma^{1/2} + \widehat{\Sigma}^{1/2}\|_2^{-1}.$$

Since $\widehat{\Sigma}$ is an average of *i.i.d.* sub-Gaussian variables, we have

$$\mathbb{P}\left(\|\Sigma - \widehat{\Sigma}\|_2 \geq c\sqrt{t/n}\right) \leq \exp(-ct).$$

As $\|cov(\alpha(y_i))\Sigma^{-1/2}\|_2 \leq C < \infty$, for the first term in (42),

$$\mathbb{P}\left(2\|\Sigma^{-1/2} - \widehat{\Sigma}^{-1/2}\|_2\|cov(\alpha(y_i))\Sigma^{-1/2}\|_2 \geq c_1\sqrt{t/n}\right) \leq \exp(-c_2 t). \tag{43}$$

To bound the second term in (42), for binary $y_i$, it holds that

$$\alpha(1) = \mathbb{E}[(w_i^\intercal, v_i)|y_i = 1] \quad \hat{\alpha}(1) = \frac{1}{\sum_{i=1}^{n}\mathbb{1}(y_i = 0)}\sum_{i=1}^{n}(w_i^\intercal, \hat{v}_i)\mathbb{1}(y_i = 1)$$

$$\alpha(0) = \mathbb{E}[(w_i^\intercal, v_i)|y_i = 0] \quad \hat\alpha(0) = \frac{1}{\sum_{i=1}^n \mathbb{1}(y_i = 0)} \sum_{i=1}^n (w_i^\intercal, \hat v_i)\mathbb{1}(y_i = 0).$$

By some simple algebra, we can show that

$$cov(\alpha(y_i)) = \mathbb{P}(y_i = 1)\mathbb{P}(y_i = 0)(\alpha(1) - \alpha(0))(\alpha(1) - \alpha(0))^\intercal.$$

The following decomposition holds

$$\left\| \frac{1}{n}\sum_{i=1}^n \hat\alpha(y_i)\hat\alpha(y_i)^\intercal - \frac{cov(\alpha(y_i))}{\mathbb{P}(y_i = 1)\mathbb{P}(y_i = 0)} \right\|_2$$

$$\leq 2\|(\hat\alpha(1) - \hat\alpha(0) - \alpha(1) + \alpha(0))(\alpha(1) - \alpha(0))^\intercal\|_2 + \|\alpha(1) - \alpha(0) - \hat\alpha(1) + \hat\alpha(0)\|_2^2$$

$$\leq 4\|\alpha(1) - \alpha(0)\|_2 \max_{k\in\{0,1\}} \|\hat\alpha(k) - \alpha(k)\|_2 + 4 \max_{k\in\{0,1\}} \|\hat\alpha(k) - \alpha(k)\|_2^2.$$

First notice that
$$\alpha(k) = \frac{\mathbb{E}\left[(w_i^\intercal, v_i)\mathbb{1}(y_i = k)\right]}{\mathbb{P}(y_i = k)}.$$

$$\|\hat\alpha(k) - \alpha(k)\|_2 \leq \left| \frac{1}{\mathbb{P}(y_i = k)} - \frac{n}{\sum_{i=1}^n \mathbb{1}(y_i = k)} \right| |\mathbb{E}\left[(w_i^\intercal, v_i)\mathbb{1}(y_i = k)\right]|$$

$$+ \frac{1}{\mathbb{P}(y_i = k)}\left\| \frac{1}{n}\sum_{i=1}^n (w_i^\intercal, \hat v_i)\mathbb{1}(y_i = k) - \mathbb{E}\left[(w_i^\intercal, v_i)\mathbb{1}(y_i = k)\right] \right\|_2.$$

$$|\frac{1}{n}\sum_{i=1}^n (\hat v_i - v_i)\mathbb{1}(y_i = k)| = |\frac{1}{n}\sum_{i=1}^n \mathbb{1}(y_i = k)w_i^\intercal(\hat\gamma - \gamma)|$$

$$\leq \|\frac{1}{n}\sum_{i=1}^n \mathbb{1}(y_i = k)w_i^\intercal\|_2\|\hat\gamma - \gamma\|_2.$$

By Condition 5.1(a), $\mathbb{1}(y_i = k)w_i^\intercal$ are independent sub-Gaussian variables with sub-Gaussian norm no larger than the sub-Gaussian norm of $w_i$. Hence,

$$\mathbb{P}\left( |\frac{1}{n}\sum_{i=1}^n (\hat v_i - v_i)\mathbb{1}(y_i = k)| \geq c_1\sqrt{t/n} \right) \leq \exp(-c_2 t).$$

39

Moreover, $\mathbb{1}(y_i = k)$ and $w_i, v_i$ are all sub-Gaussian. Hence, it is straight forward to show that

$$\mathbb{P}\left( \left\| \frac{1}{n} \sum_{i=1}^{n} \hat{\alpha}(y_i)\hat{\alpha}(y_i)^\intercal - \frac{cov(\alpha(y_i))}{\mathbb{P}(y_i = 1)\mathbb{P}(y_i = 0)} \right\|_2 \geq c_3\sqrt{t/n} \right) \leq \exp(-c_4 t)$$

for sufficiently large constants $c_3$ and $c_4$.

In view of (42), we have shown

$$\mathbb{P}\left( \|\widehat{\Omega} - \Omega\|_2 \geq c_5\sqrt{\frac{t}{n}} \right) \leq \exp(-c_6 t) \tag{44}$$

for sufficiently large constants $c_5$ and $c_6$.

Next, we show the the eigenvalues of $\widehat{\Omega}$ converges to the eigenvalues of $\Omega$. In fact,

$$\max_{1\leq k\leq p} \left| \hat{\lambda}_k - \hat{\lambda}_k \right| \leq \max_{\|u\|_2=1} |u^\intercal(\widehat{\Omega} - \Omega)u| \leq \|\widehat{\Omega} - \Omega\|_2.$$

For the eigenvectors, we use Theorem 5 of Karoui (2008). Under Condition 5.1(b), we have

$$\|\widehat{\Phi}_{.,m} - \Phi_{.,m}\|_2 \leq \frac{\|\widehat{\Omega} - \Omega\|_2}{\lambda_m(\Omega)} \ \forall \ 1 \leq m \leq M_\Omega.$$

In view of (44), we have shown

$$\mathbb{P}\left( \max_{1\leq m\leq M_\Omega} \|\widehat{\Phi}_{.,m} - \Phi_{.,m}\|_2 \geq C_1\sqrt{\frac{t}{n}} \right) \leq \exp(-C_2 t).$$

This implies that $\widehat{\Theta}$ defined in (21) is a consistent estimator of

$$\tilde{\Theta} = \begin{cases} (\Phi_{1:p,i^*}, \Phi_{1:p,j^*}) & \text{if (20) exists,} \\ \Phi_{1:p,1} & \text{otherwise.} \end{cases}$$

Define an event

$$E_0 = \left\{ \tilde{\Theta} \text{ spans } \Phi_{1:p,} \text{ and } \widehat{M} = \dim(\Phi_{1:p,.}) \right\}. \tag{45}$$

In event $E_0$, $\Theta$ defined in (40) equals $\tilde{\Theta}$. It left to show that $\mathbb{P}(E_0) \to 1$. By Proposition A.1, we know that the dimension of $\Phi_{1:p,.}$ is at most 2. Moreover,

$$\max_{i,j \leq M_\Omega} |\langle \widehat{\Phi}_{1:p,i}, \widehat{\Phi}_{1:p,j} \rangle - \langle \Phi_{1:p,i}, \Phi_{1:p,j} \rangle| \leq \max_{i \leq M_\Omega} \|\widehat{\Phi}_{1:p,i} - \Phi_{1:p,i}\|_2 \max_{j \leq M_\Omega} \|\Phi_{1:p,j}\|_2.$$

Hence, when $|\mathrm{cor}(\phi_{1:p,i}, \phi_{1:p,j})| = 1$,

$$\mathbb{P}\left(\max_{i,j \leq M_\Omega} |\mathrm{cor}(\widehat{\Phi}_{1:p,i}, \widehat{\Phi}_{1:p,j})| \leq 1 - c\sqrt{\frac{\log n}{n}}\right)$$

$$\leq \mathbb{P}\left(\max_{i \leq M_\Omega} \|\widehat{\Phi}_{1:p,i} - \Phi_{1:p,i}\|_2 \geq c_1\sqrt{\frac{\log n}{n}}\right) \leq \exp(-c_2 \log n).$$

When $|\mathrm{cor}(\Phi_{1:p,i}, \Phi_{1:p,j})| \leq c_0 < 1$,

$$\mathbb{P}\left(\max_{i,j \leq M_\Omega} |\mathrm{cor}(\widehat{\Phi}_{1:p,i}, \widehat{\Phi}_{1:p,j})| \geq 1 - c\sqrt{\frac{\log n}{n}}\right)$$

$$\leq \mathbb{P}\left(\max_{i \leq M_\Omega} \|\widehat{\Phi}_{1:p,i} - \Phi_{1:p,i}\|_2 \geq 1 - c_0 - c_1\sqrt{\frac{\log n}{n}}\right) \leq \exp(-c_2 \log n).$$

Hence,

$$\mathbb{P}(E_0) \geq 1 - \exp(-c_3 \log n) \to 1.$$

*Proof of Lemma 5.1.* For $\widehat{\gamma}$ computed via (19), under Condition 5.1, it is easy to show that

$$\sqrt{n}(\hat{\gamma} - \gamma) \xrightarrow{D} N\left(0, \sigma_v^2 \mathbb{E}^{-1}[w_i w_i^\intercal]\right). \tag{46}$$

Define an event

$$E_1 = \left\{\widehat{\mathcal{S}} = \mathcal{S} \text{ and } (50) \text{ holds}\right\} \cap E_0. \tag{47}$$

Let $\widehat{\omega}_j = \hat{\sigma}_v^2 \{\widehat{\Sigma}^{-1}\}_{j,j}$. It is easy to show

$$|\hat{\sigma}_v^2 - \sigma_v^2| = O_P(n^{-1/2}).$$

41

We first show that $\mathbb{P}(E_1) \to 1$ as $n \to \infty$. For $j \in \mathcal{S}$, we have

$$\mathbb{P}\left(|\widehat{\gamma}_j| \geq \sqrt{\widehat{\omega}_j}\sqrt{\frac{2.01 \log n}{n}}\right) \geq \mathbb{P}\left(|\gamma_j| - |\widehat{\gamma}_j - \gamma_j| \geq \sqrt{\widehat{\omega}_j}\sqrt{\frac{2.01 \log n}{n}}\right)$$

$$= \mathbb{P}\left(|\widehat{\gamma}_j - \gamma_j| \leq |\gamma_j| - \sqrt{\widehat{\omega}_j}\sqrt{\frac{2.01 \log n}{n}}\right) \to 1,$$

where the convergence follows from (46) and $|\gamma_j| \geq c_0 > 0$ for $j \in \mathcal{S}$. For $j \in \mathcal{S}^c$, we have

$$\mathbb{P}\left(|\widehat{\gamma}_j| > \sqrt{\widehat{\omega}_j}\sqrt{\frac{2.01 \log n}{n}}\right)$$

$$= \mathbb{P}\left(|\widehat{\gamma}_j - \gamma_j| > \sqrt{\widehat{\omega}_j}\sqrt{\frac{2.01 \log n}{n}}\right) = o(1),$$

where the last step is due to $\|\hat{\gamma} - \gamma\|_2 = O_P(n^{-1/2})$. Combining above two expressions, we have establish that

$$\mathbf{P}\left(\widehat{\mathcal{S}} = \mathcal{S}\right) \to 1. \tag{48}$$

It suffices to prove the rest of the results conditioning on the event $\{\widehat{\mathcal{S}} = \mathcal{S}\}$. By the sub-Gaussian property of observed data,

$$\mathbb{P}\left(\max_{j \in \mathcal{S}}\left|\frac{\widehat{\Theta}_{j,m}}{\hat{\gamma}_j} - \frac{\Theta_{j,m}}{\gamma_j}\right| \geq c_1\sqrt{\frac{t}{n}}\right) \leq \exp(-c_2 t) \tag{49}$$

for some positive constants $c_1$ and $c_2$. We have shown in Proposition 3.1 that for $j \in \mathcal{V}$, $\frac{\Theta_{j,m}}{\gamma_j} = b_m$ and for $j \notin \mathcal{V}$,

$$\frac{\Theta_{j,m}}{\gamma_j} = \beta_m + \frac{\kappa_j T_{1,m} + \eta_j T_{2,m}}{\gamma_j}.$$

Notice that for $j \notin \mathcal{V}$, it is possible that $\frac{\Theta_{j,m}}{\gamma_j} = b_m$. It suffices to show that for

$$\left\{\frac{\widehat{\Theta}_{k,m}}{\hat{\gamma}_k} \geq \max_{j \in \mathcal{V}}\frac{\widehat{\Theta}_{j,m}}{\hat{\gamma}_j} \text{ or } \frac{\widehat{\Theta}_{k,m}}{\hat{\gamma}_k} > \min_{j \in \mathcal{V}}\frac{\widehat{\Theta}_{j,m}}{\hat{\gamma}_j}, \forall k \text{ such that} \frac{\Theta_{k,m}}{\gamma_k} \neq b_m, \forall 1 \leq m \leq M\right\}, \tag{50}$$

$\mathbb{P}(50 \text{ holds}) \rightarrow 1$. That is, any $\frac{\widehat{\Theta}_{k,m}}{\hat{\gamma}_k}$ cannot be the median if $\frac{\Theta_{k,m}}{\gamma_k} \neq b_m$. (50) can be proved by noticing that

$$\max_{j \in \mathcal{S}} |\frac{\widehat{\Theta}_{j,m}}{\hat{\gamma}_j} - \frac{\Theta_{j,m}}{\gamma_j}| = O_P(|\mathcal{S}|n^{-1/2}) = o_P(1).$$

If $\frac{\Theta_{k,m}}{\gamma_k} - b_m > 0$, then

$$\mathbb{P}\left(\frac{\widehat{\Theta}_{k,m}}{\hat{\gamma}_k} > \max_{j \in \mathcal{V}} \frac{\widehat{\Theta}_{k,m}}{\hat{\gamma}_j}\right) \geq \mathbb{P}\left(\frac{\Theta_{k,m}}{\gamma_k} - b_m > Cn^{-1/2}\right) \rightarrow 1$$

for some constant $C > 0$. If If $\frac{\Theta_{k,m}}{\gamma_k} - b_m < 0$, then

$$\mathbb{P}\left(\frac{\widehat{\Theta}_{k,m}}{\hat{\gamma}_k} < \max_{j \in \mathcal{V}} \frac{\widehat{\Theta}_{j,m}}{\hat{\gamma}_j}\right) \geq \mathbb{P}\left(\frac{\Theta_{k,m}}{\gamma_k} - b_m < Cn^{-1/2}\right) \rightarrow 1$$

for some constant $C > 0$. Hence, (50) holds. We have shown $\mathbb{P}(E_1) \rightarrow 1$. In event $E_1$, by (49),

$$\mathbb{P}\left(\|\widehat{B} - B\|_2 \geq c_1 t | E_1\right) \leq \exp(-c_2 t). \tag{51}$$

for some large enough constants $c_1$ and $c_2$. The results of Lemma 5.1 hold in view of (50).

## A.4 Verification of Condition 5.3

We provide some generic examples of $f_t$ and $q(\cdot)$ such that Condition 5.3 holds when $M = 2$. Proposition A.3 provides a sufficient condition for Condition 5.3 (a) and (c). Proposition A.4 provides a sufficient condition for Condition 5.3 (b) when $u_i$ has support $\mathbb{R}$. Proposition A.5 provides a sufficient condition for Condition 5.3 (b) when $h$ is an indicator function.

Let $t_i^* = ((d_i, w_i^\mathsf{T})B^*, v_i)$ and $s_i^* = ((d, w^\mathsf{T})B^*, v_i)$. Let $f_{t^*}$ denote the density of $t_i^*$. We use $\mathcal{T}^*$ and $\mathcal{T}_v$ to denote the support of the density functions $f_{t^*}$ and $f_v$, respectively. For a set $\mathcal{T}$, we use $\mathcal{T}^{\text{int}}$ to denote its interior.

**Proposition A.3** (A sufficient condition for Condition 5.3 (a) and (c))**.** *Suppose that the support of $t_i^*$ is $\mathcal{T}^* = [-a_1, a_1] \times [-a_2, a_2] \times [-a_3, a_3]$ and $\int_{t^* \in (\mathcal{T}^*)^{\text{int}}} f_t(t)dt = 1$, where $a_1, a_2 > 0$ can be*

$\infty$ *and* $|a_3| \leq C < \infty$. *Suppose that the density* $f_{t^*}$ *satisfies*

$$c_1 \leq \inf_{x \in \mathcal{T}_v^{\text{int}}} f_{t^*}((d, w^\intercal)B^*, x) \leq \sup_{x \in \mathcal{T}_v^{\text{int}}} f_{t^*}((d, w^\intercal)B^*, x) \leq C_1$$

*for some constants* $c_1$ *and* $C_1$. *Moreover, we assume that* $f_{t^*}(\tilde{t})$ *is differentiable and Lipschitz in* $\mathcal{T}^*$ *and* $f_v(v)$ *uniformly bounded in* $\mathcal{T}_v$.

*For any* $u \in \{(d, w^\intercal)B^*_{\cdot,1} \pm Ch\} \times \{(d, w^\intercal)B^*_{\cdot,2} \pm Ch\}$ *with some sufficiently large constant* $C$, *it holds that* $|u_1| < a_1$ *and* $|u_2| < a_2$. *Then Condition 5.3 (a) and (c) hold true.*

*Proof of Proposition A.3.* We first verify Condition 5.3 (a). As $M = 2$, $T$ is invertible. Because $T$ is a $2 \times 2$ constant matrix, we know that $c \leq |T^{-1}| < C < \infty$. Hence, by the linear transformation of density

$$f_t(\tilde{t}) = f_{t^*}(\tilde{t}_{1:2}T^{-1}, \tilde{t}_3)|T^{-1}|. \tag{52}$$

As $\mathcal{T}^*$ defined in (A.3) is convex, above expression implies that $\mathcal{T}$ is also convex, no matter $a_1, a_2 = \infty$ or not. Moreover,

$$\min_i f_t(s_i) \geq \inf_{x \in \mathcal{T}_v^{\text{int}}} f_t((d, w^\intercal)B, x) = \inf_{x \in \mathcal{T}_v^{\text{int}}} f_{t^*}(s_i^*)|T^{-1}| \geq c_0 > 0$$

for some constant $c_0 > 0$. Similarly, one can show that

$$\sup_{x \in \mathcal{T}_v^{\text{int}}} f_t((d, w^\intercal)B, x) \leq C_0 < \infty.$$

For the derivative of $f_t$, by (52),

$$\max_{1 \leq i \leq n} \sup_{t_0 \in \mathcal{N}_h(s_i) \cap \mathcal{T}} \|\nabla f_t(t_0)\|_\infty = \max_{1 \leq i \leq n} \sup_{t_0 \in \mathcal{N}_h(s_i) \cap \mathcal{T}} \left\| \nabla f_{t^*}((t_0)_{1:2}T^{-1}, (t_0)_3) \begin{pmatrix} T^{-1} & 0 \\ 0 & 1 \end{pmatrix} \right\|_\infty |T^{-1}|$$

$$\leq \left\| \sup_{\tilde{t}^*_{1:2} \in (d, w^\intercal)B^* \pm hT^{-1}, \tilde{t}^*_3 \in \mathcal{T}_v} \frac{\partial f_{t^*}(\tilde{t}^*)}{\partial \tilde{t}^*} \right\|_\infty (1 + 2\|T^{-1}\|_{\max})|T^{-1}|.$$

As $T^{-1}$ has bounded norms, the interval $(d, w^\intercal)B^* \pm hT^{-1}$ is inside $\{(d, w^\intercal)B^*_{\cdot,1} \pm Ch\} \times \{(d, w^\intercal)B^*_{\cdot,2} \pm Ch\}$. As $\{(d, w^\intercal)B^*_{\cdot,1} \pm Ch\} \times \{(d, w^\intercal)B^*_{\cdot,2} \pm Ch\}$ is a subset of $[-a_1, a_1] \times [-a_2, a_2]$

and $f_{t^*}$ is differentiable and Lipschitz in $\mathcal{T}^*$, we have

$$\max_{1 \leq i \leq n} \sup_{t_0 \in \mathcal{N}_h(s_i) \cap \mathcal{T}} \|\nabla f_t(t_0)\|_\infty \leq C_3 < \infty.$$

The convexity of $\mathcal{T}_v = [-a_3, a_3]$ is obvious.

For Condition 5.3 (c), since the evaluation point $|(d, w^{\mathsf{T}})B^*_{\cdot,1} \pm Ch| \leq a_1$ and $|(d, w^{\mathsf{T}})B^*_{\cdot,2} \pm Ch| \leq a_2$, we know that $((d, w^{\mathsf{T}})B + \Delta^{\mathsf{T}}, v)^{\mathsf{T}} \in \mathcal{T}$ for any $\Delta \in \mathbb{R}^2$ satisfying $\|\Delta\|_\infty \leq h$ and for any $v \in \mathcal{T}_v$.

**Proposition A.4** (A sufficient condition for Condition 5.3 (b)). *Assume that $v_i$ has a compact support $\mathcal{T}_v$. The function $q(\cdot, \cdot) : \mathbb{R}^2 \to [0, 1]$ is twice differentiable and its first two derivatives are uniformly bounded. The random variable $q(d\beta + w^{\mathsf{T}}\kappa, u_i)$ is away from zero and one at some point $u_0$ such that $f(u_0|w_i^{\mathsf{T}}\eta, v_i) > 0$ for any $v_i \in \mathcal{T}_v$. Moreover, assume that the conditional density $f_u(u|z^{\mathsf{T}}\eta, v)$ comes from a location-scale family such that*

$$f_u(u|w^{\mathsf{T}}\eta, v) = \frac{1}{\sigma(w^{\mathsf{T}}\eta, v)} f_0\left(\frac{u - \mu(w^{\mathsf{T}}\eta, v)}{\sigma(w^{\mathsf{T}}\eta, v)}\right),$$

*where $f_0$, $\mu(w^{\mathsf{T}}\eta, v) = \mathbb{E}[u|w^{\mathsf{T}}\eta, v]$, and $\sigma^2(w^{\mathsf{T}}\eta, v) = Var(u|w^{\mathsf{T}}\eta, v)$ are all twice differentiable and their first two derivatives are uniformly bounded. Then Condition 5.3 (b) holds true.*

*Proof of Proposition A.4.* We first show that $g(s_i)$ is uniformly bounded away from zero and one. By (3) and (8),

$$g(s_i) = \mathbb{E}[y_i|d_i = d, w_i = w, v_i = v_i] = \int q(d\beta + w^{\mathsf{T}}\kappa, u_i) f_u(u_i|w^{\mathsf{T}}\eta, v_i) du_i.$$

Since $q(d\beta + w^{\mathsf{T}}\eta, u_i)$ is Lipschitz in $u_i$,

$$|q(d\beta + w^{\mathsf{T}}\eta, u_i) - q(d\beta + w^{\mathsf{T}}\eta, u_0)| \leq C|u_i - u_0|$$

for some constant $C$. Hence, for any

$$|u_i - u_0| \leq \frac{1 - q(d\beta + w^{\mathsf{T}}\eta, u_0)}{2C},$$

45

$$q(d\beta + w^\intercal\eta, u_i) \leq q(d\beta + w^\intercal\eta, u_0) + \frac{1 - q(d\beta + w^\intercal\eta, u_0)}{2} \leq c_1 < 1. \tag{53}$$

Therefore,

$$\int q(d\beta + w^\intercal\kappa, u_i) f_u(u_i|w^\intercal\eta, v_i) du_i \leq \int_{|u_i - u_0| > \frac{1 - q(d\beta + w^\intercal\eta, u_0)}{2C}} f_u(u_i|w^\intercal\eta, v_i) du_i$$

$$+ c_1 \int_{|u_i - u_0| \leq \frac{1 - q(d\beta + w^\intercal\eta, u_0)}{2C}} f_u(u_i|w^\intercal\eta, v_i) du_i,$$

where the last step is due to $q(\cdot) \leq 1$ and (53).

Because $f_u(u_0|w^\intercal\eta, v_i) > 0 \ \forall v_i \in \mathcal{T}_v$ and $\mathcal{T}_v$ is compact, there exists a constant $c_0$ such that $f_u(u_0|w^\intercal\eta, v_i) \geq c_0 > 0 \ \forall v \in \mathcal{T}_v$. Using the Lipschitz property of $f_u(u_i|w^\intercal\eta, v_i)$ in $u_i$, it is easy to show that

$$\int_{|u_i - u_0| \leq \frac{1 - q(d\beta + w^\intercal\eta, u_0)}{2C}} f_u(u_i|w^\intercal\eta, v_i) du_i \geq c_2 > 0$$

and hence

$$g(s_i) = \int q(d\beta + w^\intercal\kappa, u_i) f_u(u_i|w^\intercal\eta, v_i) du_i$$

$$\leq 1 - \int_{|u_i - u_0| \leq \frac{1 - q(d\beta + w^\intercal\eta, u_0)}{2C}} f_u(u_i|w^\intercal\eta, v_i) du_i + c_1 \int_{|u_i - u_0| \leq \frac{1 - q(d\beta + w^\intercal\eta, u_0)}{2C}} f_u(u_i|w^\intercal\eta, v_i) du_i$$

$$\leq 1 - (1 - c_1)c_2 < 1$$

uniformly in $v_i$. Similarly one can show that $g(s_i)$ is bounded away from zero uniformly in $s_i$.

Next, we show the Lipschitz property of $g$. Let $s_i^* = ((d, w^\intercal)B^*, v_i)^\intercal$. We first show the Lipschitz property of $g$ at $s_i$ is implied by the Lipschitz property of $g^*$ and $s_i^*$. For $M = 2$, $T$ is invertible and

$$\frac{\partial g(s_i)}{\partial s_i} = \frac{\partial g^*(s_i^*)}{\partial s_i} = \frac{\partial g^*(s_i^*)}{\partial s_i^*}\frac{\partial s_i^*}{\partial s_i} = \frac{\partial g^*(s_i^*)}{\partial s_i^*}T^{-1}.$$

As the columns of $B$ and $B^*$ are normalized,

$$\|\frac{\partial s_i^*}{\partial s_i}\|_2 \leq C < \infty.$$

Same arguments hold for $\partial g(s_i)/\partial(s_i)_2$. Using the above arguments, we arrive at

$$\|\frac{\partial g(s_i)}{\partial s_i}\|_2 \leq \|\frac{\partial g^*(s_i^*)}{\partial s_i}\|_2 C$$

for some constant $C > 0$.

We are left to establish the Lipschitz property of $g^*$ at $s_i^*$. Notice that $q((s_i^*)_1, u_i) f_u(u_i|(s_i^*)_2, (s_i^*)_3)$ is Lebesgue-integrable because $q(\cdot, \cdot) \in [0, 1]$ and $f_u(\cdot)$ is a density function. In addition, $\sup_{x \in \mathbb{R}^2} |q'(x)| \leq C < \infty$ and $C f_u(u_i|(s_i^*)_2, (s_i^*)_3)$ is Lebesgue-integrable with respect to $u_i$. Hence, we change the order of differentiation and integration to get that

$$\frac{\partial g^*(s_i^*)}{\partial (s_i^*)_1} = \int q'((s_i^*)_1, u_i) f_u(u_i|(s_i^*)_2, (s_i^*)_3) du_i$$

and hence

$$\sup_{s_i^*} |\frac{\partial g^*(s_i^*)}{\partial (s_i^*)_1}| \leq C < \infty.$$

Similarly, we can show that

$$\sup_{s_i^*} \left| \frac{\partial^2 g^*(s_i^*)}{\{\partial(s_i^*)_1\}^2} \right| \leq C < \infty.$$

For the partial derivatives with respect to $((s_i^*)_2, (s_i^*)_3)$, by our assumption on $f_u(u|w^\intercal \eta, v)$, we can use change of variable to arrive at

$$g(s_i^*) = \int q((s_i^*)_1, \sigma_i x + \mu_i) f_0(x) dx,$$

where $\mu(w_i^\intercal \eta, v_i)$ is abbreviated as $\mu_i$ and $\sigma(w_i^\intercal \eta, v_i)$ is abbreviated as $\sigma_i$, and

$$\int f_0(x) dx = \int f_u(u|w_i^\intercal \eta, v_i) du = 1.$$

Using similar arguments as above, the conditions of Proposition A.3 imply that

$$|\frac{\partial g^*(s_i^*)}{\partial (w_i^\intercal \eta, v_i)}| \leq C \int (|x| + 1) f_0(x) dx \leq C' < \infty.$$

As a result, we can change the order of differentiation and integration to get

$$\sup_{s_i^*} \|\frac{\partial g^*(s_i^*)}{\partial(w_i^\mathsf{T}\eta, v_i)}\|_2 \leq C' < \infty.$$

Similarly, we can show that

$$\sup_{s_i^*} \left\|\frac{\partial^2 g^*(s_i^*)}{\{\partial(w_i^\mathsf{T}\eta, v_i)\}^{\otimes 2}}\right\|_2 \leq C'' < \infty \text{ and } \sup_{s_i^*} \left\|\frac{\partial^2 g^*(s_i^*)}{\partial(s_i^*)_1 \partial(w_i^\mathsf{T}\eta, v_i)}\right\|_2 \leq C'' < \infty.$$

**Proposition A.5** (Second sufficient condition for Condition 5.3 (b)). *Assume that $v_i$ has a compact support $\mathcal{T}_v$ and*

$$q(d\beta + w^\mathsf{T}\kappa, u_i) = \mathbb{1}(d\beta + w^\mathsf{T}\kappa + u_i \geq c)$$

*for fixed some constant $c$. Then*

$$g^*(s_i^*) = \mathbb{P}(u_i \geq c - d\beta - w^\mathsf{T}\kappa | w_i^\mathsf{T}\eta = w^\mathsf{T}\eta, v_i = v).$$

*If $g^*$ satisfies Condition 5.3(b), then $g$ satisfies Condition 5.3(b).*

*Proof of Proposition A.5.* The proof is obvious and is omitted here.

### A.5 Proof of Theorem 5.1

It follows from the condition $h = n^{-c}$ for $0 < c < 1/4$ that $nh^4 \gg \log n$ and $h \log n \to 0$. We recall the following definitions,

$$t_i = ((d_i, w_i^\mathsf{T})B, v_i)^\mathsf{T}, \quad \widehat{t}_i = ((d_i, w_i^\mathsf{T})\widehat{B}, \widehat{v}_i)^\mathsf{T}, \quad s_i = ((d, w^\mathsf{T})B, v_i)^\mathsf{T}, \quad \widehat{s}_i = ((d, w^\mathsf{T})\widehat{B}, \widehat{v}_i)^\mathsf{T}.$$

Since we take $M = 2$, on the event $E_0$ defined in (45), we have $\widehat{M} = 2$. Hence, the kernel is defined in three dimensions, that is, for $a, b \in \mathbb{R}^3$,

$$K_H(a, b) = \prod_{l=1}^{3} \frac{1}{h} k\left(\frac{a_l - b_l}{h}\right)$$

where $h$ is the bandwidth and $k(x) = \mathbf{1}\,(|x| \leq 1/2)$. We define the events

$$\mathcal{A}_1 = \left\{ \|\widehat{B} - B\|_2 \leq C\sqrt{\frac{\log n}{n}}, \|\widehat{\gamma} - \gamma\|_2 \leq C\sqrt{\frac{\log n}{n}} \right\}, \quad \mathcal{A}_2 = \max\{\|w_i\|_\infty, |d_i|\} \lesssim \sqrt{\log n}$$

By Lemma 5.1 and $w_i$ and $v_i$ being sub-gaussian, we establish that $\mathbf{P}(\mathcal{A}_1 \cap \mathcal{A}_3) \geq 1 - n^{-c} - P(E_1)$. On the event $\mathcal{A}_1 \cap \mathcal{A}_2$, we have

$$\max_{1 \leq i \leq n} \max \left\{ \|\widehat{s}_i - s_i\|_2, \|\widehat{t}_i - t_i\|_2 \right\} \leq C \log n/\sqrt{n}$$

for a large positive constant $C > 0$.

We start with the decomposition

$$\widehat{\mathrm{ASF}}(d, w) - \mathrm{ASF}(d, w) = \frac{1}{n} \sum_{i=1}^{n} [\widehat{g}(\widehat{s}_i) - g(\widehat{s}_i)] + \frac{1}{n} \sum_{i=1}^{n} g(\widehat{s}_i) - \int g(s_i) f_v(v_i) dv_i \qquad (54)$$

where $f_v$ is the density of $v_i$. By (17) in the main paper, we define

$$\epsilon_i = y_i - \mathbb{E}[y_i | (d_i, w_i^\intercal)B, v_i] = y_i - g((d_i, w_i^\intercal)B, v_i) \quad \text{for } 1 \leq i \leq n. \qquad (55)$$

We plug in the expression of $\widehat{g}(\widehat{s}_i)$ and decompose the error $\frac{1}{n} \sum_{i=1}^{n} [\widehat{g}(\widehat{s}_i) - g(\widehat{s}_i)]$ as

$$
\begin{aligned}
&\frac{1}{n} \sum_{i=1}^{n} \frac{\sum_{j=1}^{n} [y_j - g(\widehat{s}_i)] K_H(\widehat{s}_i, \widehat{t}_j)}{\sum_{j=1}^{n} K_H(\widehat{s}_i, \widehat{t}_j)} = \frac{1}{n} \sum_{i=1}^{n} \frac{\sum_{j=1}^{n} \epsilon_j K_H(\widehat{s}_i, \widehat{t}_j)}{\sum_{j=1}^{n} K_H(\widehat{s}_i, \widehat{t}_j)} \\
&+ \frac{1}{n} \sum_{i=1}^{n} \frac{\sum_{j=1}^{n} [g(\widehat{t}_j) - g(\widehat{s}_i)] K_H(\widehat{s}_i, \widehat{t}_j)}{\sum_{j=1}^{n} K_H(\widehat{s}_i, \widehat{t}_j)} + \frac{1}{n} \sum_{i=1}^{n} \frac{\sum_{j=1}^{n} [g(t_j) - g(\widehat{t}_j)] K_H(\widehat{s}_i, \widehat{t}_j)}{\sum_{j=1}^{n} K_H(\widehat{s}_i, \widehat{t}_j)}.
\end{aligned}
\qquad (56)
$$

Since

$$\left| g(\widehat{t}_j) - g(t_j) \right| \cdot K_H(\widehat{s}_i, \widehat{t}_j) \leq \|\nabla g(t_j + c(\widehat{t}_j - t_j))\|_2 \|\widehat{t}_j - t_j\|_2 \cdot K_H(\widehat{s}_i, \widehat{t}_j),$$

we apply the boundedness assumption on $\nabla g$ imposed in Condition 5.3 (b) and obtain that $\left| g(\widehat{t}_j) - g(t_j) \right| \lesssim \log n/\sqrt{n}$ on the event $\mathcal{A}$. Here, we use the fact that, if $K_H(\widehat{s}_i, \widehat{t}_j) > 0$ and $C \log n/\sqrt{n} \leq h/2$, then $\|\widehat{t}_j - s_i\|_\infty \leq \|\widehat{t}_j - \widehat{s}_i\|_\infty + \|\widehat{s}_i - s_i\|_\infty \leq h$.

Hence, we have

$$\left| \frac{1}{n} \sum_{i=1}^{n} \frac{\sum_{j=1}^{n} [g(t_j) - g(\widehat{t}_j)] K_H(\widehat{s}_i, \widehat{t}_j)}{\sum_{j=1}^{n} K_H(\widehat{s}_i, \widehat{t}_j)} \right| \lesssim \log n / \sqrt{n}.$$

Then following from (54) and (56), it is sufficient to control the following terms,

$$\underbrace{\frac{1}{n} \sum_{i=1}^{n} g(\widehat{s}_i) - \int g(s_i) f_v(v_i) dv_i}_{T_1} + \underbrace{\frac{1}{n} \sum_{i=1}^{n} \frac{\sum_{j=1}^{n} \epsilon_j K_H(\widehat{s}_i, \widehat{t}_j)}{\sum_{j=1}^{n} K_H(\widehat{s}_i, \widehat{t}_j)}}_{T_2} + \underbrace{\frac{1}{n} \sum_{i=1}^{n} \frac{\sum_{j=1}^{n} [g(\widehat{t}_j) - g(\widehat{s}_i)] K_H(\widehat{s}_i, \widehat{t}_j)}{\sum_{j=1}^{n} K_H(\widehat{s}_i, \widehat{t}_j)}}_{T_3}. \tag{57}$$

We now control the three terms $T_1$, $T_2$ and $T_3$ separately.

**Control of $T_1$.** The term $T_1$ is controlled by the following lemma, whose proof is presented in Section B.1.

**Lemma A.1.** *Suppose the assumptions of Theorem 5.1 hold, then with probability larger than* $1 - n^{-c} - \frac{1}{t^2}$,

$$\left| \frac{1}{n} \sum_{i=1}^{n} g(\widehat{s}_i) - \int g(s_i) f_v(v_i) dv_i \right| \lesssim \frac{t + \log n}{\sqrt{n}} \tag{58}$$

**Control of $T_2$.** We approximate $T_2$ by $\frac{1}{n} \sum_{i=1}^{n} \frac{\frac{1}{n} \sum_{j=1}^{n} \epsilon_j K_H(s_i, t_j)}{\frac{1}{n} \sum_{j=1}^{n} K_H(s_i, t_j)}$, which can be expressed as $\frac{1}{n} \sum_{j=1}^{n} \epsilon_j a_j$ with

$$a_j = \frac{1}{n} \sum_{i=1}^{n} \frac{K_H(s_i, t_j)}{\frac{1}{n} \sum_{j=1}^{n} K_H(s_i, t_j)}. \tag{59}$$

Then the approximation error is

$$\frac{1}{n} \sum_{i=1}^{n} \frac{\frac{1}{n} \sum_{j=1}^{n} \epsilon_j K_H(\widehat{s}_i, \widehat{t}_j)}{\frac{1}{n} \sum_{j=1}^{n} K_H(\widehat{s}_i, \widehat{t}_j)} - \frac{1}{n} \sum_{i=1}^{n} \frac{\frac{1}{n} \sum_{j=1}^{n} \epsilon_j K_H(s_i, t_j)}{\frac{1}{n} \sum_{j=1}^{n} K_H(s_i, t_j)}$$

$$= \frac{1}{n} \sum_{i=1}^{n} \frac{\frac{1}{n} \sum_{j=1}^{n} \epsilon_j [K_H(\widehat{s}_i, \widehat{t}_j) - K_H(s_i, t_j)]}{\frac{1}{n} \sum_{j=1}^{n} K_H(\widehat{s}_i, \widehat{t}_j)} + \frac{1}{n} \sum_{i=1}^{n} \frac{\frac{1}{n} \sum_{j=1}^{n} \epsilon_j K_H(s_i, t_j)}{\frac{1}{n} \sum_{j=1}^{n} K_H(s_i, t_j)} \left( \frac{\frac{1}{n} \sum_{j=1}^{n} K_H(s_i, t_j)}{\frac{1}{n} \sum_{j=1}^{n} K_H(\widehat{s}_i, \widehat{t}_j)} - 1 \right) \tag{60}$$

The following two lemmas are needed to control $T_2$. The proofs of Lemma A.2 and A.3 are presented in Section B.2 and B.3, respectively.

**Lemma A.2.** *Suppose the assumptions of Theorem 5.1 hold, then with probability larger than $1 - n^{-C}$ for some positive constant $C > 1$, for all $1 \leq i \leq n$,*

$$\frac{1}{2} f_t(s_i) - C\sqrt{f_t(s_i)\frac{\log n}{nh^3}} \leq \frac{1}{n}\sum_{j=1}^{n} K_H(s_i, t_j) \leq f_t(s_i) + C\sqrt{f_t(s_i)\frac{\log n}{nh^3}} \tag{61}$$

$$\frac{1}{n}\sum_{j=1}^{n} \left| K_H(s_i, t_j) - K_H(\widehat{s}_i, \widehat{t}_j) \right| \lesssim \frac{\log n}{\sqrt{nh}} \tag{62}$$

$$\left| \frac{1}{n}\sum_{j=1}^{n} \epsilon_j K_H(s_i, t_j) \right| \lesssim \sqrt{\frac{\log n}{nh^3}} \tag{63}$$

$$\left| \frac{1}{n}\sum_{j=1}^{n} \epsilon_j [K_H(\widehat{s}_i, \widehat{t}_j) - K_H(s_i, t_j)] \right| \lesssim \frac{\log n}{n^{3/4}h^2}. \tag{64}$$

**Lemma A.3.** *Suppose the assumptions of Theorem 5.1 hold, then*

$$\frac{\frac{1}{n}\sum_{j=1}^{n} \epsilon_j a_j}{\sqrt{\frac{1}{n^2}\sum_{j=1}^{n} \text{Var}(\epsilon_j \mid d_j, w_j)a_j^2}} \to N(0, 1) \tag{65}$$

*where $\epsilon_j$ is defined in (55) and $a_j$ is defined in (59). With probability larger than $1 - n^{-C}$,*

$$\sqrt{\frac{1}{n^2}\sum_{j=1}^{n} \text{Var}(\epsilon_j \mid d_j, w_j)a_j^2} \asymp \frac{1}{\sqrt{nh^2}} \tag{66}$$

A combination of (61) and (62) leads to

$$\frac{1}{8} f_t(s_i) - C\sqrt{|f_t(s_i)|\frac{\log n}{nh^3}} \leq \frac{1}{n}\sum_{j=1}^{n} K_H(\widehat{s}_i, \widehat{t}_j) \leq f_t(s_i) + C\sqrt{|f_t(s_i)|\frac{\log n}{nh^3}}. \tag{67}$$

Together with (61), (62), (63) and $\min_i f_t(s_i) \geq c_0$ for some positive constant $c_0 > 0$,

$$\mathbf{P}\left( \left| \frac{1}{n}\sum_{i=1}^{n} \frac{\frac{1}{n}\sum_{j=1}^{n} \epsilon_j K_H(s_i, t_j)}{\frac{1}{n}\sum_{j=1}^{n} K_H(s_i, t_j)} \left( \frac{\frac{1}{n}\sum_{j=1}^{n} K_H(s_i, t_j)}{\frac{1}{n}\sum_{j=1}^{n} K_H(\widehat{s}_i, \widehat{t}_j)} - 1 \right) \right| \gtrsim \frac{(\log n)^{3/2}}{nh^{5/2}} \right) \leq n^{-C}$$

By (67), (64) and $\min_i f_t(s_i) \geq c_0$, we have

$$\mathbf{P}\left( \left| \frac{1}{n} \sum_{i=1}^{n} \frac{\frac{1}{n} \sum_{j=1}^{n} \epsilon_j [K_H(\widehat{s}_i, \widehat{t}_j) - K_H(s_i, t_j)]}{\frac{1}{n} \sum_{j=1}^{n} K_H(\widehat{s}_i, \widehat{t}_j)} \right| \gtrsim \frac{\log n}{n^{3/4} h^2} \right) \leq n^{-C}.$$

Since $nh^4 \gg (\log n)^2$, we have

$$\sqrt{nh^2} \left| \frac{1}{n} \sum_{i=1}^{n} \frac{\frac{1}{n} \sum_{j=1}^{n} \epsilon_j K_H(\widehat{s}_i, \widehat{t}_j)}{\frac{1}{n} \sum_{j=1}^{n} K_H(\widehat{s}_i, \widehat{t}_j)} - \frac{1}{n} \sum_{i=1}^{n} \frac{\frac{1}{n} \sum_{j=1}^{n} \epsilon_j K_H(s_i, t_j)}{\frac{1}{n} \sum_{j=1}^{n} K_H(s_i, t_j)} \right| = o_p(1).$$

Together with Lemma A.3, we establish that

$$\frac{\frac{1}{n} \sum_{i=1}^{n} \frac{\frac{1}{n} \sum_{j=1}^{n} \epsilon_j K_H(\widehat{s}_i, \widehat{t}_j)}{\frac{1}{n} \sum_{j=1}^{n} K_H(\widehat{s}_i, \widehat{t}_j)}}{\sqrt{\frac{1}{n^2} \sum_{j=1}^{n} \mathrm{Var}(\epsilon_j) a_j^2}} \to N(0, 1). \tag{68}$$

**Control of $T_3$.** We decompose $T_3$ as

$$\frac{1}{n} \sum_{i=1}^{n} \frac{\sum_{j=1}^{n} [\nabla g(\widehat{s}_i)]^\intercal (\widehat{t}_j - \widehat{s}_i) K_H(\widehat{s}_i, \widehat{t}_j)}{\sum_{j=1}^{n} K_H(\widehat{s}_i, \widehat{t}_j)} + \frac{1}{n} \sum_{i=1}^{n} \frac{\sum_{j=1}^{n} (\widehat{t}_j - \widehat{s}_i)^\intercal \triangle g(\widehat{s}_i + c_{ij}(\widehat{t}_j - \widehat{s}_i))(\widehat{t}_j - \widehat{s}_i) K_H(\widehat{s}_i, \widehat{t}_j)}{\sum_{j=1}^{n} K_H(\widehat{s}_i, \widehat{t}_j)} \tag{69}$$

for some constant $c_{ij} \in (0, 1)$. We show that the second term of (69) is the higher order term, controlled as,

$$\left| \frac{1}{n} \sum_{i=1}^{n} \frac{\sum_{j=1}^{n} (\widehat{t}_j - \widehat{s}_i)^\intercal \triangle g(\widehat{s}_i + c(\widehat{t}_j - \widehat{s}_i))(\widehat{t}_j - \widehat{s}_i) K_H(\widehat{s}_i, \widehat{t}_j)}{\sum_{j=1}^{n} K_H(\widehat{s}_i, \widehat{t}_j)} \right| \leq h^2$$

To establish the above inequality, we apply the boundedness assumption on the hessian $\triangle g$ imposed in Condition 5.3 (b) and and we use the fact that, if $K_H(\widehat{s}_i, \widehat{t}_j) > 0$ and $C \log n / \sqrt{n} \leq h/2$, then $\|\widehat{t}_j - s_i\|_\infty \leq \|\widehat{t}_j - \widehat{s}_i\|_\infty + \|\widehat{s}_i - s_i\|_\infty \leq h$.

Now we control the first term of (69) as

$$
\begin{aligned}
\frac{1}{n} \sum_{i=1}^{n} & \frac{\sum_{j=1}^{n} [\nabla g(\widehat{s}_i)]^{\mathsf{T}} (\widehat{t}_j - \widehat{s}_i) K_H(\widehat{s}_i, \widehat{t}_j)}{\sum_{j=1}^{n} K_H(\widehat{s}_i, \widehat{t}_j)} \\
= {}& \frac{1}{n} \sum_{i=1}^{n} \frac{\sum_{j \neq i} [\nabla g(\widehat{s}_i)]^{\mathsf{T}} (\widehat{t}_j - \widehat{s}_i) K_H(\widehat{s}_i, \widehat{t}_j)}{\sum_{j=1}^{n} K_H(\widehat{s}_i, \widehat{t}_j)} + \frac{1}{n} \sum_{i=1}^{n} \frac{[\nabla g(\widehat{s}_i)]^{\mathsf{T}} (\widehat{t}_i - \widehat{s}_i) K_H(\widehat{s}_i, \widehat{t}_i)}{\sum_{j=1}^{n} K_H(\widehat{s}_i, \widehat{t}_j)} \\
= {}& \frac{1}{n} \sum_{i=1}^{n} \frac{\sum_{j \neq i} [\nabla g(s_i)]^{\mathsf{T}} (t_j - s_i) K_H(s_i, t_j)}{\sum_{j=1}^{n} K_H(s_i, t_j)} + \frac{1}{n} \sum_{i=1}^{n} \frac{[\nabla g(\widehat{s}_i)]^{\mathsf{T}} (\widehat{t}_i - \widehat{s}_i) K_H(\widehat{s}_i, \widehat{t}_i)}{\sum_{j=1}^{n} K_H(\widehat{s}_i, \widehat{t}_j)} \\
& + \frac{1}{n} \sum_{i=1}^{n} \frac{\sum_{j \neq i} [\nabla g(\widehat{s}_i)]^{\mathsf{T}} (\widehat{t}_j - \widehat{s}_i) K_H(\widehat{s}_i, \widehat{t}_j)}{\sum_{j=1}^{n} K_H(\widehat{s}_i, \widehat{t}_j)} - \frac{1}{n} \sum_{i=1}^{n} \frac{\sum_{j \neq i} [\nabla g(s_i)]^{\mathsf{T}} (t_j - s_i) K_H(s_i, t_j)}{\sum_{j=1}^{n} K_H(s_i, t_j)}
\end{aligned}
\tag{70}
$$

We introduce the following lemma to control (70), whose proof can be found in Section B.4.

**Lemma A.4.** *Suppose the assumptions of Theorem 5.1 hold, then with probability larger than* $1 - n^{-C}$ *for some positive constant* $C > 0$,

$$
\left| \frac{1}{n} \sum_{i=1}^{n} \frac{\sum_{j \neq i} [\nabla g(s_i)]^{\mathsf{T}} (t_j - s_i) K_H(s_i, t_j)}{\sum_{j=1}^{n} K_H(s_i, t_j)} \right| \lesssim h^2 + \sqrt{\frac{\log n}{nh}}
\tag{71}
$$

*and*

$$
\left| \frac{1}{n} \sum_{i=1}^{n} \frac{[\nabla g(\widehat{s}_i)]^{\mathsf{T}} (\widehat{t}_i - \widehat{s}_i) K_H(\widehat{s}_i, \widehat{t}_i)}{\sum_{j=1}^{n} K_H(\widehat{s}_i, \widehat{t}_j)} \right| \lesssim \frac{1}{nh^2}
\tag{72}
$$

*and*

$$
\left| \frac{1}{n} \sum_{i=1}^{n} \frac{\sum_{j \neq i} [\nabla g(\widehat{s}_i)]^{\mathsf{T}} (\widehat{t}_j - \widehat{s}_i) K_H(\widehat{s}_i, \widehat{t}_j)}{\sum_{j=1}^{n} K_H(\widehat{s}_i, \widehat{t}_j)} - \frac{1}{n} \sum_{i=1}^{n} \frac{\sum_{j \neq i} [\nabla g(s_i)]^{\mathsf{T}} (t_j - s_i) K_H(s_i, t_j)}{\sum_{j=1}^{n} K_H(s_i, t_j)} \right| \lesssim \frac{\log n}{\sqrt{n}}
\tag{73}
$$

By applying Lemma A.4, we have

$$
\left| \frac{1}{n} \sum_{i=1}^{n} \frac{\sum_{j=1}^{n} [g(\widehat{t}_j) - g(\widehat{s}_i)] K_H(\widehat{s}_i, \widehat{t}_j)}{\sum_{j=1}^{n} K_H(\widehat{s}_i, \widehat{t}_j)} \right| \lesssim h^2 + \frac{1}{nh^2} + \sqrt{\frac{\log n}{nh}}
\tag{74}
$$

53

By combining (58), (68) and (74), we establish that, with probability larger than $1 - \frac{1}{t^2} - n^{-C} - P(E_1)$ for some positive constant $C > 0$,

$$\left| \widehat{\mathrm{ASF}}(d, w) - \mathrm{ASF}(d, w) \right| \lesssim \frac{t}{\sqrt{nh^2}} + h^2 + \frac{\log n}{\sqrt{n}} + \sqrt{\frac{\log n}{nh}}$$

This implies (30) in the main paper under the bandwidth condition $h = n^{-\mu}$ for $0 < \mu < 1/4$. Together with (66), (68) and the bandwidth condition that $h = n^{-\mu}$ for $0 < \mu < 1/6$, we establish the asymptotic normality and the asymptotic variance level in Theorem 5.1.

### A.6  Proof of Corollary 5.1

The proof is similar to that of Theorem 5.1. The main extra step is to establish the asymptotic variance (31) in the main paper. We introduce the following lemma as a modification of Lemma A.3 and present its proof in Section B.3.

**Lemma A.5.** *Suppose the assumptions of Corollary 5.1 hold, then*

$$\frac{\frac{1}{n} \sum_{j=1}^{n} \epsilon_j c_j}{\sqrt{\frac{1}{n^2} \sum_{j=1}^{n} \mathrm{Var}(\epsilon_j \mid d_j, w_j) c_j^2}} \to N(0, 1) \tag{75}$$

*where $\epsilon_j$ is defined in (55) and*

$$c_j = \frac{1}{n} \sum_{i=1}^{n} \frac{K_H(s_i, t_j)}{\frac{1}{n} \sum_{j=1}^{n} K_H(s_i, t_j)} - \frac{1}{n} \sum_{i=1}^{n} \frac{K_H(r_i, t_j)}{\frac{1}{n} \sum_{j=1}^{n} K_H(r_i, t_j)}.$$

*With probability larger than $1 - n^{-C}$,*

$$\sqrt{\frac{\mathrm{V_{CATE}}}{n}} = \sqrt{\frac{1}{n^2} \sum_{j=1}^{n} \mathrm{Var}(\epsilon_j \mid d_j, w_j) c_j^2} \asymp \frac{1}{\sqrt{nh^2}} \tag{76}$$

Then we apply the above lemma together with the same arguments as Theorem 5.1 to establish Corollary 5.1.

# B Proof of Lemmas

## B.1 Proof of Lemma A.1

The error $\frac{1}{n}\sum_{i=1}^n g(\widehat{s}_i) - \int g(s_i)f_v(v_i)dv_i$ can be decomposed as

$$\frac{1}{n}\sum_{i=1}^n g(s_i) - \int g(s_i)f_v(v_i)dv_i + \frac{1}{n}\sum_{i=1}^n [g(\widehat{s}_i) - g(s_i)]$$

Since $\frac{1}{n}\sum_{i=1}^n g(s_i) - \int g(s_i)f_v(v_i)dv_i$ has mean zero and variance

$$\frac{1}{n}\int \left( g(s_i) - \int g(s_i)f_v(v_i)dv_i \right)^2 f_v(v_i)dv_i \leq \frac{1}{n}\int g^2(s_i)f_v(v_i)dv_i,$$

we establish that, with probability larger than $1 - \frac{1}{t^2}$, $\left|\frac{1}{n}\sum_{i=1}^n g(s_i) - \int g(s_i)f_v(v_i)dv_i\right| \lesssim t/\sqrt{n}$. Together with the fact that $|g(\widehat{s}_i) - g(s_i)| \leq \|\nabla g(s_i + c(\widehat{s}_i - s_i)\|_2\|\widehat{s}_i - s_i\|_2 \leq \max_s \|\nabla g(s)\|_2\|\widehat{s}_i - s_i\|_2 \lesssim \log n/\sqrt{n}$, we establish (58).

## B.2 Proof of Lemma A.2

We use $\mathcal{T} \in \mathbb{R}^3$ to denote the support of $t_j$ and assume that $\min_{1\leq i\leq n} f_t(s_i) \geq c_0$ for a given positive constant $c_0 > 0$. For $j \neq i$, $s_i$ is independent of $t_j$ and we use $\mathbb{E}_j K_H(s_i, t_j)$ to denote the expectation taken with respect to $t_j$ conditioning on $s_i$.

We now show that $\mathbb{E}_j K_H(s_i, t_j)$ for $j \neq i$ is close to $f_t(s_i)$ by expressing $\mathbb{E}_j K_H(s_i, t_j)$ as

$$
\begin{aligned}
&\mathbb{E}_j K_H(s_i, t_j) \\
&= \int_{\|t-s_i\|_\infty \leq h/2} \frac{1}{h^3} f_t(t)\mathbf{1}_{\{t\in\mathcal{T}\}}dt \\
&= \int_{\|t-s_i\|_\infty \leq h/2} \frac{1}{h^3} \left( f_t(s_i) + [\nabla f_t(s_i + c(t - s_i))]^\mathsf{T}(t - s_i) \right) \mathbf{1}_{\{t\in\mathcal{T}\}}dt \\
&= f_t(s_i)c_* + \int_{\|t-s_i\|_\infty \leq h/2} \frac{1}{h^3}[\nabla f_t(s_i + c(t - s_i))]^\mathsf{T}(t - s_i)\mathbf{1}_{\{t\in\mathcal{T}\}}dt
\end{aligned}
$$

where $0 < c < 1$ is a positive constant and $c_* = \int_{\|t-s_i\|_\infty \leq h/2} \frac{1}{h^3}\mathbf{1}_{\{t\in\mathcal{T}\}}dt$. Note that

$$\int_{\|t-s_i\|_\infty \leq h/2} \frac{1}{h^3}\mathbf{1}_{\{t\in\mathcal{T}\}}dt = 1 - \int \frac{1}{h^3}\mathbf{1}_{\{t\notin\mathcal{T},\|t-s_i\|_\infty\leq h/2\}}dt.$$

Under the Condition 5.3 (c), the event $\{t \notin \mathcal{T}, \|t - s_i\|_\infty \leq h/2\}$ implies that the third entry $v$ of the vector $t \in \mathbb{R}^3$ does not belong to the support $\mathcal{T}_v$ of $f_v$. Hence

$$\int \frac{1}{h^3} \mathbf{1}_{\{t \notin \mathcal{T}, \|t - s_i\|_\infty \leq h/2\}} dt \leq \int \frac{1}{h^3} \mathbf{1}_{\{v \notin \mathcal{T}_v, \|t - s_i\|_\infty \leq h/2\}} dt = \frac{1}{h} \int \mathbf{1}_{\{v \notin \mathcal{T}_v, \|v - v_i\|_\infty \leq h/2\}} dv.$$

Define $v_{\min} = \inf_v \{v : f_v > 0\}$ and $v_{\max} = \sup_v \{v : f_v > 0\}$. We adopt the notation that $v_{\min} = -\infty$ and $v_{\max} = \infty$ when the support $\mathcal{T}_v$ is unbounded from below and above, respectively. We have

$$\frac{1}{h} \int \mathbf{1}_{\{v \notin \mathcal{T}_v, \|v - v_i\|_\infty \leq h/2\}} dv \leq \frac{1}{h} \max \left\{ \int_{v_{\min} - h/2}^{v_{\min}} dv, \int_{v_{\max}}^{v_{\max} + h/2} dv \right\} = 1/2.$$

Hence we have $c_* \in [1/2, 1]$. Since $\|\nabla f_t(s_i + c(t - s_i))\|_2 \leq C$, we establish that, for $j \neq i$,

$$|\mathbb{E}_j K_H(s_i, t_j) - c_* f_t(s_i)| \leq Ch \quad \text{for} \quad c_* \in [1/2, 1]. \tag{77}$$

**Proof of** (61). We state the Bernstein inequality (Bennett 1962) in the following lemma.

**Lemma B.1.** *Suppose that $\{X_i\}_{1 \leq i \leq n}$ are independent zero mean random variables and $|X_i| \leq M$ almost surely. Then we have*

$$\mathbf{P} \left( \left| \sum_{i=1}^n X_i \right| \geq t \right) \leq 2 \exp \left( -\frac{t^2/2}{\sum_{i=1}^n \mathbb{E} X_i^2 + Mt/3} \right).$$

We decompose

$$\frac{1}{n} \sum_{j=1}^n K_H(s_i, t_j) = \left(1 - \frac{1}{n}\right) \frac{1}{n-1} \sum_{j \neq i} K_H(s_i, t_j) + \frac{1}{n} K_H(s_i, t_j).$$

We fix $1 \leq i \leq n$ and take $j \neq i$. Since $|f(s_i)| \geq c_0$ for some positive constant $c_0 > 0$ and $\mathbb{E} K_H(s_i, t_j)^2 = \mathbb{E} K_H(s_i, t_j)/h^3$, it follows from (77) that $\sum_{j \neq i} \mathbb{E} K_H(s_i, t_j)^2 \lesssim |f_t(s_i)| n/h^3$. We

now apply Lemma B.1 with $M = 1/h^3$ and obtain

$$\mathbf{P}\left(\left|\frac{1}{n-1}\sum_{j\neq i}(K_H(s_i, t_j) - \mathbb{E}_j K_H(s_i, t_j))\right| \gtrsim \sqrt{|f_t(s_i)|\frac{\log n}{nh^3}}\right) \leq n^{-C} \tag{78}$$

for some large positive constant $C > 1$. Together with $\frac{1}{n}K_H(s_i, t_j) \leq \frac{1}{nh^3}$ and $|f(s_i)| \geq c_0$ for some positive constant $c_0 > 0$, we establish (61).

**Proof of** (62). Define $h_a = h - 2C_0 \log n/\sqrt{n}$ and $h_b = h + 2C_0 \log n/\sqrt{n}$ for some large constant $C_0 > 0$. Define the set $B_a = \{t \in \mathbb{R}^3 : \|t - s_i\|_\infty \leq h_a\}$ $B_b = \{t \in \mathbb{R}^3 : \|t - s_i\|_\infty \leq h_b\}$ and define the kernel functions

$$K_H^a(s_i, t_j) = \frac{1}{h^3}\prod_{l=1}^{3} k\left(\frac{s_{il} - t_{jl}}{h_a}\right) \quad \text{and} \quad K_H^b(s_i, t_j) = \frac{1}{h^3}\prod_{l=1}^{3} k\left(\frac{s_{il} - t_{jl}}{h_b}\right) \tag{79}$$

where $k(x) = \mathbf{1}(|x| \leq 1/2)$. On the event $\mathcal{A}_1 \cap \mathcal{A}_2$, we have $\max_{1\leq l\leq 3}\left|[\widehat{s}_i - \widehat{t}_j]_l - (s_i - t_j)_l\right| \leq 2C_0 \log n/\sqrt{n}$, and hence

$$K_H[\widehat{s}_i, \widehat{t}_j] \leq K_H^b(s_i, t_j) \quad \text{and} \quad K_H[\widehat{s}_i, \widehat{t}_j] \geq K_H^a(s_i, t_j).$$

Then we establish that, on the event $\mathcal{A}_1 \cap \mathcal{A}_2$,

$$\left|\frac{1}{n}\sum_{j\neq i}K_H(\widehat{s}_i, \widehat{t}_j) - \frac{1}{n}\sum_{j\neq i}K_H(s_i, t_j)\right| \leq \frac{1}{n}\sum_{j\neq i}\left|K_H(\widehat{s}_i, \widehat{t}_j) - K_H(s_i, t_j)\right| \leq \frac{1}{n}\sum_{j\neq i}(K_H^b - K_H^a)(s_i, t_j) \tag{80}$$

Conditioning on the $i$-th observation, we have

$$\mathbb{E}_j\left[(K_H^b - K_H^a)(s_i, t_j)\right] \lesssim \frac{1}{h^3}\mathbb{E}\left(\mathbf{1}_{t_j \in B_b} - \mathbf{1}_{t_j \in B_a}\right) \lesssim \frac{1}{h^3}(h_b^3 - h_a^3) \lesssim \frac{\log n}{\sqrt{n}h} \tag{81}$$

where the last inequality follows from the fact that $h_b^3 - h_a^3 \lesssim h^2 \log n/\sqrt{n}$. Since

$$\left|(K_H^b - K_H^a)(s_i, t_j) - \mathbb{E}_j(K_H^b - K_H^a)(s_i, t_j)\right| \leq \frac{1}{h^3}$$

and

$$\sum_{j \neq i} \mathbb{E}_j \left( \left[ (K_H^b - K_H^a)(s_i, t_j) - \mathbb{E}_j (K_H^b - K_H^a)(s_i, t_j) \right] \right)^2$$

$$\leq n \mathbb{E}_j \left[ (K_H^b - K_H^a)(s_i, t_j) \right]^2 \leq \frac{n}{h^3} \mathbb{E}_j \left[ (K_H^b - K_H^a)(s_i, t_j) \right] \lesssim \frac{\sqrt{n} \log n}{h^4},$$

we apply Lemma B.1 and establish

$$\mathbf{P} \left( \left| \frac{1}{n-1} \sum_{j \neq i} \left[ (K_H^b - K_H^a)(s_i, t_j) - \mathbb{E}_j (K_H^b - K_H^a)(s_i, t_j) \right] \right| \gtrsim \frac{\log n}{nh^3} \cdot (\sqrt{n} h^2 \log n)^{1/2} \right) \leq n^{-C}$$

Together with (81), we establish (62).

**Proof of** (63). Note that, conditioning on the $i$-th data point, $\{\epsilon_j K_H(s_i, t_j)\}_{j \neq i}$ are independent mean zero random variable with $|\epsilon_j K_H(s_i, t_j)| \leq 1/h^3$ and

$$\sum_{j=1}^n \mathbb{E}_j \left( \epsilon_j K_H(s_i, t_j) \right)^2 \lesssim |f_t(s_i)| \, n/h^3,$$

where the inequality follows from (77) and boundedness of $\epsilon_j$. We apply Lemma B.1 and establish

$$\mathbf{P} \left( \left| \frac{1}{n-1} \sum_{j \neq i} \epsilon_j K_H(s_i, t_j) \right| \gtrsim \sqrt{|f_t(s_i)| \frac{\log n}{nh^3}} \right) \leq n^{-C}$$

Together with $\frac{1}{n} \epsilon_j K_H(s_i, t_j) \leq \frac{1}{nh^3}$, we establish (63).

**Proof of** (64). For $\widehat{s}_i$ and $\widehat{t}_i$ where $1 \leq i \leq n$, we express them in terms of the difference matrix $\widehat{\Delta}^B = \widehat{B} - B \in \mathbb{R}^{(p+1) \times 2}$ and the difference vector $\widehat{\Delta}^\gamma = \widehat{\gamma} - \gamma \in \mathbb{R}^p$,

$$\widehat{s}_i - s_i = \left( (d, w^\intercal) \widehat{\Delta}^B, w_i^\intercal \widehat{\Delta}^\gamma \right)^\intercal, \quad \widehat{t}_i - t_i = \left( (d_i, w_i^\intercal) \widehat{\Delta}^B, w_i^\intercal \widehat{\Delta}^\gamma \right)^\intercal.$$

Define the general difference matrix $\Delta^B \in \mathbb{R}^{(p+1) \times 2}$, the difference vector $\Delta^\gamma \in \mathbb{R}^p$ and $\Delta = ((\Delta_{\cdot 1}^B)^\intercal, (\Delta_{\cdot 2}^B)^\intercal, (\Delta^\gamma)^\intercal)^\intercal \in \mathbb{R}^{3p+2}$. We introduce general functions $s_i : \mathbb{R}^{3p+2} \to \mathbb{R}$ and $t_i : \mathbb{R}^{3p+2} \to \mathbb{R}$,

$$s_i(\Delta) = s_i + \left( (d, w^\intercal) \Delta^B, w_i^\intercal \Delta^\gamma \right)^\intercal, \quad t_j(\Delta) = t_j + \left( (d_i, w_i^\intercal) \Delta^B, w_i^\intercal \Delta^\gamma \right)^\intercal$$

and have $\widehat{s}_i = s_i(\widehat{\Delta})$ and $\widehat{t}_i = t_i(\widehat{\Delta})$ and $s_i = s_i(0)$ and $t_j = t_j(0)$. On the event $\mathcal{A}_1 \cap \mathcal{A}_2$, we have $\widehat{\Delta} = ((\widehat{\Delta}_{\cdot 1}^B)^\intercal, (\widehat{\Delta}_{\cdot 2}^B)^\intercal, (\widehat{\Delta}^\gamma)^\intercal)^\intercal \in B^{3p+2}\left(C\sqrt{\log n/n}\right)$ where $B^{3p+2}\left(C\sqrt{\log n/n}\right)$ denotes the ball in $\mathbb{R}^{3p+2}$ with radius $C\sqrt{\log n/n}$ for a large constant $C > 0$. We use $\{\Delta_1, \cdots, \Delta_{L_n}\}$ to denote a $\tau_n$-net of $B^{3p+2}\left(C\sqrt{\log n/n}\right)$ such that for any $\Delta \in B^{3p+2}\left(C\sqrt{\log n/n}\right)$, there exists $1 \le l \le L_n$ such that $\|\Delta - \Delta_l\|_2 \le \tau_n$, where $\tau_n > 0$ is a positive number, $L_n$ is a positive integer and both $\tau_n$ and $L_n$ are allowed to grow with the sample size $n$. It follows from Lemma 5.2 of Vershynin (2010) that

$$L_n \le \left(1 + \frac{2C\sqrt{\log n/n}}{\tau_n}\right)^{3p+2}. \tag{82}$$

For $\widehat{\Delta}$, there exists $1 \le \widehat{l} \le L_n$ such that $\|\widehat{\Delta} - \Delta_{\widehat{l}}\|_2 \le \tau_n$ and hence

$$\left| \frac{1}{n} \sum_{j=1}^n \epsilon_j [K_H(\widehat{s}_i, \widehat{t}_j) - K_H(s_i, t_j)] \right|$$

$$\le \left| \frac{1}{n} \sum_{j=1}^n \epsilon_j [K_H(s_i(\Delta_{\widehat{l}}), t_j(\Delta_{\widehat{l}})) - K_H(s_i, t_j)] \right| + \left| \frac{1}{n} \sum_{j=1}^n \epsilon_j [K_H(\widehat{s}_i, \widehat{t}_j) - K_H(s_i(\Delta_{\widehat{l}}), t_j(\Delta_{\widehat{l}}))] \right| \tag{83}$$

We shall control the the two terms on the right hand side of (83). Regarding the first term on the right hand side of (83), we have

$$\left| \frac{1}{n} \sum_{j=1}^n \epsilon_j [K_H(s_i(\Delta_{\widehat{l}}), t_j(\Delta_{\widehat{l}})) - K_H(s_i, t_j)] \right| \le \frac{1}{nh^3} + \max_{1 \le l \le L_n} \left| \frac{1}{n} \sum_{j \ne i} \epsilon_j [K_H(s_i(\Delta_l), t_j(\Delta_l)) - K_H(s_i, t_j)] \right| \tag{84}$$

In the following, we control (84) by the maximal inequality and a similar argument as the proof of (63). Note that, conditioning on the $i$-th data point, we have $\{\epsilon_j [K_H(s_i(\Delta_l), t_j(\Delta_l)) - K_H(s_i, t_j)]\}_{j \ne i}$ are independent mean zero random variable with

$$|\epsilon_j [K_H(s_i(\Delta_l), t_j(\Delta_l)) - K_H(s_i, t_j)]| \le 1/h^3$$

and

$$\sum_{j=1}^n \mathbb{E}_j \left(\epsilon_j [K_H(s_i(\Delta_l), t_j(\Delta_l)) - K_H(s_i, t_j)]\right)^2 \lesssim \frac{n}{h^3} \mathbb{E}_j \left(K_H^b(s_i, t_j) - K_H^a(s_i, t_j)\right) \le \frac{\sqrt{n}\log n}{h^4}$$

59

where the last inequality follows from (81). We apply Lemma B.1 and establish

$$\mathbf{P}\left(\left|\frac{1}{n}\sum_{j\neq i}\epsilon_j[K_H(s_i(\Delta_l),t_j(\Delta_l))-K_H(s_i,t_j)]\right|\gtrsim\frac{1}{n}\sqrt{\log(n\cdot L_n)\frac{\sqrt{n}\log n}{h^4}}\right)\leq(nL_n)^{-C}$$

for some positive constant $C>1$. By the maximal inequality, we have

$$\mathbf{P}\left(\max_{1\leq l\leq L_n}\left|\frac{1}{n}\sum_{j\neq i}\epsilon_j[K_H(s_i(\Delta_l),t_j(\Delta_l))-K_H(s_i,t_j)]\right|\gtrsim\frac{1}{n}\sqrt{\log(n\cdot L_n)\frac{\sqrt{n}\log n}{h^4}}\right)\leq L_n\cdot(nL_n)^{-C}$$

Together with (84) and $nh^4(\log n)^2\to\infty$, we establish

$$\mathbf{P}\left(\left|\frac{1}{n}\sum_{j=1}^{n}\epsilon_j[K_H(s_i(\Delta_{\widehat{l}}),t_j(\Delta_{\widehat{l}}))-K_H(s_i,t_j)]\right|\gtrsim\frac{1}{n}\sqrt{\log(n\cdot L_n)\frac{\sqrt{n}\log n}{h^4}}\right)\leq L_n\cdot(nL_n)^{-C}$$

$$(85)$$

Regarding the second term on the right hand side of (83), it follows from boundedness of $\epsilon_i$ that

$$\left|\frac{1}{n}\sum_{j=1}^{n}\epsilon_j[K_H(\widehat{s}_i,\widehat{t}_j)-K_H(s_i(\Delta_{\widehat{l}}),t_j(\Delta_{\widehat{l}}))]\right|\lesssim\frac{1}{n}\sum_{j=1}^{n}\left|K_H(\widehat{s}_i,\widehat{t}_j)-K_H(s_i(\Delta_{\widehat{l}}),t_j(\Delta_{\widehat{l}}))\right|$$

$$\leq\frac{1}{nh^3}+\frac{1}{n}\sum_{j\neq i}\left|K_H(\widehat{s}_i,\widehat{t}_j)-K_H(s_i(\Delta_{\widehat{l}}),t_j(\Delta_{\widehat{l}}))\right|$$

$$(86)$$

Define $h_c=h-C\sqrt{\log n}\tau_n$ and $h_d=h+C\sqrt{\log n}\tau_n$ for some large positive constant $C>0$ and define the kernel functions

$$K_H^c(s_i,t_j)=\frac{1}{h^3}\prod_{l=1}^{3}k\left(\frac{s_{il}-t_{jl}}{h_c}\right)\quad\text{and}\quad K_H^d(s_i,t_j)=\frac{1}{h^3}\prod_{l=1}^{3}k\left(\frac{s_{il}-t_{jl}}{h_d}\right)$$

On the event $\mathcal{A}_2$, we have $\|\widehat{s}_i-s_i(\Delta_{\widehat{l}})\|_2\leq C\sqrt{\log n}\tau_n$ and $\|\widehat{t}_j-t_j(\Delta_{\widehat{l}})\|_2\leq C\sqrt{\log n}\tau_n$. As a consequence, we have $K_H[\widehat{s}_i,\widehat{t}_j]\leq K_H^d(s_i(\Delta_{\widehat{l}}),t_j(\Delta_{\widehat{l}}))$ and $K_H[\widehat{s}_i,\widehat{t}_j]\geq K_H^c(s_i(\Delta_{\widehat{l}}),t_j(\Delta_{\widehat{l}}))$ and then obtain

$$\frac{1}{n}\sum_{j\neq i}\left|K_H(\widehat{s}_i,\widehat{t}_j)-K_H(s_i(\Delta_{\widehat{l}}),t_j(\Delta_{\widehat{l}}))\right|\leq\frac{1}{n}\sum_{j\neq i}(K_H^d-K_H^c)(s_i(\Delta_{\widehat{l}}),t_j(\Delta_{\widehat{l}}))$$

Together with (86), we establish

$$\left| \frac{1}{n} \sum_{j=1}^{n} \epsilon_j [K_H(\widehat{s}_i, \widehat{t}_j) - K_H(s_i(\Delta_{\widehat{l}}), t_j(\Delta_{\widehat{l}}))] \right| \leq \frac{1}{nh^3} + \max_{1 \leq l \leq L_n} \left| \frac{1}{n} \sum_{j \neq i} (K_H^d - K_H^c)(s_i(\Delta_l), t_j(\Delta_l)) \right|$$

(87)

Now we control (87) using the similar argument as that for (62). Similar to (81), we have

$$\mathbb{E}_j \left[ (K_H^d - K_H^c)(s_i(\Delta_l), t_j(\Delta_l) \right] \lesssim \frac{1}{h^3}(h_d^3 - h_c^3) \lesssim \tau_n/h(1 + \tau_n/h) \lesssim \tau_n/h$$

where the last inequality follows for $\tau_n \lesssim h$.

Since

$$\left| (K_H^d - K_H^c)(s_i(\Delta_l), t_j(\Delta_l) - \mathbb{E}_j(K_H^d - K_H^c)(s_i(\Delta_l), t_j(\Delta_l) \right| \leq \frac{1}{h^3}$$

and

$$\sum_{j \neq i} \mathbb{E}_j \left[ (K_H^d - K_H^c)(s_i(\Delta_l), t_j(\Delta_l) - \mathbb{E}_j(K_H^d - K_H^c)(s_i(\Delta_l), t_j(\Delta_l) \right]^2$$

$$\leq n \mathbb{E}_j \left[ (K_H^d - K_H^c)(s_i(\Delta_l), t_j(\Delta_l) \right]^2 \leq \frac{n}{h^3} \mathbb{E}_j \left[ (K_H^d - K_H^c)(s_i(\Delta_l), t_j(\Delta_l) \right] \lesssim \frac{n\tau_n}{h^4}(1 + \tau_n/h).$$

we apply Lemma B.1 and establish that, with probability larger than $1 - (nL_n)^{-C}$ for some large constant $C > 1$,

$$\left| \frac{1}{n} \sum_{j \neq i} \left[ (K_H^d - K_H^c)(s_i(\Delta_l), t_j(\Delta_l) - \mathbb{E}_j(K_H^d - K_H^c)(s_i(\Delta_l), t_j(\Delta_l) \right] \right| \lesssim \frac{\log(nL_n)}{nh^3} \left( 1 + \sqrt{\frac{nh^2\tau_n}{\log(nL_n)}} \right)$$

and hence we have

$$\mathbf{P} \left( \max_{1 \leq l \leq L_n} \frac{1}{n} \sum_{j \neq i} \left[ (K_H^d - K_H^c)(s_i(\Delta_l), t_j(\Delta_l)) \right] \gtrsim \frac{\tau_n}{h} + \frac{\log(nL_n)}{nh^3} \left( 1 + \sqrt{\frac{nh^2\tau_n}{\log(nL_n)}} \right) \right) \leq (nL_n)^{-C}$$

(88)

We take $\tau_n = \frac{1}{\sqrt{n}} \cdot \sqrt{\log n/n}$ and then use (82) to establish $\log L_n \lesssim (3q + 2) \log n$ and hence

$$\mathbf{P} \left( \max_{1 \leq l \leq L_n} \frac{1}{n} \sum_{j \neq i} \left[ (K_H^d - K_H^c)(s_i(\Delta_l), t_j(\Delta_l)) \right] \gtrsim \frac{\log n}{nh^3} \right) \leq L_n \cdot (nL_n)^{-C}$$

61

Together with (87), we establish

$$\mathbf{P}\left(\left|\frac{1}{n}\sum_{j=1}^{n}\epsilon_j[K_H(\widehat{s}_i,\widehat{t}_j) - K_H(s_i(\Delta_{\widehat{l}}), t_j(\Delta_{\widehat{l}}))]\right| \gtrsim \frac{\log n}{nh^3}\right) \leq L_n \cdot (nL_n)^{-C}$$

Together with (83), (85) and $nh^4(\log n)^2 \to \infty$, we establish (64).

### B.3  Proof of Lemmas A.3 and A.5

**Proof of Lemma A.3.**   We note that, conditioning on $\{d_j, w_j\}_{1\leq j\leq n}$, $a_j\epsilon_j = a_j(y_j - g(t_j))$ are independent random variables and

$$\mathbb{E}\frac{1}{n}\sum_{j=1}^{n}\epsilon_j a_j = \mathbb{E}\frac{1}{n}\sum_{j=1}^{n}\mathbb{E}(\epsilon_j \mid d_j, w_j)a_j = 0.$$

We now check the Lyapunov condition by calculating

$$V = \sum_{j=1}^{n}\mathrm{Var}(\epsilon_j \mid d_j, w_j)a_j^2 = \sum_{j=1}^{n}g(t_j)(1 - g(t_j))a_j^2.$$

We can express the weight $a_j = \frac{1}{n}\sum_{i=1}^{n}\frac{K_H(s_i,t_j)}{\frac{1}{n}\sum_{j=1}^{n}K_H(s_i,t_j)}$ as,

$$\frac{1}{h^2}\mathbf{1}(|t_{j1} - s_{i1}| \leq h/2)\mathbf{1}(|t_{j2} - s_{i2}| \leq h/2)\frac{1}{nh}\sum_{i=1}^{n}\mathbf{1}(|v_i - v_j| \leq h/2)\frac{1}{\frac{1}{n}\sum_{j=1}^{n}K_H(s_i,t_j)} \quad (89)$$

since $s_{i1}$ and $s_{i2}$ remain the same for all $1 \leq i \leq n$. We define two events $\mathcal{A}_3$ and $\mathcal{A}_4$ as

$$\mathcal{A}_3 = \left\{c_1 \leq \frac{1}{nh}\sum_{i=1}^{n}\mathbf{1}(|v_i - v_j| \leq h/2)\frac{1}{\frac{1}{n}\sum_{j=1}^{n}K_H(s_i,t_j)} \leq C_1, \quad \text{for } 1 \leq j \leq n.\right\}$$

$$\mathcal{A}_4 = \left\{\frac{1}{h^{2\tau}} \lesssim \frac{1}{n}\sum_{j=1}^{n}\left(\frac{1}{h^2}\mathbf{1}(|t_{j1} - s_{i1}| \leq h/2)\mathbf{1}(|t_{j2} - s_{i2}| \leq h/2)\right)^{1+\tau} \lesssim \frac{1}{h^{2\tau}}\right\}$$

for any positive constant $\tau > 0$. At the end of this subsection, we show that

$$\mathbf{P}(\mathcal{A}_3 \cap \mathcal{A}_4) \geq 1 - n^{-C}, \quad \text{for some large constant } C > 0. \quad (90)$$

On the event $\mathcal{A}_3$, it follows from (89) that

$$c_1 \leq \frac{a_j}{\frac{1}{h^2}\mathbf{1}(|t_{j1} - s_{i1}| \leq h/2)\mathbf{1}(|t_{j2} - s_{i2}| \leq h/2)} \leq C_1, \quad \text{for } 1 \leq j \leq n.$$

and hence

$$V^2 \asymp \sum_{j=1}^{n} g(t_j)(1 - g(t_j)) \left( \frac{1}{h^2}\mathbf{1}(|t_{j1} - s_{i1}| \leq h/2)\mathbf{1}(|t_{j2} - s_{i2}| \leq h/2) \right)^2 \tag{91}$$

On the event $\mathcal{A}_3 \cap \mathcal{A}_4$, we have

$$\sum_{j=1}^{n} |a_j|^{2+c} \lesssim n\frac{1}{h^{2(1+c)}} \quad \text{for any positive constant} \quad c > 0. \tag{92}$$

By Condition 5.3(b), since $g(s_i)$ is bounded away from zero and one and the gradient $\nabla g$ is bounded near $s_i$, we establish that

$$g(t_j)(1 - g(t_j))\mathbf{1}(|t_{j1} - s_{i1}| \leq h/2)\mathbf{1}(|t_{j2} - s_{i2}| \leq h/2) \geq c$$

for a positive constant $c > 0$. Hence, on the event $\mathcal{A}_4$, we have

$$V^2 \asymp n/h^2 \tag{93}$$

Then for any positive constant $c > 0$, we have

$$\frac{1}{V^{1+\frac{c}{2}}} \sum_{j=1}^{n} \mathbb{E}\left[ (\epsilon_j a_j)^{(2+c)} \mid d_j, w_j \right] \cdot \mathbf{1}_{\mathcal{A}_3 \cap \mathcal{A}_4} \lesssim \frac{1}{(n/h^2)^{1+\frac{c}{2}}} \cdot n \cdot \frac{1}{h^{2(1+c)}} \leq \frac{1}{(nh^2)^{c/2}}$$

where the second inequality follows from (92) and bounded $\epsilon_j$. Hence, we have checked Lyapunov condition and shown that

$$\frac{\sum_{j=1}^{n} \epsilon_j a_j}{\sqrt{V}} \mid \{d_j, w_j\}_{1 \leq j \leq n} \in \mathcal{A}_3 \cap \mathcal{A}_4 \xrightarrow{d} N(0, 1)$$

Together with (90), we establish (65). We establish (66) by (93).

63

**Proof of Lemma A.5.** The proof of Lemma A.5 is similar to that of Lemma A.3. We define
$a'_j = \frac{1}{n} \sum_{i=1}^{n} \frac{K_H(r_i, t_j)}{\frac{1}{n} \sum_{j=1}^{n} K_H(r_i, t_j)}$ and then $c_j = a_j - a'_j$. Similar to (89), we have

$$a'_j = \frac{1}{h^2} \mathbf{1}(|t_{j1} - r_{i1}| \leq h/2) \mathbf{1}(|t_{j2} - r_{i2}| \leq h/2) \frac{1}{nh} \sum_{i=1}^{n} \mathbf{1}(|v_i - v_j| \leq h/2) \frac{1}{\frac{1}{n} \sum_{j=1}^{n} K_H(r_i, t_j)} \tag{94}$$

Since $|d - d'| \cdot \max\{|B_{11}|, |B_{21}|\} \geq h$, we have $\max\{|r_{i1} - s_{i1}|, |r_{i2} - s_{i2}|\} \geq h$ and hence it follows from (89) and (94) that

$$a_j \cdot a'_j = 0 \quad \text{for} \quad 1 \leq j \leq n. \tag{95}$$

Similar to (91), we apply (95) to establish that

$$\begin{aligned}
V_{\text{CATE}}^2 &\asymp \sum_{j=1}^{n} g(t_j)(1 - g(t_j)) \left( \frac{1}{h^2} \mathbf{1}(|t_{j1} - s_{i1}| \leq h/2) \mathbf{1}(|t_{j2} - s_{i2}| \leq h/2) \right)^2 \\
&\quad + \sum_{j=1}^{n} g(t_j)(1 - g(t_j)) \left( \frac{1}{h^2} \mathbf{1}(|t_{j1} - r_{i1}| \leq h/2) \mathbf{1}(|t_{j2} - r_{i2}| \leq h/2) \right)^2
\end{aligned} \tag{96}$$

We apply the same argument as (93) to establish that

$$V_{\text{CATE}}^2 \asymp n/h^2 \tag{97}$$

Similar to (92), we apply (95) to establish that

$$\sum_{j=1}^{n} |c_j|^{2+c} \leq \sum_{j=1}^{n} |a_j|^{2+c} + \sum_{j=1}^{n} |a'_j|^{2+c} \lesssim n \frac{1}{h^{2(1+c)}} \quad \text{for any positive constant} \quad c > 0. \tag{98}$$

Then for any positive constant $c > 0$, we have

$$\frac{1}{V_{\text{CATE}}^{1+\frac{c}{2}}} \sum_{j=1}^{n} \mathbb{E}\left[ (\epsilon_j c_j)^{(2+c)} \mid d_j, w_j \right] \cdot \mathbf{1}_{\mathcal{A}_3 \cap \mathcal{A}_4} \lesssim \frac{1}{(n/h^2)^{1+\frac{c}{2}}} \cdot n \cdot \frac{1}{h^{2(1+c)}} \leq \frac{1}{(nh^2)^{c/2}}$$

where the second inequality follows from (98) and bounded $\epsilon_j$. Hence, we have checked Lyapunov condition and shown that

$$\frac{\sum_{j=1}^n \epsilon_j c_j}{\sqrt{V_{\mathrm{CATE}}}} \mid \{d_j, w_j\}_{1 \leq j \leq n} \in \mathcal{A}_3 \cap \mathcal{A}_4 \xrightarrow{d} N(0,1)$$

Together with (90), we establish (75). We establish (76) by (97).

**Proof of** (90)   It follows from (61) and the condition that $\log n/(nh^3) \to 0$ that

$$\frac{1}{nh} \sum_{i=1}^n \mathbf{1}(|v_i - v_j| \leq h/2) \frac{1}{\frac{1}{n}\sum_{j=1}^n K_H(s_i, t_j)} \asymp \frac{1}{nhf_t(s_i)} \sum_{i=1}^n \mathbf{1}(|v_i - v_j| \leq h/2)$$

$$\asymp \frac{1}{nh} \sum_{i=1}^n \mathbf{1}(|v_i - v_j| \leq h/2)$$

where the last part holds since $f_t(s_i)$ is uniformly bounded from above and below across all $1 \leq i \leq n$. Note that for any fixed $1 \leq j \leq n$, we have

$$\left| \frac{1}{nh} \sum_{i=1}^n \mathbf{1}(|v_i - v_j| \leq h/2) - \frac{1}{nh} \sum_{i \neq j} \mathbf{1}(|v_i - v_j| \leq h/2) \right| \leq \frac{1}{nh} \tag{99}$$

and

$$c \leq \left| \mathbb{E}_{-j} \frac{1}{nh} \sum_{i \neq j} \mathbf{1}(|v_i - v_j| \leq h/2) \right| \leq C$$

where $\mathbb{E}_{-j}$ denotes the expectation conditioning on the $j$-th observation and some positive constants $c > 0$ and $C > 0$. We apply Lemma B.1 and establish that, with probability larger than $1 - n^{-C}$

$$\max_{1 \leq j \leq n} \left| \frac{1}{nh} \sum_{i \neq j} \mathbf{1}(|v_i - v_j| \leq h/2) - \mathbb{E}_{-j} \frac{1}{nh} \sum_{i \neq j} \mathbf{1}(|v_i - v_j| \leq h/2) \right| \lesssim \sqrt{\frac{\log n}{nh}}$$

Combined with (99), we have established $\mathbf{P}(\mathcal{A}_3) \geq 1 - n^{-C}$.

Since $\mathbb{E}_j \mathbf{1}(|t_{j1} - s_{i1}| \leq h/2)\mathbf{1}(|t_{j2} - s_{i2}| \leq h/2) \lesssim h^2$, we have

$$\mathbb{E}_j \left( \frac{1}{h^2} \mathbf{1}(|t_{j1} - s_{i1}| \leq h/2)\mathbf{1}(|t_{j2} - s_{i2}| \leq h/2) \right)^{1+\tau} \lesssim \frac{1}{h^{2\tau}}$$

and

$$\sum_{j=1}^{n} \mathbb{E}\left(\left(\frac{1}{h^2}\mathbf{1}(|t_{j1} - s_{i1}| \leq h/2)\mathbf{1}(|t_{j2} - s_{i2}| \leq h/2)\right)^{1+\tau}\right)^2 \lesssim n\frac{1}{h^{2+4\tau}}$$

Together with

$$\left(\frac{1}{h^2}\mathbf{1}(|t_{j1} - s_{i1}| \leq h/2)\mathbf{1}(|t_{j2} - s_{i2}| \leq h/2)\right)^{1+\tau} \lesssim \frac{1}{h^{2(1+\tau)}},$$

we apply Lemma B.1 and establish that, with probability larger than $1 - n^{-C}$,

$$\frac{1}{n}\sum_{j=1}^{n}\left(\frac{1}{h^2}\mathbf{1}(|t_{j1} - s_{i1}| \leq h/2)\mathbf{1}(|t_{j2} - s_{i2}| \leq h/2)\right)^{1+\tau} \lesssim \frac{1}{h^{2\tau}} + \frac{\sqrt{\log n}}{\sqrt{nh^{2+4\tau}}}$$

By the fact $\log n/(nh^3) \to 0$, we establish $\mathbf{P}(\mathcal{A}_4) \geq 1 - n^{-C}$.

## B.4   Proof of Lemma A.4

We use $\mathcal{T} \subset \mathbb{R}^3$ to denote the support of $f_t$ and define

$$\mathcal{T}^h = \left\{t \in \mathcal{T} : \mathcal{N}_{h/2}(t) \subset \mathcal{T}\right\}, \quad \text{with} \quad \mathcal{N}_{h/2}(t) = \left\{r \in \mathbb{R}^3 : \|r - t\|_\infty \leq h/2\right\}.$$

Here, $\mathcal{N}_{h/2}(t)$ denotes a specific $h/2$ neighborhood of $t$ and $\mathcal{T}^h$ denotes the set of $t$ such that it is not close to the boundary of $\mathcal{T}$.

**Proof of** (71).   We start with the decomposition

$$\frac{1}{n}\sum_{i=1}^{n}\frac{\sum_{j\neq i}[\nabla g(s_i)]^\intercal(t_j - s_i)K_H(s_i, t_j)}{\sum_{j=1}^{n}K_H(s_i, t_j)} = \frac{1}{n}\sum_{i=1}^{n}\frac{\sum_{j\neq i}[\nabla g(s_i)]^\intercal(t_j - s_i)K_H(s_i, t_j)}{\sum_{j=1}^{n}K_H(s_i, t_j)}\mathbf{1}_{s_i\in\mathcal{T}^h}$$
$$+ \frac{1}{n}\sum_{i=1}^{n}\frac{\sum_{j\neq i}[\nabla g(s_i)]^\intercal(t_j - s_i)K_H(s_i, t_j)}{\sum_{j=1}^{n}K_H(s_i, t_j)}\mathbf{1}_{s_i\notin\mathcal{T}^h}$$

$$(100)$$

We have

$$\left|\frac{1}{n}\sum_{i=1}^{n}\frac{\sum_{j\neq i}[\nabla g(s_i)]^\intercal(t_j - s_i)K_H(s_i, t_j)}{\sum_{j=1}^{n}K_H(s_i, t_j)}\mathbf{1}_{s_i\notin\mathcal{T}^h}\right| \lesssim h\frac{1}{n}\sum_{i=1}^{n}\mathbf{1}_{s_i\notin\mathcal{T}^h} \qquad (101)$$

Under the Condition 5.3 (c), the event $s_i = ((d, w^\intercal)B, v_i)^\intercal \notin \mathcal{T}^h$ implies that $[v_i - h/2, v_i + h/2]$ does not belong to the support $\mathcal{T}_v$ of $f_v$. That is,

$$\mathbb{E}[\mathbf{1}_{s_i \notin \mathcal{T}^h}] = \mathbb{P}([v_i - h/2, v_i + h/2] \not\subset \mathcal{T}_v) \leq \int_{v_{\min}-h/2}^{v_{\min}} f_v(v)dv + \int_{v_{\max}}^{v_{\max}+h/2} f_v(v)dv \leq h$$

where $v_{\min} = \inf_v \{v : f_v > 0\}$ and $v_{\max} = \sup_v \{v : f_v > 0\}$ denote the lower and upper boundaries of the support of $f_v$. If the support of $f_v$ is unbounded, we adopt the notation that $v_{\min} = -\infty$ and $v_{\max} = \infty$.

By applying the Bernstein inequality (Lemma B.1) to $\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{s_i \notin \mathcal{T}^h}$, we have

$$\mathbf{P}\left( \left| \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{s_i \notin \mathcal{T}^h} - \mathbb{E}\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{s_i \notin \mathcal{T}^h} \right| \gtrsim \sqrt{\frac{h \log n}{n}} \right) \leq n^{-C}$$

Since $\left| \mathbb{E}\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{s_i \notin \mathcal{T}^h} \right| \lesssim h$, we have $\mathbf{P}\left( \left| \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{s_i \notin \mathcal{T}^h} \right| \gtrsim h \right) \leq n^{-C}$. Hence, we further upper bound (101) by

$$\mathbf{P}\left( \left| \frac{1}{n} \sum_{i=1}^n \frac{\sum_{j \neq i}[\nabla g(s_i)]^\intercal(t_j - s_i)K_H(s_i, t_j)}{\sum_{j=1}^n K_H(s_i, t_j)} \mathbf{1}_{s_i \notin \mathcal{T}^h} \right| \gtrsim h^2 \right) \leq n^{-C}. \tag{102}$$

Now, we control the first term on the right hand side of (100),

$$\left| \frac{1}{n} \sum_{i=1}^n \frac{\sum_{j \neq i}[\nabla g(s_i)]^\intercal(t_j - s_i)K_H(s_i, t_j)}{\sum_{j=1}^n K_H(s_i, t_j)} \mathbf{1}_{s_i \in \mathcal{T}^h} \right| \leq \max_{1 \leq i \leq n} \left| \frac{\frac{1}{n}\sum_{j \neq i}[\nabla g(s_i)]^\intercal(t_j - s_i)K_H(s_i, t_j)}{\frac{1}{n}\sum_{j=1}^n K_H(s_i, t_j)} \mathbf{1}_{s_i \in \mathcal{T}^h} \right|$$

Now we fix $i \in \{1, \cdots, n\}$ and condition on the $i$-th observation. For $j \neq i$, we use $\mathbb{E}_j$ denotes the expectation is taken with respect to the $j$-th observation, conditioning on the $i$-th observation. We can focus on the case $s_i \in \mathcal{T}^h$ since, otherwise, we have the trivial upper bound $0$. We define

$b = t_j - s_i$ and then we obtain

$$\mathbb{E}_j[\nabla g(s_i)]^\intercal [t_j - s_i] K_H[s_i, t_j] \mathbf{1}_{s_i \in \mathcal{T}^h}$$

$$= \int_{\|b\|_\infty \leq h/2} [\nabla g(s_i)]^\intercal b \frac{1}{h^3} f_t(s_i + b) db \cdot \mathbf{1}_{s_i \in \mathcal{T}^h}$$

$$= \int_{\|b\|_\infty \leq h/2} [\nabla g(s_i)]^\intercal b \frac{1}{h^3} [f_t(s_i) + b^\intercal \nabla f_t(s_i + cb)] db \cdot \mathbf{1}_{s_i \in \mathcal{T}^h}$$

$$= \frac{1}{h^3} \int_{\|b\|_\infty \leq h/2} [\nabla g(s_i)]^\intercal b b^\intercal \nabla f_t(s_i + cb) db \cdot \mathbf{1}_{s_i \in \mathcal{T}^h}$$

where the last equality follows from the fact that $\int_{\|b\|_\infty \leq h/2} b \, db = 0$.

Since $|[\nabla g(s_i)]^\intercal b b^\intercal \nabla f_t(s_i + cb)| \cdot \mathbf{1}_{s_i \in \mathcal{T}^h} \leq Ch^2$, we have

$$|\mathbb{E}_j[\nabla g(s_i)]^\intercal [t_j - s_i] K_H[s_i, t_j] \mathbf{1}_{s_i \in \mathcal{T}^h}| \lesssim Ch^2. \tag{103}$$

Now, it is sufficient to control $\left| \frac{1}{n} \sum_{j \neq i} ([\nabla g(s_i)]^\intercal [t_j - s_i] K_H[s_i, t_j] - \mathbb{E}_j[\nabla g(s_i)]^\intercal [t_j - s_i] K_H[s_i, t_j]) \right|$

Since $|[\nabla g(s_i)]^\intercal [t_j - s_i] K_H[s_i, t_j] - \mathbb{E}_j[\nabla g(s_i)]^\intercal [t_j - s_i] K_H[s_i, t_j]| \lesssim h \cdot \frac{1}{h^3}$ and

$$\sum_{j \neq i} \mathbb{E}_j |[\nabla g(s_i)]^\intercal [t_j - s_i] K_H[s_i, t_j] - \mathbb{E}_j[\nabla g(s_i)]^\intercal [t_j - s_i] K_H[s_i, t_j]|^2$$

$$\leq n \mathbb{E}_j [[\nabla g(s_i)]^\intercal [t_j - s_i] K_H[s_i, t_j]]^2 \lesssim nh^2/h^3$$

By applying Lemma B.1, we establish that, with probability larger than $1 - n^{-C}$,

$$\left| \frac{1}{n} \sum_{j \neq i} ([\nabla g(s_i)]^\intercal [t_j - s_i] K_H[s_i, t_j] - \mathbb{E}_j[\nabla g(s_i)]^\intercal [t_j - s_i] K_H[s_i, t_j]) \right| \lesssim \sqrt{\frac{\log n}{nh}} \tag{104}$$

Combined with (102) and (103), we establish (71).

**Proof of** (72). The proof of (72) follows from

$$\left| [\nabla g(\widehat{s}_i)]^\intercal (\widehat{t}_i - \widehat{s}_i) K_H(\widehat{s}_i, \widehat{t}_i) \right| \leq \frac{1}{h^2}$$

and (67).

**Proof of** (73). We define

$$\mathcal{A}_{ij} = \left\{ h - 2C_0\sqrt{\frac{\log p}{n}} \leq \min_{1 \leq l \leq 3} |t_{j,l} - s_{i,l}| \leq \max_{1 \leq l \leq 3} |t_{j,l} - s_{i,l}| \leq h + 2C_0\sqrt{\frac{\log p}{n}} \right\} \cap \mathcal{A}_1 \cap \mathcal{A}_2$$

On the event $\mathcal{A}_{ij}$, we have

$$K_H(\widehat{s}_i, \widehat{t}_j) = K_H(s_i, t_j) \tag{105}$$

We start with

$$\left| \frac{1}{n} \sum_{i=1}^{n} \frac{\sum_{j \neq i} [\nabla g(\widehat{s}_i)]^\intercal (\widehat{t}_j - \widehat{s}_i) K_H(\widehat{s}_i, \widehat{t}_j)}{\sum_{j=1}^{n} K_H(\widehat{s}_i, \widehat{t}_j)} - \frac{1}{n} \sum_{i=1}^{n} \frac{\sum_{j \neq i} [\nabla g(s_i)]^\intercal (t_j - s_i) K_H(s_i, t_j)}{\sum_{j=1}^{n} K_H(s_i, t_j)} \right|$$

$$\leq \max_{1 \leq i \leq n} \left| \frac{\frac{1}{n} \sum_{j \neq i} [\nabla g(\widehat{s}_i)]^\intercal (\widehat{t}_j - \widehat{s}_i) K_H(\widehat{s}_i, \widehat{t}_j)}{\frac{1}{n} \sum_{j=1}^{n} K_H(\widehat{s}_i, \widehat{t}_j)} - \frac{\frac{1}{n} \sum_{j \neq i} [\nabla g(s_i)]^\intercal (t_j - s_i) K_H(s_i, t_j)}{\frac{1}{n} \sum_{j=1}^{n} K_H(s_i, t_j)} \right| \tag{106}$$

Now we fix $i \in \{1, 2, \cdots, n\}$. We define

$$\widehat{b}_j = [\nabla g(\widehat{s}_i)]^\intercal (\widehat{t}_j - \widehat{s}_i) K_H(\widehat{s}_i, \widehat{t}_j) \quad b_j = [\nabla g(s_i)]^\intercal (t_j - s_i) K_H(s_i, t_j)$$

and then it is equivalent to control

$$\frac{\frac{1}{n} \sum_{j \neq i} \widehat{b}_j}{\frac{1}{n} \sum_{j=1}^{n} K_H(\widehat{s}_i, \widehat{t}_j)} - \frac{\frac{1}{n} \sum_{j \neq i} b_j}{\frac{1}{n} \sum_{j=1}^{n} K_H(s_i, t_j)}$$

$$= \frac{\frac{1}{n} \sum_{j \neq i} (\widehat{b}_j - b_j)}{\frac{1}{n} \sum_{j=1}^{n} K_H(\widehat{s}_i, \widehat{t}_j)} + \frac{\frac{1}{n} \sum_{j \neq i} b_j}{\frac{1}{n} \sum_{j=1}^{n} K_H(s_i, t_j)} \left( \frac{\frac{1}{n} \sum_{j=1}^{n} K_H(s_i, t_j)}{\frac{1}{n} \sum_{j=1}^{n} K_H(\widehat{s}_i, \widehat{t}_j)} - 1 \right) \tag{107}$$

We now control $\frac{1}{n} \sum_{j \neq i} (\widehat{b}_j - b_j)$ as

$$\frac{1}{n} \sum_{j \neq i} (\widehat{b}_j - b_j) = \frac{1}{n} \sum_{j \neq i} (\widehat{b}_j - b_j) \cdot \mathbf{1}_{\mathcal{A}_{ij}^c} + \frac{1}{n} \sum_{j \neq i} (\widehat{b}_j - b_j) \cdot \mathbf{1}_{\mathcal{A}_{ij}} \tag{108}$$

It follows from (105) that

$$(\widehat{b}_j - b_j) \cdot \mathbf{1}_{\mathcal{A}_{ij}^c} = \left( [\nabla g(\widehat{s}_i)]^\intercal (\widehat{t}_j - \widehat{s}_i) - [\nabla g(s_i)]^\intercal (t_j - s_i) \right) \cdot \mathbf{1}_{\mathcal{A}_{ij}} \cdot K_H(s_i, t_j) \tag{109}$$

and hence

$$\left| \frac{1}{n} \sum_{j \neq i} (\widehat{b}_j - b_j) \cdot \mathbf{1}_{\mathcal{A}_{ij}} \right| \leq \max_{1 \leq i \leq n} \left| [\nabla g(\widehat{s}_i)]^\mathsf{T}(\widehat{t}_j - \widehat{s}_i) - [\nabla g(s_i)]^\mathsf{T}(t_j - s_i) \right| \frac{1}{n} \sum_{j \neq i} \mathbf{1}_{\mathcal{A}_{ij}} \cdot K_H(s_i, t_j)$$

On the event $\mathcal{A}_1 \cap \mathcal{A}_2$, we have $\max_{1 \leq i \leq n} \left| [\nabla g(\widehat{s}_i)]^\mathsf{T}(\widehat{t}_j - \widehat{s}_i) - [\nabla g(s_i)]^\mathsf{T}(t_j - s_i) \right| \lesssim \log n / \sqrt{n}$ and $\frac{1}{n} \sum_{j \neq i} \mathbf{1}_{\mathcal{A}_{ij}} \cdot K_H(s_i, t_j) \leq \frac{1}{n} \sum_{j \neq i} K_H(s_i, t_j)$, we apply (61) and establish that, with probability larger than $1 - n^{-C}$,

$$\left| \frac{1}{n} \sum_{j \neq i} (\widehat{b}_j - b_j) \cdot \mathbf{1}_{\mathcal{A}_{ij}} \right| \lesssim \frac{\log n}{\sqrt{n}} \cdot |f_t(s_i)| . \tag{110}$$

Note that

$$\left| \frac{1}{n} \sum_{j \neq i} (\widehat{b}_j - b_j) \cdot \mathbf{1}_{\mathcal{A}_{ij}^c} \right| \lesssim \frac{1}{h^2} \frac{1}{n} \sum_{j \neq i} \mathbf{1}_{\mathcal{A}_{ij}^c} \leq \frac{1}{nh^2} \sum_{j \neq i} h^3 (K_H^b - K_H^a)(s_i, t_j)$$

where the kernels $K_H^b$ and $K_H^a$ are defined in (79). By combining (81), (B.2) and $nh^4 \gg (\log n)^2$, we establish that with probability larger than $1 - n^{-C}$,

$$\left| \frac{1}{n} \sum_{j \neq i} (\widehat{b}_j - b_j) \cdot \mathbf{1}_{\mathcal{A}_{ij}^c} \right| \lesssim \frac{\log n}{\sqrt{n}} .$$

Combined with (110), we establish that, with probability larger than $1 - n^{-C}$,

$$\left| \frac{1}{n} \sum_{j \neq i} (\widehat{b}_j - b_j) \right| \lesssim \frac{\log n}{\sqrt{n}} . \tag{111}$$

By (61), (62), (67) (104) and (103), we have

$$\left| \frac{\frac{1}{n} \sum_{j \neq i} b_j}{\frac{1}{n} \sum_{j=1}^n K_H(s_i, t_j)} \left( \frac{\frac{1}{n} \sum_{j=1}^n K_H(s_i, t_j)}{\frac{1}{n} \sum_{j=1}^n K_H(\widehat{s}_i, \widehat{t}_j)} - 1 \right) \right| \lesssim \left( Ch^2 + \sqrt{\frac{\log n}{nh}} \right) \frac{\log n}{\sqrt{n}h}$$

Combined with (111), we establish (73).

70

# C  Simulations and data applications

## C.1  Implementations in Section 6

For the "Valid-CF" method, we first estimate the control variable $v_i$ as in (19). As we have no observed confounders here, the "Valid-CF" identifies

$$\mathbb{E}[y_i|d_i, w_i, v_i] = \tilde{g}(d_i, v_i)$$

for some unknown function $\tilde{g}$. Hence,

$$\mathbb{E}[y_i^{(d)}|w_i = w, v_i = v] = \tilde{g}(d, v) \ \text{ and } \ \text{ASF}(d, w) = \int \tilde{g}(d, v_i)f_v(v_i)dv_i.$$

We implement "Valid-CF" by estimating $\tilde{g}$ by a two-dimensional kernel estimators and apply partial mean to estimate the causal effects.

For the "Logit-Median" method, we estimate $\hat{\gamma}$ as in (19) and estimate $\widehat{\mathcal{S}}$ as in (23). Define

$$(\check{\Phi}, \check{\rho}) = \arg\max_{\Phi, \rho} \sum_{i=1}^{n} \{y_i \log \text{logit}(w_i\Phi + \widehat{v}_i\rho) + (1 - y_i) \log(1 - \text{logit}(w_i\Phi + \widehat{v}_i\rho))\}.$$

Then we estimate $\beta$ via

$$\check{\beta} = \text{Median} \left( \{\check{\Phi}_j/\hat{\gamma}_j\}_{j\in\widehat{\mathcal{S}}} \right).$$

We estimate the invalid effects as $\check{\pi} = \check{\Phi} - \check{\beta}\hat{\gamma}$. Then we estimate $\text{CATE}(d, d'|w)$ with

$$\text{logit}(d\check{\beta} + w^\intercal\check{\pi}) - \text{logit}(d'\check{\beta} + w^\intercal\check{\pi}).$$

The standard deviation of the estimated $\text{CATE}(d, d'|w)$ is based on 50 bootstrap resampling.

For the "TSHT" method, we use the R code from Guo et al. (2018), which deals with invalid IVs in linear models.

In Table 5, we report the inference results for $\text{CATE}(-2, 2|w)$ in binary outcome models (i) and (ii) with Logit-Median. For Logit-Median, its coverage probabilities are also close to the nominal level. This implies a mild effect of misspecified model (i) as logistic. It can be partially seen from Figure 3 of the main paper that that the functional form of ASF is close to the logistic function. In

this setting, we see that the Logit-Median method has coverage probabilities close to the nominal level. In model (ii), the logistic model is severely misspecified. The coverage probabilities of Logit-Median decrease as sample sizes get larger and as IVs get stronger. This demonstrates the bias caused by model misspecification.

| | | Binary(i) | | | | | | Binary (ii) | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | $N(0, I_p)$ | | | $U[-1.73, 1.73]$ | | | $N(0, I_p)$ | | | $U[-1.73, 1.73]$ | | |
| $n$ | $c_\gamma$ | MAE | COV | SE | MAE | COV | SE | MAE | COV | SE | MAE | COV | SE |
| 500 | 0.4 | 0.090 | 0.977 | 0.15 | 0.081 | 0.973 | 0.15 | 0.085 | 0.950 | 0.14 | 0.088 | 0.963 | 0.14 |
| 500 | 0.6 | 0.057 | 0.967 | 0.10 | 0.064 | 0.980 | 0.10 | 0.064 | 0.963 | 0.10 | 0.067 | 0.950 | 0.10 |
| 500 | 0.8 | 0.047 | 0.973 | 0.08 | 0.040 | 0.977 | 0.08 | 0.054 | 0.940 | 0.07 | 0.057 | 0.923 | 0.07 |
| 1000 | 0.4 | 0.069 | 0.963 | 0.10 | 0.059 | 0.967 | 0.10 | 0.071 | 0.927 | 0.10 | 0.065 | 0.953 | 0.10 |
| 1000 | 0.6 | 0.049 | 0.963 | 0.07 | 0.042 | 0.963 | 0.07 | 0.050 | 0.930 | 0.07 | 0.055 | 0.943 | 0.07 |
| 1000 | 0.8 | 0.033 | 0.963 | 0.05 | 0.035 | 0.967 | 0.05 | 0.043 | 0.877 | 0.05 | 0.054 | 0.850 | 0.05 |
| 2000 | 0.4 | 0.041 | 0.966 | 0.07 | 0.046 | 0.973 | 0.07 | 0.056 | 0.933 | 0.07 | 0.049 | 0.940 | 0.07 |
| 2000 | 0.6 | 0.031 | 0.960 | 0.05 | 0.033 | 0.943 | 0.05 | 0.043 | 0.880 | 0.05 | 0.050 | 0.850 | 0.05 |
| 2000 | 0.8 | 0.041 | 0.966 | 0.07 | 0.020 | 0.973 | 0.04 | 0.045 | 0.777 | 0.04 | 0.047 | 0.777 | 0.04 |

Table 5: Inference of CATE$(-2, 2|w)$ in binary outcome models (i) and (ii) with Logit-Median. We report the median absolute errors (MAE) for $\widehat{\text{CATE}}(-2, 2|w)$ and average coverage probabilities (COV) and average standard error (SE) for the confidence intervals of $\mu$ where $w_i$ are *i.i.d.* Gaussian or uniform with range $[-1.73, 1.73]$. Each setting is replicated with 300 independent experiments.

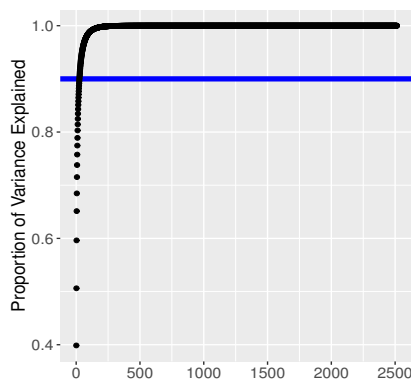## C.2 The results of PCA in Section 7



Figure 5: The cumulative proportion of explained variance by the 2514 principle components (PCs) for HDL exposure.
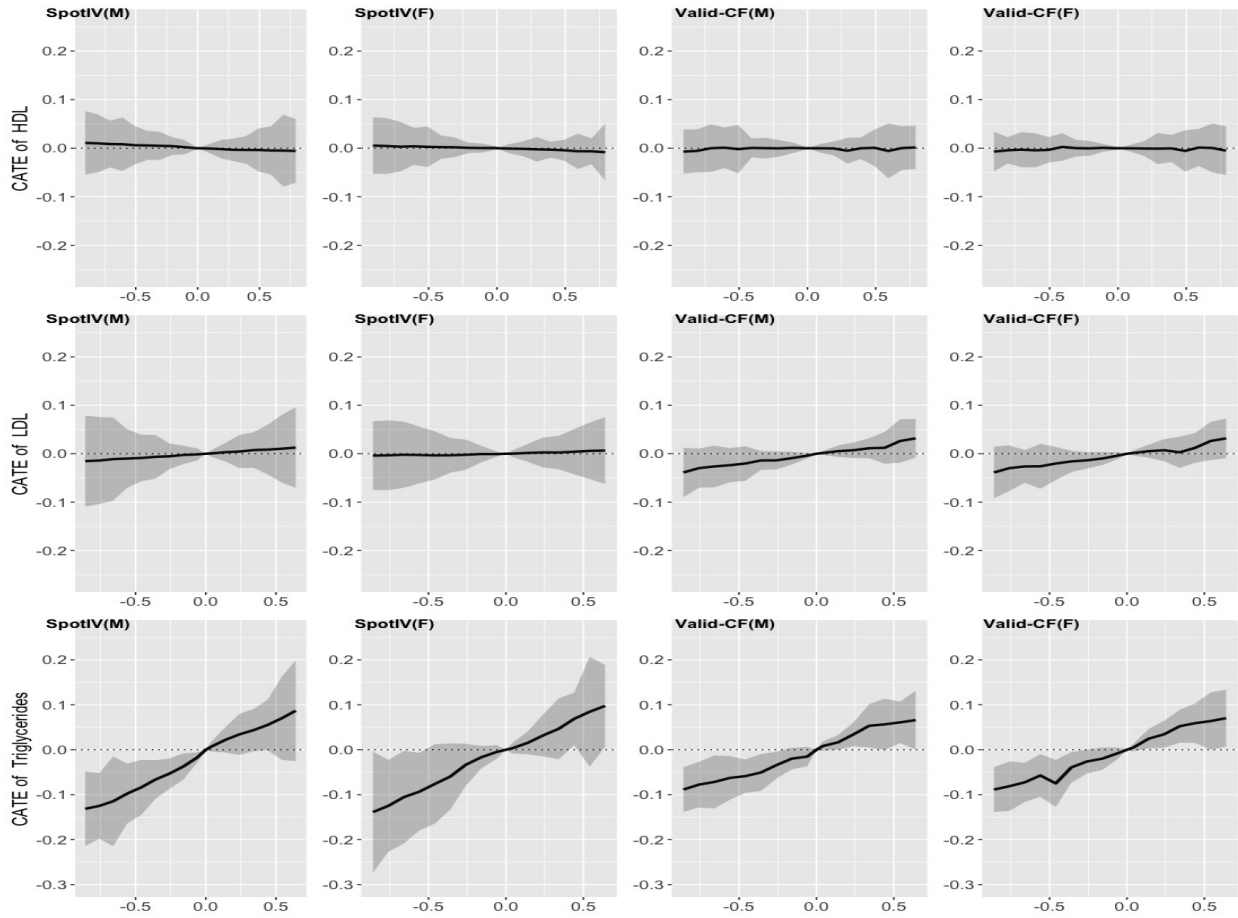
Figure 6: The constructed 95% CIs for CATE$(d, 0|x_M)$ and for CATE$(d, 0|x_F)$ with HDL, LDL, and Triglycerides exposures at different levels of $d$. The first and third columns report the results given by spotIV and Valid-CF for CATE$(d, 0|x_M)$, respectively. The second and fourth columns report the results given by spotIV and Valid-CF for CATE$(d, 0|x_F)$, respectively.